

1 Leveraging distant relatedness to quantify human
2 mutation and gene conversion rates

3 Pier Francesco Palamara^{1,2,*}, Laurent Francioli³, Giulio Genovese², Peter Wilton⁴,
4 Alexander Gusev^{1,2}, Hilary Finucane^{1,2}, Sriram Sankararaman^{2,5}, The Genome of the
5 Netherlands Consortium, Shamil Sunyaev^{2,5}, Paul I.W. de Bakker^{3,6}, John Wakeley⁴,
6 Itsik Pe'er⁷, and Alkes L. Price^{1,2,8}

7 ¹*Department of Epidemiology, Harvard T. H. Chan School of Public Health, Boston, MA, 02115, U.S.A.*

8 ²*Program in Medical and Population Genetics, Broad Institute of Harvard and MIT, Cambridge, MA,
9 02142, U.S.A.*

10 ³*Department of Epidemiology, Julius Center for Health Sciences and Primary Care, University Medical
11 Center Utrecht, Utrecht, The Netherlands.*

12 ⁴*Department of Organismic and Evolutionary Biology, Harvard University, Cambridge, MA, 02138,
13 U.S.A.*

14 ⁵*Department of Genetics, Harvard Medical School, Boston, MA, 02115, U.S.A.*

15 ⁶*Department of Medical Genetics, Center for Molecular Medicine, University Medical Center Utrecht,
16 Utrecht, The Netherlands.*

17 ⁷*Department of Computer Science, Columbia University, New York City, NY, 10027, U.S.A.*

18 ⁸*Department of Biostatistics, Harvard T. H. Chan School of Public Health, Boston, MA, 02115, U.S.A.*

19 * *Correspondence: ppalama@hsph.harvard.edu*

Abstract

The rate at which human genomes mutate is a central biological parameter that has many implications for our ability to understand demographic and evolutionary phenomena. We present a method for inferring mutation and gene conversion rates using the number of sequence differences observed in identical-by-descent (IBD) segments together with a reconstructed model of recent population size history. This approach is robust to, and can quantify, the presence of substantial genotyping error, as validated in coalescent simulations. We applied the method to 498 trio-phased Dutch individuals from the Genome of the Netherlands (GoNL) project, sequenced at an average depth of 13x. We infer a point mutation rate of $1.66 \pm 0.04 \times 10^{-8}$ per base per generation, and a rate of $1.26 \pm 0.06 \times 10^{-9}$ for < 20 bp indels. Our estimated average genome-wide mutation rate is higher than most pedigree-based estimates reported thus far, but lower than estimates obtained using substitution rates across primates. By quantifying how estimates vary as a function of allele frequency, we infer the probability that a site is involved in non-crossover gene conversion as $5.99 \pm 0.69 \times 10^{-6}$, consistent with recent reports. We find that recombination does not have observable mutagenic effects after gene conversion is accounted for, and that local gene conversion rates reflect recombination rates. We detect a strong enrichment for recent deleterious variation among mismatching variants found within IBD regions, and observe summary statistics of local IBD sharing to closely match previously proposed metrics of background selection, but find no

41 significant effects of selection on our estimates of mutation rate. We detect no ev-
42 idence for strong variation of mutation rates in a number of genomic annotations
43 obtained from several recent studies.

44 **Introduction**

45 Germline mutations represent a fundamental evolutionary force that shapes phenotypic
46 variation and has a profound impact on heritable diversity. Precise estimation of muta-
47 tion rates has several applications, including the interpretation of mutations implicated in
48 diseases [1, 2, 3, 4], studies of natural selection [5, 6, 7], the timing of demographic events
49 inferred using genetic analysis [8, 9, 10, 11], and the study of several aspects of human
50 mutagenesis [12]. High throughput sequencing technologies have recently enabled the
51 quantification of germline mutation rates, but the estimates obtained using these meth-
52 ods are inconsistent with previous studies. The source of these inconsistencies, whether
53 biological or due to methodological biases, is at the center of recent debate [10], and new
54 methods are required to gain additional insight into germline mutation rates.

55 Several aspects of methods for inferring mutation rates were recently reviewed in
56 [13, 14]. Approaches based on sequence divergence among primates [15, 16] (phylogenetic
57 methods) estimate mutation rates that range between $2.0 - 2.5 \times 10^{-8}$. These meth-
58 ods depend on a number of population-genetic assumptions, and further uncertainty is
59 contributed by the need to map inferred per-year mutation rates to a per-generation

60 scale. Pedigree-based estimates of mutation rates ranging between $0.9 - 1.3 \times 10^{-8}$
61 [17, 18, 19, 3, 20, 21], on the other hand, are based on direct observation of mutation
62 events, but issues related to sequencing error complicate the inference, requiring the use
63 of stringent genotyping filters. Several statistical inconsistencies, such as underdispersion
64 of the reported estimates, have been outlined in [13]. This motivates the development of
65 new methods for estimating mutation rates. Two methods based on the pairwise sequen-
66 tially Markovian coalescent (PSMC [9]) recently estimated slightly higher mutation rates
67 than pedigree-based studies ([22, 23]).

68 In this work, we propose a new method for estimating mutation rates using muta-
69 tions occurring within identical-by-descent haplotype blocks (IBD [24, 25, 26, 27, 28, 29,
70 30, 31, 32]) transmitted through recent common ancestors that lived in the past ~ 100
71 generations ($\sim 3,000$ years) before present. Our approach is robust to, and can quantify,
72 the presence of substantial amounts of genotyping error in the analyzed sequences. By
73 quantifying how estimates vary as a function of allele frequency, we correct our estimate
74 for additional biases created by the occurrence of gene conversion events on the ancestral
75 lineages leading to shared common ancestors. We apply this methodology to analyze
76 250 trio families from the Netherlands sequenced at an average of $\sim 13\times$, obtaining a
77 genome-wide average point mutation rate estimate of $1.66 \pm 0.04 \times 10^{-8}$ per base, per
78 generation. We further apply our methodology to infer the rate of < 20 bp indels, which
79 we estimated to be $1.26 \pm 0.06 \times 10^{-9}$. We analyze the relationship between recombina-

80 tion rates and mutation rates [33], and find that after accounting for the occurrence of
81 gene conversion events, no significant association is detected. In addition to estimating
82 the rate of mutation events, we derived an estimate for the rate at which a genomic lo-
83 cus is involved in a non-crossover gene conversion event of $5.99 \pm 0.69 \times 10^{-6}$. We find
84 that mismatching variants within IBD regions - representing mutations of recent origin
85 - are strongly enriched for deleterious variation. We show that the length of IBD shar-
86 ing along the genome closely reflects widely adopted summary statistics of background
87 selection [5]. Background selection, however, is observed not to have significant effects
88 on our estimates. Finally, we explore enrichment or depletion of several specific genomic
89 annotations, and find no evidence for substantial differences compared to the average
90 genome-wide mutation rate.

91 **Materials and Methods**

92 *Overview of methods*

93 Pairs of purportedly unrelated individuals from a population often share long stretches of
94 chromosomal regions inherited identical-by-descent (IBD) from recent common ancestors
95 that lived in the past tens to few hundreds of generations. These IBD segments can be
96 detected using several available methods [34, 25, 35, 36], and reflect genetic relationships
97 that are typically not known to the affected individuals, but are found to be ubiquitous
98 even in outbred populations [37, 30]. IBD segments are defined in our work as contiguous

99 chromosomal regions for which two sampled chromosomes share the same most recent
100 common ancestor (MRCA). In this work, we are interested in IBD segments co-inherited
101 from ancestors that lived within ~ 100 generations before present, which can be reliably
102 detected in trio-phased real data. Occasional mutations segregating along the lineages
103 connecting a pair of IBD haplotypes to their MRCA will create mismatched sites on the
104 shared haplotypes that can be used to infer the rate at which new germline mutations
105 appear. If the exact number of generations separating the IBD segments (via their MRCA)
106 is known, one may infer the mutation rate by dividing the number of observed sequence
107 mismatches by the number of generations and the physical length, for all segments. A
108 special case of this approach is used in trio-based analyses, where transmitted parental
109 haplotypes and IBD offspring haplotypes are separated by a single generation.

110 We briefly describe solutions to three challenges. First, to estimate the number of
111 generations separating two IBD segments (twice the time to most recent common ancestor;
112 tMRCA), we use a recently developed method [29] that relies on the spectrum of observed
113 IBD segment lengths to infer demographic history, which is then used to obtain a posterior
114 mean estimate of the average tMRCA for pools of IBD segments of different lengths,
115 as detailed in the real data description and the Appendix. Second, to deal with the
116 presence of genotyping errors, rather than relying on stringent filtering criteria (as in
117 trio-based analyses [17, 18, 19, 3, 20, 21]), we regress the observed sequence mismatches
118 for several IBD length thresholds on the estimated tMRCA; the slope of this regression

119 reflects the rate at which new mutations accumulate per generation time unit, while
120 the genotyping error rate is captured by the intercept. We refer to this procedure as
121 tMRCA regression (illustrated in Figure 1). Finally, we correct for the occurrence of non-
122 crossover gene conversion events along the lineages leading to the MRCA, exploiting the
123 relationship between an allele's frequency and the probability that it is involved in a gene
124 conversion event. To do this, we repeatedly perform tMRCA regression, only considering
125 mismatching sites with frequency less than a specified maximum allele frequency (MaAF)
126 value, and regress the obtained mutation rate estimates on the MaAF threshold. We show
127 that the intercept of this regression provides a gene conversion-corrected mutation rate
128 estimate. We refer to this procedure as MaAF-threshold regression (illustrated in Figure
129 2). This also allows us to estimate the rate at which a genomic locus is involved in a non-
130 crossover gene conversion event, which is proportional to the difference between corrected
131 and uncorrected estimates for mutation rates. We have released open-source software
132 (IBDMUT) implementing these methods (see Web Resources).

133 *IBD detection and demographic inference*

134 Shared IBD segments can be detected in large samples of purportedly unrelated indi-
135 viduals using several available algorithms [34, 25, 35, 36]. These approaches typically
136 rely on the similarity of haplotypes within pairs of individuals (identity-by-state), and
137 probabilistic models that enable confidently detecting long (e.g. > 1 cM) shared IBD

138 segments. The accuracy of IBD detection is substantially improved when trio-phased in-
139 dividuals are available, as is the case in the analyzed data set. For both simulations and
140 real data analysis reported in this paper, we used the GERMLINE [25] IBD detection
141 algorithm. Inferring demographic history for the past ~ 100 generations can be achieved
142 by analyzing the length distributions of shared IBD segments, which provide information
143 on the distribution of recent coalescent events within a population [29, 38, 30]. Details
144 of the IBD detection and demographic inference methods for the reported real data anal-
145 ysis are described in the “GoNL data set” section. In order to test the accuracy of the
146 proposed methodology, we relied on “ground truth” IBD sharing in some simulations, i.e.
147 IBD segments extracted from the sampled ancestral recombination graphs (ARG), and
148 made use of the simulated demographic history to inform mutation and gene conversion
149 rate inference in several synthetic scenarios.

150 *Estimating the mutation rate via tMRCA regression*

151 The proposed methodology for the inference of mutation rates requires the availability of
152 haploid genotype data, a list of IBD segments between pairs of haploid individuals that
153 are longer than a specified Morgan length threshold, including start/end position, and a
154 demographic model, which may be inferred from the spectrum of IBD shared segments
155 as described in [29]. For each IBD segment i , we obtain an observed mismatch rate by
156 counting the number of sequence differences m_i in the haploid genotypes within the region,

157 and dividing by the region size s_i in base pairs: $\theta_i = m_i/s_i$. The observed mismatch rate
158 is then obtained by averaging all observations $\hat{\theta}_u = n_u^{-1} \sum_{i=1}^{n_u} \theta_i$, for n_u segments longer
159 than u Morgans. We repeat this measurement for several thresholds u , obtaining a vector
160 of observed mismatch rates $\hat{\theta}$. Due to the lack of detailed pedigree structures at deep
161 time scales, the exact number of meiotic events separating two individuals that share IBD
162 segments is generally unknown. Using the reconstructed demographic model, we therefore
163 infer the posterior mean age t_u of pooled IBD segments longer than a known genetic length
164 threshold u , using coalescent theory recently developed in [29, 30], the details of which
165 are summarized in the Appendix. Finally, we regress the observed mismatch rates $\hat{\theta}_u$
166 on twice the posterior mean age (in generations) to the MRCA of the IBD segments
167 $\hat{\theta} = \alpha + 2\mu\mathbf{t} + \epsilon$. We refer to this regression as the tMRCA regression. Older segments
168 will tend to harbor a larger number of sequence differences, due to the higher chance
169 of mutation events occurring along the lineages connecting extant individuals to their
170 most recent common ancestor. The slope μ of this regression will capture the rate at
171 which mutations arise per unit of time. Note that we are neglecting the uncertainty on
172 the measurement in the regressor \mathbf{t} , i.e. the inferred age of the pooled IBD segments.
173 As shown in simulations, however, this only results in negligible biases for the estimated
174 slope coefficient due to the large number of pooled segments.

175 Assuming a genotyping error model for which false positive/negative genotype calls
176 are independent of the average coalescent time of pairs of individuals at a locus, the

177 intercept α of this regression is expected to capture the rate at which genotyping errors
178 occur on the considered range of IBD segments.

179 *Controlling for gene conversion via MaAF-threshold regression*

180 Non-crossover gene conversion events occur at a rate that is correlated to recombination,
181 and have been observed to have frequency higher than recombination events [39]. In the
182 coalescent process, gene conversion may be modeled as two consecutive recombination
183 events that occur very close to each other [40], at an average distance of ~ 300 base pairs
184 [41, 42]. These events introduce the possibility that polymorphisms that are segregating
185 in the population may be assimilated into haplotypes within IBD regions. These poly-
186 morphisms can create sequence differences between IBD individuals that are not due to
187 newly arising mutations. Note that, whereas gene conversion events change the MRCA
188 of the ~ 300 bp converted segment, here we do not consider this to break an IBD block.
189 Furthermore, because the number of gene conversion events is related to the number of
190 meiotic events, short IBD regions will tend to exhibit more gene conversion-driven mis-
191 matches than longer, more recent IBD segments, therefore resulting in an upward bias
192 when mutation rate is estimated via the slope of tMRCA regression. The mismatching
193 variants observed on IBD segments, therefore, will be due to at least two distinct sources
194 of heterozygosity. The first, θ_p , which we call population heterozygosity in the remainder,
195 represents the effect of gene conversion events which introduce standing genetic varia-

196 tion onto IBD blocks. The second source of heterozygosity is due to newly arising point
197 mutations on IBD blocks, and will be referred to as θ_μ . For IBD segments of a chosen
198 length, we can express the total observed mismatch rate as $\theta = \theta_\mu + \theta_p$. To estimate the
199 mutation rate due to point mutations only, we need to exclude the effects of θ_p from our
200 calculations. We make the following two observations:

- 201 1. The frequency of mutations that arise on long (e.g. ≥ 1 cM) IBD segments is
202 typically low in the population (Figure S1), so that θ_μ is mostly due to rare variants.
- 203 2. If we divide the allele frequency spectrum into bins of equal width, we find an
204 approximately uniform contribution to θ_p for each frequency. This implies that
205 if we compute the frequency-bounded population heterozygosity $\theta_{p,f}$, using only
206 variants of frequency at most f , we observe an approximately linear relationship
207 between $\theta_{p,f}$ and f (Figure S2, see additional calculations in the Appendix).

208 Observation (1) implies that if we exclude high frequency variants when we compute μ
209 using the proposed regression approach, the contribution of θ_μ to the observed mismatch
210 rate on IBD segments will be largely unaffected. Furthermore, observation (2) suggests
211 that if we estimate a frequency-bounded value of μ_f by ignoring variants of frequency
212 higher than a threshold f , the contribution of population heterozygosity due to gene
213 conversion events, $\theta_{p,f}$, will be decreased to an extent that is approximately linear in
214 f . Assuming that the contribution of θ_μ to μ_f is unaffected for values of f in the range

215 $F = [F_{min}, F_{max}]$, we may therefore regress μ_F on F , and observe a linear relationship. We
216 refer to this regression as the MaAF-threshold regression (Figure 2). The intercept of this
217 regression will then reflect an estimate of μ without the confounding effects of θ_p , while
218 the contribution of θ_μ is left unchanged. We avoid computing values of μ_F corresponding
219 to $F \in [0, F_{min})$, for a sufficiently large F_{min} (e.g. > 0.1), as this may result in removing
220 variants that are due to new point mutation events on the IBD segments, which we use
221 to estimate μ . Finally, note that we neglect the possibility that point mutations arising
222 on IBD segments are removed via gene conversion, as this does not substantially affect
223 the estimates.

224 *Estimating the gene conversion rate*

225 The difference between the mutation rate computed without correcting for gene conversion
226 events and the estimate obtained after removing the effects of gene conversion can be used
227 to quantify the probability that a base pair within IBD segments is involved in a gene
228 conversion event during meiosis. This difference, which we indicate as μ_{GC} , represents the
229 probability of observing a heterozygous site due to existing polymorphisms introduced via
230 gene conversion in a single generation. This rate can be expressed as $\mu_{GC} = p(GC) \times$
231 $p(\theta_p|GC)$, i.e. the product of the probability of a base pair being involved in a gene
232 conversion event, multiplied by the probability of assimilating a heterozygous site given
233 the gene conversion occurs at the locus. The quantity $p(\theta_p|GC)$ can be estimated using

234 the genome-wide heterozygosity of the analyzed sample, and the value of μ_{GC} may be
235 estimated using the previously described correction method. An estimate of $p(GC)$ is
236 therefore obtained as $\hat{p}(GC) = \hat{\mu}_{GC} \times \hat{p}(\theta_p|GC)^{-1}$, and a confidence interval is obtained
237 having computed the standard error of these estimates via weighted block jackknife [43].

238 *Coalescent simulations*

239 We used extensive coalescent simulation to evaluate the proposed methodology. To this
240 end, we used a publicly available coalescent simulator, COSI2 [44], which allows simulat-
241 ing gene conversion events, and our implementation of a coalescent simulator, inspired
242 by the existing GENOME algorithm [45], that enables simulations of a large number of
243 samples and efficiently extract information on IBD segments. The algorithm proceeds
244 backwards in time and, for each individual at generation g , samples a parent at the dis-
245 crete time $g + 1$ in the past, occasionally resulting in coalescent events and sampling a
246 new parent when a recombination event occurs. To speed up computation, the GENOME
247 approach divides the simulated region into relatively large chunks that are not allowed
248 to recombine, discretizing the recombination process and resulting in approximate LD
249 structure at short genomic intervals. The version we developed enables substantial im-
250 provements of memory and run-time requirements, while circumventing the simplifying
251 assumption of non-recombining LD blocks made by the original GENOME algorithm.
252 Briefly, the speed-up over the original algorithm was obtained by sampling recombination

253 breakpoints from an exponential distribution, and by only storing chromosomal regions
254 and individuals that are relevant for the calculation of the ARG at each simulated gener-
255 ation. In addition to this, several improvements to data structures and other algorithmic
256 details were applied. To evaluate our methodology, we further developed an extension
257 of the program that allows efficiently extracting IBD segments from the ancestral re-
258 combination graph without requiring to test differences in shared common ancestors for
259 each marginal tree in the ARG, as done in previous works [29, 38, 31]. We have released
260 open-source software (ARGON) implementing the simulator (see Web Resources).

261 To assess the impact of demographic history on our estimates, we simulated three plau-
262 sible demographic scenarios, in addition to the reconstructed GoNL demographic history.
263 The simulated populations comprised an expanding population that experienced a severe
264 founding event 30 generations before present, a population that undergoes severe expo-
265 nential contraction (referred to as Ashkenazi and Maasai, respectively, due to resemblance
266 with recently studied groups [29]), and an exponentially expanding population (referred
267 to as Europeans, see Figure S3). We used two types of recombination maps to simulate
268 non-uniform recombination rates along the genome. A map where the recombination rate
269 alternates between 1cM/Mb and 2cM/Mb with intervals of 1 Mb, and one where one Mb
270 every five has a six-fold increase in recombination rate over a baseline of 1cM/Mb (Figure
271 S4). To assess the impact of genotyping errors on our methodology, we simulated three
272 types of genotyping errors. We first simulated errors for which a previously unobserved

273 variant is created (“de-novo” errors), or false-positive/negative calls on existing variants.
274 To model frequency-dependent genotyping error rates, we used a beta distribution as a
275 prior for sampling the frequency of planted genotyping errors [46]. For “de-novo” false
276 positive errors, the frequency determines the number of individuals that are affected by
277 an erroneous genotype call. For false-positive/negative genotyping errors, the sampled
278 frequency corresponds to the frequency of the allele that is chosen to add/remove er-
279 roneous genotype calls. Three shape parameters were tested for the beta distribution:
280 $\alpha = 0.01$, $\alpha = 0.5$, resulting in a strong preference for rare variants being erroneously
281 called, and $\alpha = 1$, resulting in a uniform distribution. In both cases the β parameter was
282 set to 1 (see Figure S5). For all simulations, posterior mean estimates for the age of IBD
283 segments were obtained using the coalescent distributions of the simulated models.

284 *GoNL data set*

285 We analyzed sequence data from a recent study of 250 trio families from the Netherlands
286 (the Genome of the Netherlands project [47], GoNL Release 4). The data set consists of
287 748 individuals that passed quality control. The samples were sequenced at an average
288 of 13x using Illumina HiSeq 2000 technology, and variants were called using the GATK
289 UnifiedGenotyper v1.6 software [48]. Trio phasing was performed using MVNcall [49], and
290 additional quality control filters were applied as detailed in [47]. Indels (GoNL Release 5)
291 were detected combining the output of several detection algorithms, as detailed in [47].

292 In addition to the quality control filters applied in the original analysis of the data, we
293 further excluded regions that did not meet several quality criteria derived from the 1000
294 Genomes Project phase 1, as described in [50]. Specifically, we excluded from the analysis
295 (1) low complexity regions; (2) markers that did not pass Hardy Weinberg equilibrium
296 tests; (3) sites with excess coverage; (4) regions where common large insertions were
297 detected; (5) regions not in the strict mask of the 1000 Genomes Project phase 1; (6)
298 segmental duplications of the human genome.

299 Trio-phasing is expected to result in accurate estimation of haploid sequences in the
300 GoNL data. Low frequency variants, in particular, are unlikely to result in doubly het-
301 erozygous parents, so that phasing of rare polymorphisms is generally trivial. Occasional
302 phasing mistakes are however to be expected even in the presence of trios. For all analyses
303 shown in the remainder, we have used a threshold of 1.0 for the MVNCall phasing and
304 genotype calling posterior. This choice resulted in minimized estimated genotyping error,
305 with minimal effects on the estimated mutation rate (see Results, figures S6, S7, S8, and
306 S9)

307 IBD segments and an inferred demographic model were obtained from the analysis
308 described in [47]. For these analyses, only informative variants with very high trio-
309 phasing and genotype calling quality were retained. For IBD detection, markers with
310 frequency below 1% and with MVNCall [49] posterior less than 1.0 were excluded from
311 IBD calculations, resulting in a total of $\sim 3,500,000$ high quality variants. IBD detection

312 was performed on haploid trio-phased individuals using GERMLINE [25] with parameters
313 “-err_hom 2 -err_het 0 -bits 75 -haploid”, i.e. using windows of 75 markers, allowing
314 a maximum of 2 mismatching sites per window to accommodate gene conversion and
315 possible residual genotyping and trio-phasing errors, and retaining only segments of at
316 least 1 cM for the demographic analysis. These parameters were chosen to be the
317 most conservative parameters such that all transmitted haplotypes were detected intact
318 in parent-child pairs along the genome. We further excluded from the analysis genomic
319 regions outside 5 standard deviations from the mean genome-wide sharing, retaining 26
320 chromosomal regions, reported in [47], each longer than 45 cM, for a total of 2,160 cM,
321 and IBD density of 3.07×10^{-3} per site per pair.

322 Demographic inference was performed using the DoRIS software tool [29, 38]. The
323 resulting demographic history is one of exponential expansion, with an ancestral pop-
324 ulation size of 11,500 haploid individuals 150 generations in the past. Two periods of
325 exponential expansion were inferred. The expansion rate between generations 150 and
326 10 was inferred to be 0.0146, followed by a strong expansion in the recent generations at
327 rate 0.479 per generation. Due to the scarcity of extremely recent coalescent events, the
328 magnitude of the latter expansion period is inferred with a high degree of uncertainty,
329 however this was observed to not have appreciable effects on the analysis described in the
330 remainder (see Results).

331 IBD detection is expected to be noisy at the boundaries of detected segments, so that

332 non-IBD regions will be occasionally included in the estimated segments (and sometimes
333 excluded). Because short IBD segments tend to harbor a larger fraction of miscalled non-
334 IBD regions, which increase the observed mismatch rate, an upward bias in the tMRCA
335 regression slope is expected. To cope with this, we have excluded 0.5 cM on either side
336 of the IBD segments from the analysis of mutations and gene conversion rates, as we
337 observed inflation due to noisy boundary estimation plateaus for values larger than this
338 threshold (Figure S10).

339 Note that, due to violations of independence and homoscedasticity assumptions in
340 the performed regressions, all reported standard errors were computed using weighted
341 block jackknife [43], using the 26 independent chromosomal regions obtained as previously
342 described.

343 *Enrichment of deleterious variation in IBD regions*

344 We tested whether mutations arising between the present generation and the MRCA of
345 IBD segments are enriched for deleterious variation. To this end, we ran the ANNOVAR
346 software tool (version “2015Mar22” [51]) on the GoNL variants, and obtained numeric
347 scores for the Polyphen 2 (“ljb23_pp2hvar” [52]) and Gerp++ (“gerp++gt2” [53]) an-
348 notations, restricting the analysis to scores > 2 for the latter. To test for enrichment,
349 we compared the average score of genome-wide variants to the average score of variants
350 found mismatching within IBD regions, treating all variants as independent and reporting

351 Z-test p-values.

352 *Analysis of annotated genomic regions*

353 Several sites along the genome were excluded from the analysis following the filtering cri-
354 teria previously described. In addition, we analyzed mutation rates in specific regions de-
355 scribed in several annotations (e.g. DNase I hypersensitive sites [54, 55, 56], histone mod-
356 ifications [57, 56, 58], constrained genes [4], and several others [59, 60, 61, 62, 61, 63, 64],
357 see Table S1). It is sufficient to neglect regions that fall outside the genomic annotation
358 at hand when computing the observed mismatch rate in the tMRCA regression. Anno-
359 tations that are too small or too clustered in specific regions of the genome may result
360 in downward biases of the estimated mutation rate, due to the “inspection paradox” of
361 the Poisson process underlying the IBD sharing model [29] (Figure S11). This bias was
362 computed and corrected using a permutation procedure (Table S1).

363 Sequence context is an important determinant of mutation rates, and trinucleotide
364 context is often used to account for context-dependent mutation rate variation [65, 66].
365 When analyzing mutation rates within different genomic regions, we computed annotation-
366 specific correction factors to account for the differences in mutation rates that are ex-
367 pected as a result of trinucleotide context variation, using the trinucleotide context-specific
368 mutation-rate matrix of Kryukov [66] (details in Table S1).

369 We additionally derived mutation rates for different mutation categories: CpG/non-

370 CpG and transition/transversions. First, we identified the ancestral allele for all GoNL
371 variants. To do this, we downloaded the ancestral alignment used in the 1000 Genomes
372 project ([67] see Web Resources). The ancestral allele for loci that were not present in
373 this sequence (545,279 out of 12,181,714) was set to the major allele found in the 1000
374 Genomes data set ($N = 300,503$), or set to the allele found in the human reference genome
375 hg19 (see Web Resources) if monomorphic in the 1000 Genomes data set ($N = 244,776$).
376 We then computed mutation rates using MaAF-threshold regression, excluding variants
377 that did not match the analyzed mutation type (e.g. CpG transition), and scaled the
378 resulting rate by the fraction of genome that may harbor the specific kind of mutation
379 (e.g. CpG/non-CpG).

380 **Results**

381 *Simulations*

382 We evaluated the accuracy and robustness of the method via extensive coalescent sim-
383 ulation (see Materials and Methods). To assess the impact of demographic history on
384 our estimates, we simulated several plausible demographic scenarios, and modeled geno-
385 typing errors using a beta distribution with different parameters, specifying error rate at
386 different allele frequencies (Figure S5). We extracted ground truth IBD shared segments
387 from the synthetic ancestral recombination graph, and simulated three types of errors,
388 referred to as de-novo, false positive and false negative errors. De novo errors create

389 erroneous variants in loci that are not truly polymorphic in the sequenced samples. False
390 positive/negative errors affect existing variants, adding or removing derived alleles. To
391 simulate frequency-dependent error rates, we sampled affected alleles using a beta distri-
392 bution to select frequency of the derived allele for all kinds of errors. Inferred mutation
393 rates under all three types of errors are displayed in Figure 3. We observed that tMRCA
394 regression is robust to the presence of substantial levels of de-novo genotyping errors,
395 consistent with the fact that IBD segments of different lengths are equally affected by the
396 spurious sequence mismatches that result from errors of this kind. When false positive
397 genotyping errors were simulated, we observed our approach to be robust to errors up
398 to a rate of $\sim 10^{-5}$ per base pair. False negatives were tolerated up to a frequency of
399 $\sim 10^{-6}$. Very large values of false positive/negative genotyping error rates resulted in a
400 downward bias of the estimates, which is due to the fact that IBD segments of different
401 lengths harbor a slightly different spectrum of mismatching sites, and are therefore not
402 equally likely to be affected by spurious genotype calls (see Figure S1). Similar results
403 were observed for several kinds of genotyping error distribution, demographic model, and
404 recombination map, although the approach proved more robust for error distributions
405 that are less concentrated on very rare variants (figures S3, S4, S5, and S12). The inter-
406 cept of the tMRCA was observed to reflect genotyping error, with average values between
407 1 and 2 times the simulated error rate (Figure S13), depending on the type of error and
408 the parameters of the distribution used to select the frequency of affected alleles.

409 To compare the power of the proposed method to the power of trio-based mutation
410 rate inference, we simulated data at various sample sizes using the GoNL demographic
411 model. Due to the quadratic increase in IBD sharing pairs as sample size increases, the
412 proposed method results in smaller standard errors than the trio-based approach, except
413 at very small sample sizes (Figure 4). However, for demographic models that result in
414 substantial IBD sharing due to a small recent effective population size, higher sample size
415 did not substantially decrease the standard error (Figure S14). This is due to the fact that
416 as new samples are added, early coalescent events result in overlapping ancestral lineages
417 across pairs of individuals, so that limited new information is obtained from increasing
418 the sample size.

419 We finally tested the MaAF-threshold regression approach to correct biases introduced
420 by non-crossover gene conversion events and estimate the probability that a base pair
421 is involved in gene conversion. We simulated realistic mutation and gene conversion
422 rates and used GERMLINE to detect IBD sharing after subsampling synthetic SNPs in
423 order to match the allele frequencies observed in the GoNL data. We observed good
424 performance of the MaAF-threshold-regression in recovering the simulated mutation rate
425 value (Figure 5), and a small downward bias when recovering the gene conversion rate
426 using the GERMLINE IBD discovery parameters used in the real data analysis (Figure
427 S15).

428 *Average genome-wide mutation rate and gene conversion rate in the GoNL data set*

429 We analyzed 498 founders that passed quality control in 250 trio families sequenced within
430 the Genome of the Netherlands (GoNL) project (see Materials and Methods). Due to
431 the trio design of the GoNL study, the average $\sim 13\times$ sequencing depth is effectively
432 doubled to $\sim 26\times$ for the transmitted haplotypes in the 498 analyzed founders. 248 trios
433 and 2 duos passed sequencing quality control. In the remainder, we report results for the
434 analysis of transmitted haplotypes only.

435 We estimated a mutation rate of $2.08 \pm 0.06 \times 10^{-8}$ (Figure 6) before correcting
436 for gene conversion events. For all analyses of mutation and gene conversion rates, we
437 report results for a minimum IBD segment length of 1.6 cM, discarding 0.5 cM on either
438 edge of the segments, and ignoring variants with a trio-phasing and genotyping posterior
439 value less than 1.0. Choosing more conservative values for the minimum length and
440 edge exclusion cutoffs resulted in compatible estimates (Figures S10, S16, and S17). As
441 expected, including variants with lower trio-phasing and genotyping posterior resulted in
442 higher estimates of genotyping error, but negligible effects were observed on the estimates
443 of mutation rate (figures S6, S7, S8, and S9). The tMRCA regression intercept, which
444 reflects genotyping and phasing error rate (see Materials and Methods), was estimated
445 to be $2.21 \pm 0.09 \times 10^{-6}$, within a range that is not expected to result in biases in the
446 tMRCA regression slope based on simulations (figures 3, S6, S9, S12, and S13).

447 We then performed MaAF-threshold regression to correct for gene conversion events
448 (see Materials and Methods, Figure 7). Using this approach, we estimated a genome-
449 wide average mutation rate of $1.66 \pm 0.04 \times 10^{-8}$ per base, per generation. The difference
450 between the corrected and uncorrected estimates, $4.18 \pm 0.48 \times 10^{-9}$, reflects the chance of
451 a heterozygous base pair entering IBD segments during meiosis, as a result of non cross-
452 over gene conversion. This quantity can be used to estimate the chance that a base pair
453 is involved in a gene-conversion tract (see Materials and Methods). Heterozygosity in the
454 GoNL data is observed to be $\sim 6.98 \times 10^{-4}$, consistent with a diploid long-term effective
455 population size in humans of approximately 10,500 individuals. Using this quantity, a
456 base pair is estimated to be involved in a gene conversion event at a rate of $5.99 \pm 0.69 \times$
457 10^{-6} per meiotic event. This rate is in good agreement with a recently published estimate
458 of $5.9 \pm 0.71 \times 10^{-6}$ [39]. We further computed estimates of the rate of mutations at CpG
459 and non-CpG sites (Table S3). The obtained estimates for CpG sites were higher than in
460 previous reports based on trio analysis, consistent with a higher genome-wide rate (see
461 Table 2 in [20]).

462 Because our analysis relies on a reconstructed demographic model, which is used
463 to determine IBD segment age, we performed sensitivity analysis to assess the impact
464 of potential inaccuracies of the reconstructed model on our results (Table S2). The
465 demographic history of the analyzed Dutch samples comprises two consecutive periods
466 of exponential expansion (Figure S3). When a genome-wide average mutation rate was

467 inferred for a demographic model with ancestral population size perturbed by 10%, we
468 observed a $\sim 1.8\%$ difference in the inferred average mutation rate. We observed very
469 limited effects on the mutation rate estimate when perturbing the present-day population
470 size, which is inferred with uncertainty due to the scarcity of very recent coalescent events
471 (Table S2).

472 *Average genome-wide indel rate*

473 Besides measuring the rate of germline point mutation events, the proposed approach
474 allows quantifying additional evolutionary parameters associated with the transmission
475 of IBD chromosomal regions. We applied the same procedure used for inferring mutation
476 rates to infer the rate of $< 20\text{bp}$ indels, which we estimated to be $1.26 \pm 0.06 \times 10^{-9}$.
477 This rate is higher than a recent estimate of 0.68×10^{-9} reported in [68], but compatible
478 with a second recent estimate of $1.5 \pm 0.18 \times 10^{-9}$ [21], both obtained via observation of
479 de-novo events in trios. We additionally used our method to estimate the gene conversion
480 rate based on indels, obtaining a rate of $9.02 \pm 2.91 \times 10^{-6}$ per meiotic event, compatible
481 with the rate obtained from point mutations.

482 *Recombination does not strongly impact mutation rate*

483 We used our approach to analyze annotation-specific mutation rates (see Materials and
484 Methods). We looked for association between recombination rates and mutation rates,

485 a relationship that has been previously detected and attributed to mutagenic properties
486 of recombination [33]. Indeed, we found our tMRCA regression estimates of mutation
487 rate to be strongly associated with recombination rate ($\beta = 0.38 \pm 0.04$ mut/rec, $p =$
488 5.27×10^{-6} , R-squared = 0.9, Figure 8). As previously mentioned, however, increased
489 sequence mismatch rate at loci that undergo frequent recombination may be a result
490 of polymorphic variants introduced by gene conversion events, which may increase the
491 slope of the tMRCA regression. Consistently, after controlling for gene conversion, we
492 observed no significant association between recombination rate and mutation rate ($\beta =$
493 -0.04 ± 0.03 , $p = 0.17$), suggesting the lack of observable mutagenic effects associated with
494 recombination hotspots (figures 8 and S18). A recent study reached similar conclusions
495 [7]. Repeating the same analysis with indels, we detect no significant association between
496 indel rate and recombination rate, for both tMRCA regression slope and MaAF-threshold
497 regression intercept ($\beta = -0.003 \pm 0.003$, $p = 0.37$ after gene conversion correction).

498 *Effects of background selection*

499 Natural selection affecting new mutations may reduce genomic variation, leading to down-
500 ward bias in our mutation rate estimates. Because our analysis is limited to mutation
501 events that occurred in the past ~ 100 generations, due to the length of IBD segments, we
502 expect the effects of natural selection on our genome-wide average mutation rate estimate
503 to be small. Genomic regions with functional or regulatory roles, however, may be under

504 selective pressures that may result in measurable impact even at these short time scales.

505 To estimate the impact of selective pressures on our estimates, we divided the genome
506 based on the B statistic proposed in [5]. The B statistic measures the impact of back-
507 ground selection on a genomic region by estimating the ratio between local effective
508 population size and the effective population size expected under neutrality, so that small
509 values of the B statistic correspond to higher selective pressures (See [5], page 11 for
510 details on the computation of the B statistic). Similarly, a local reduction in effective
511 population size impacts the spectrum of IBD shared segments, which are expected to be
512 longer on average, as a result of early coalescent events in populations of smaller effective
513 size [69, 37]. Indeed, we observed a strong correspondence between small values of the
514 B statistic and the average length of IBD segments ($p = 8.43 \times 10^{-7}$, Figure 9). The
515 effect is, as expected, such that smaller values of the B statistic correspond to longer
516 average IBD shared segments, due to reduced local effective population size. This effect
517 is remarkably strong up to the measured genome-wide average value of the B statistic.
518 We observed longer average IBD segments for large values of the B statistic, a result that
519 may be explained by biases in either of the two measures, or by the fact that additional
520 evolutionary forces, such as selection acting on standing genetic variation [69], are being
521 captured by IBD segment lengths. When we measured the impact of different values of
522 the B statistic on our estimates of mutation rate, however, we found the effect to not
523 be significant ($\beta = 2.17 \pm 1.55 \times 10^{-9}$ mutations per generation, per unit of B statistic,

524 $p = 0.19$, Figure S19). The average genome-wide value for the B statistic was estimated
525 to be 0.78 in the analyzed regions, suggesting a moderate amount of background selection
526 is acting at the average analyzed locus (a value of 1 reflects absence of background selec-
527 tion). If we were to correct the estimated average genome-wide mutation rate to account
528 for this, we would obtain an updated average mutation rate of $1.7 \pm 0.05 \times 10^{-8}$, which
529 is however not significantly different from the estimate obtained without accounting for
530 background selection.

531 *Sequence differences in IBD segments are enriched for deleterious variation*

532 Mutation events occurring within the analyzed IBD regions are expected to have arisen
533 within the past ~ 100 generations, and are therefore on average substantially younger
534 than variants randomly sampled along the genome. Several recent studies have outlined
535 the recent origin of a large fraction of functionally relevant variants [70, 71, 72, 73, 74].
536 We therefore tested whether the presence of recent mutations on IBD segments resulted
537 in an enrichment of deleterious variants compared to the average genome-wide locus,
538 contrasting average scores obtained using Polyphen 2 [52] and Gerp++ [53] annotations
539 (see Materials and Methods). Of the analyzed GoNL variants, 54,960 were annotated
540 using Polyphen 2, and 948,782 were annotated using Gerp++. Of these, 1,843 and 27,900
541 were found mismatching on IBD segments of 1 cM or longer, respectively. When average
542 scores were compared, we found that mismatching sites within IBD regions were strongly

543 enriched for higher scores in both annotations (Table 1, Polyphen2 Z-test $p = 2.8 \times 10^{-5}$;
544 Gerp++ Z-test $p = 9.03 \times 10^{-10}$). We further found a marginal association between
545 Ployphen 2 scores and the B statistic of background selection ($\beta = -0.074 \pm 0.025$,
546 $p = 0.014$, R-squared = 0.39) and a strong association between Gerp++ scores and
547 regional B statistics ($\beta = -0.734 \pm 0.068$, $p = 3.55 \times 10^{-7}$, R-squared = 0.91, Figure
548 S20), which is expected due to both measures relying on metrics related to sequence
549 conservation.

550 We finally tested for enrichment/depletion of the mutation rate in several genomic
551 annotations that have recently been extracted from several studies (e.g. DNase I hyper-
552 sensitive sites [54, 55, 56], histone modifications [57, 56, 58], constrained genes [4], and
553 several others [59, 60, 61, 62, 61, 63, 64]). Annotation-specific mutation rates after cor-
554 recting for the effects of trinucleotide context (see Materials and Methods) are reported
555 in Table S1. None of the annotations were significantly enriched or depleted for muta-
556 tion rates after controlling for trinucleotide context and multiple hypothesis testing. A
557 recent paper [75], found that cell-specific chromatin features are a strong determinant
558 of cancer mutations. Our estimated mutation rate in DNase I hypersensitive regions of
559 $1.66 \pm 0.05 \times 10^{-8}$, on the other hand, suggests that the germline mutation rate is not
560 substantially different from the genome-wide average in these regions, in line with recent
561 analyses [12].

562 Discussion

563 Because mutation events play a key role in the shaping of heritable variation, accurate
564 characterization of this evolutionary parameter is a central task of data-driven analyses
565 of genomic data. In this paper, we proposed a new method to infer mutation rates using
566 distant relatives in a large group of purportedly unrelated individuals, using long chromo-
567 somal regions co-inherited identical-by-descent by present day individuals. The proposed
568 method relies on recently developed coalescent calculations to infer population size his-
569 tory at very recent time scales, which is in turn used to reconstruct the distribution of
570 ages for IBD segments detected in the population. This procedure has several advantages
571 compared to existing other methods. First, by regressing observed sequence mismatch
572 rates on IBD segment age, the estimation is robust to substantial amounts of genotyping
573 error, which is an important confounding factor for many recent estimators of mutation
574 rates based on trio data. Second, by modeling the relationship between population het-
575 erozygosity and allele frequency, it is possible to remove the effects of non-crossover gene
576 conversion, while at the same time estimating its rate. In addition, this approach can
577 be used to estimate the rate of other kinds of mutation events, such as insertions and
578 deletions.

579 We inferred a genome-wide average point mutation rate of $1.66 \pm 0.04 \times 10^{-8}$ per base
580 per generation, or $5.71 \pm 0.14 \times 10^{-10}$ per base per year, having assumed a generation

581 length of 29 years [76]. Recent family-based estimates range within $1.0 - 1.2 \times 10^{-8}$ per
582 base per generation [10, 13, 14], being in most cases significantly lower than the estimate
583 we report here. These methods have the advantage of relying on direct observation of
584 de-novo mutation events, with minimal modeling assumptions, but are affected by the
585 need to rely on strict filtering criteria to deal with false positive/negative genotype calls,
586 which may explain the discrepancy with our results. Phylogenetic methods, on the other
587 hand, fall within $2.0 - 2.5 \times 10^{-8}$ [15, 16]. These estimates rely on several underlying
588 modeling assumptions, which provide a possible explanation for the higher inferred rates,
589 although some have suggested the possibility that these analyses may be capturing the
590 results of evolutionary changes of the mutation rate, or the effects of a varying length of
591 generation times [10, 13, 77]. The estimate of $0.4 - 0.6 \times 10^{-9}$ per year computed in [22]
592 using ancient DNA is slightly lower than our result, but the reported confidence intervals
593 are compatible. Similarly, the rate of $1.4 - 2.3 \times 10^{-8}$ per generation reported in [78],
594 and computed based on point-mutations nearby microsatellites is compatible with our
595 estimate.

596 A contemporary study [23], which is related in spirit to ours, used simulation-based
597 calibration of the decay of heterozygosity along the genome to infer an average genome-
598 wide mutation rate of $1.65 \pm 0.1 \times 10^{-8}$. This closely matches our estimated value.
599 The authors discuss several implications of this mutation rate on our ability to reconcile
600 demographic events inferred using DNA analysis and fossil records, which apply to our

601 analysis as well. Because conversion between sequence divergence and phylogenetic split
602 times across different primate species relies on a per-year mutation rate estimate, different
603 values of this rate have a direct impact on our ability to reconstruct the timing of these
604 events [10, 13]. In general, lower values of the mutation rate (e.g. the pedigree-based
605 “slow” mutation rate) result in older estimates of demographic events, while larger values
606 (e.g. the “fast” phylogenetic rates) result in earlier demographic events. Our inferred
607 “intermediate” mutation rate value of $\sim 5.71 \times 10^{-10}$ per year, consistent with that of [23],
608 implies intermediate scalings for the timing of these events which are generally compatible
609 with fossil records. Assuming no significant effects of generation time and no changes in
610 mutation rates, our estimate implies the split between humans and chimps occurred ~ 6.6
611 million years in the past, and the split between humans and orangutan about ~ 16 m.y.
612 in the past. When our estimate of mutation rate is used to interpret recently reported
613 split times across human populations [11], we find dates that are compatible with what
614 has been reconstructed using methods other than DNA-based reconstruction. The split
615 of African and non-African populations is estimated to have occurred 46 – 61 thousand
616 years in the past, while a split time of 15 thousand years is inferred for the separation
617 of East Asians and Native American populations. These estimates are lower than those
618 obtained assuming a “slow” mutation rate, but do not contradict current fossil evidence.

619 In addition to estimating the rate of point mutations, we report a gene conversion rate
620 of $5.99 \pm 0.69 \times 10^{-6}$, in close agreement with a recent report [39], and find that recom-

621 bination is not associated with mutation rates, supporting recent findings [7]. A recent
622 sperm-typing study further dissected the relationship between mutation, recombination
623 and gene conversion, finding evidence for higher mutational load in regions of high re-
624 combination, which is however contrasted by repairing mechanisms associated with gene
625 conversion [79]. These lead to a higher prevalence of GC alleles compared to AT alleles.
626 Overall, these effects may be counteracting each other in a way that results in minimal
627 differences in the total number of observed mutations in recombination-rich regions, while
628 affecting sequence composition. Interestingly, a recent study reports that recombination
629 rate affects the distribution of putatively deleterious variants along the genome, but found
630 no evidence for a role of biased gene conversion in this observation [80].

631 Finally, our method was applied to estimate the rate of short (< 20 bp) indels, which
632 have not thus far been extensively characterized. We inferred a rate of $1.26 \pm 0.06 \times 10^{-9}$,
633 compatible with two previous estimates of $1.5 \pm 0.18 \times 10^{-9}$ [21], and $1.06 \pm 0.1 \times 10^{-9}$
634 [81], but higher than the estimate of 0.68×10^{-9} reported in [68]. While these analyses are
635 likely affected by difficulties related to detection of short indels, collectively they suggest
636 that insertion and deletion mutation events occur at a significantly lower rate compared
637 to single point mutations.

638 In addition to analyzing genome-wide average rates, we looked for enrichment or
639 depletion of mutation rates in a number of genomic annotations that were recently derived
640 from several studies. Although we cannot exclude significant deviations from genome-wide

641 averages, we found no evidence for changes in overall mutation rates for the analyzed
642 regions. Notably, we observed that while the distribution of IBD shared haplotypes
643 closely reflects the effects of background selection along the genome, a negligible effect
644 is observed on our estimated mutation rates, suggesting that estimating mutation rates
645 using mutation events under the effects of ~ 100 generations of natural selection does not
646 significantly bias local mutation rate estimates in European populations. Consistent with
647 the idea that mutations on IBD segments are recent and under the effects of selective
648 forces [70, 71, 72, 73, 74], we found a strong enrichment for deleterious variants within
649 IBD regions.

650 Our method provides a new way of studying mutation and gene conversion events in
651 large samples of unrelated individuals, being robust to substantial amounts of genotyping
652 error, which constitutes a significant confounder in trio-based analyses. The main limita-
653 tion of our approach, however, is the need to rely on two fundamental components that
654 are potential sources of bias, namely detection of shared IBD segments and the need to in-
655 fer the recent demographic history for the analyzed population. Our analysis of mutation
656 rates in the GoNL dataset relies on IBD detection and demographic inference performed
657 in a previous study [47], but it is possible that additional sources of uncertainty in these
658 two components affect our results. Our conservative exclusion of substantial portions
659 of IBD segments, together with our sensitivity analysis for changes in the demographic
660 model, however, suggest that these biases, if present, should not be substantial.

661 Several potential directions for improvement of the proposed methodology and anal-
662 ysis may be outlined. First, additional developments of the coalescent calculations used
663 in this work may remove the requirement of estimating a demographic model for the
664 analyzed samples, as described in [82]. These methods need to be adapted in order to
665 deal with genotyping error and to deconvolute the contribution of mutation and gene con-
666 version to observed sequence mismatches in IBD regions. Second, it may be possible to
667 devise improved approaches for dealing with heteroscedasticity and the dependence across
668 observations in the tMRCA and MaAF-threshold regressions, which do not result in bi-
669 ases, but may reduce the efficiency of our estimators. In addition, the MaAF-threshold
670 regression currently assumes an underlying model of constant effective population size,
671 an approximation that only has small effects on our estimates (Figure S2), but may be
672 removed via additional modeling. Third, as shown in our simulation-based evaluation of
673 the method, the relationship between allele frequency and genotyping error rates limits
674 the robustness of our method to very large amounts of noise in the analyzed sequences.
675 Alternative genotype calling strategies may be employed to reduce these effects, e.g. geno-
676 type calling approaches that do not rely on whether variants are observed polymorphic
677 in other sequenced individuals. Finally, by applying the proposed tMRCA regression, it
678 may be possible to analyze multi-generation pedigrees while controlling for substantial
679 genotyping error. Detection of very long IBD segments is in fact trivial when closely
680 related individuals are analyzed, and the age of these segments is fully determined by

681 their known genealogical relationship. In addition, the proposed methods can be applied
682 to increasingly large datasets that are currently being produced, which may be reliably
683 phased to enable accurate IBD detection even in the absence of sequenced family members
684 [24].

685 **Appendix**

686 *The age of IBD segments*

687 If a pair of chromosomes find a common ancestor at time t generations before present,
688 the probability that a single site is spanned by an IBD segment of length l at least u
689 Morgans can be expressed as

$$\begin{aligned}\sigma(t) &= \int_u^\infty l(2t)^2 e^{-2tl} dl \\ &= e^{-2tu}(2tu + 1)\end{aligned}\tag{1}$$

690 The distribution $l(2t)^2 e^{-2tl}$ represents the sum of two exponential random variables
691 with parameter $2t$, which is the rate at which a recombination occurs on either side of the
692 chosen site. Note that this assumes an IBD segment is delimited by the occurrence of re-
693 combination events, which is equivalent to assuming an underlying sequentially Markovian
694 coalescent (SMC) model. For very short IBD segments (e.g. < 0.3 cM) and in popula-
695 tions that experience substantial and long-lasting isolation (e.g. $N_e < 1,000$), the slightly
696 more complex SMC' model [83] provides more accurate calculations [84, 85, 11, 86]. This

697 is however unnecessary given the demographic history and length ranges here considered.
698 It follows from the linearity of the expectation operator that the expected fraction $f(t)$
699 of genome shared IBD for a pair of individuals whose ancestral lineages coalesce at time
700 t can be obtained from the probability that a single site is spanned by an IBD segment
701 of length at least u Morgans, which we write $f(t) = \sigma(t)$. The expected length of an IBD
702 segment transmitted from a common ancestor living at time t is therefore

$$\ell(t) = \int_u^\infty l \times 2te^{2t(u-l)} dl = 1/(2t) + u \quad (2)$$

703 To obtain the expected number of IBD segments obtained if the lineages of two indi-
704 viduals coalesce at time t , we therefore divide the expected total amount of genome shared
705 IBD by the expected length of an IBD segment co-inherited from an ancestor living at
706 time t . This yields

$$\begin{aligned} n_u(t) &= Lf_u(t)/\ell(t) \\ &= \frac{Le^{-2tu}(2tu + 1)}{1/(2t) + u} \\ &= L2e^{-2tu}t, \end{aligned} \quad (3)$$

707 where L is the size, in Morgans, of the considered genomic region. To obtain the expected
708 number of IBD segments longer than u Morgans for the average pair of individuals in the
709 population, we marginalize over the distribution of pairwise coalescence times, $c(t)$, which
710 depends on the demographic history,

$$n_u = \int_0^{\infty} c(t) n_u(t) dt. \quad (4)$$

711 This quantity has a closed form expression if we assume that the population size
712 becomes constant at an arbitrarily remote point in time, and can be used to obtain the
713 posterior age distribution of IBD segment ages,

$$p_u(t) = \frac{c(t) n_u(t)}{n_u}. \quad (5)$$

714 *Contribution of individual variants to heterozygosity*

715 For a sample of K homologous sequences from a population, the heterozygosity per site
716 can be estimated by computing

$$\hat{\theta} = \frac{1}{s} \sum_{i=1}^s \frac{K}{K-1} 2 \frac{x_i}{K} \left(1 - \frac{x_i}{K}\right) \quad (6)$$

717 [87], where s is the number of sites in each sequence, x_i is the number of samples carrying
718 a derived allele at site i , and $K/(K-1)$ is a bias-correction factor. Defining $M(x)$ as the
719 total number of sites in the sample for which exactly x sequences carry a derived allele,
720 we can rewrite this equation as a sum over x :

$$\hat{\theta} = \sum_{x=1}^{K-1} \frac{M(x)}{s} \frac{2x(K-x)}{K(K-1)}. \quad (7)$$

721 The term $M(x)/s$ is the proportion of sites at which x of the K sequences carry a
722 derived allele and the term

$$\frac{2x(K-x)}{K(K-1)} \quad (8)$$

723 is the probability of discovering such a polymorphic site when just two sequences are
724 sampled without replacement from the K sequences. Note that this probability is the
725 same for sites with x copies of a derived allele as it is for sites with $K-x$ copies. Thus,
726 we may also write

$$\hat{\theta} = \sum_{x=1}^{\lfloor K/2 \rfloor} \frac{M(x) + M(K-x)}{s} \frac{2x(K-x)}{K(K-1)}, \quad (9)$$

727 where $\lfloor K/2 \rfloor$ is the largest integer that is less than or equal to $K/2$ and x is now the
728 count of the minor allele.

729 This allows us to consider the average contribution of different kinds of polymorphic
730 sites to overall heterozygosity. Under the constant population size, neutral model of [88],

$$E[M(x)] = \frac{s\theta}{x} \quad (10)$$

[89, 90] in which $\theta = 4N\mu$ is the diploid population-scaled mutation rate per site, or the expected per-site heterozygosity of the population. Using (10) together with (9) and simplifying gives

$$E[\hat{\theta}] = \sum_{x=1}^{(K-1)/2} \theta \frac{2}{K-1}, \quad (11)$$

731 in which we have assumed that K is odd for simplicity. The sum in (11) evaluates to θ ,
732 as expected for an unbiased estimator.

733 Equation (11) shows that on average the different kinds of polymorphic sites, catego-
734 rized by minor allele frequency, contribute uniformly to heterozygosity, as noted previously
735 by [91]. Another way of stating this is that polymorphisms discovered by screening in
736 samples of size two will be uniformly distributed among minor-allele frequency classes.
737 This depends on the population size being constant over time, and is also not true for
738 derived allele-frequency classes. Figure S2 shows that contributions to heterozygosity are
739 close to uniform for the GoNL site frequency spectrum even though the population size
740 has not been constant.

741 **Supplemental data**

742 Supplemental Data include 20 figures and 3 tables.

743 **Acknowledgments**

744 We express our gratitude to Nick Patterson, Priya Moorjani, Mark Lipson and Amy
745 Williams for useful scientific discussions and comments on an early draft, and to Ilya
746 Shlyakhter for support with the COSI2 simulator. This research was funded by NIH
747 grant R01 MH101244.

748 **Web Resources**

749 The developed coalescent simulator (ARGON) will be made available at <https://github.com/pierpal/ARGON>.

751 The tool to infer mutation and gene conversion rates (IBDMUT) will be made available
752 at <https://github.com/pierpal/IBDMUT>.

753 The 1000 Genomes ancestral alignments were downloaded from ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/phase1/analysis_results/supporting/ancestral_alignments/.

755 The human reference sequence (h19) was downloaded from <http://hgdownload.cse.ucsc.edu/goldenPath/hg19/chromosomes/>.

757 **Consortia**

758 The members of the Genome of the Netherlands Consortium are Laurent C Francioli,
759 Androniki Menelaou, Sara L Pulit, Freerk van Dijk, Pier Francesco Palamara, Clara C
760 Elbers, Pieter B T Neerincx, Kai Ye, Victor Guryev, Wigard P Kloosterman, Patrick
761 Deelen, Abdel Abdellaoui, Elisabeth M van Leeuwen, Mannis van Oven, Martijn Ver-
762 maat, Mingkun Li, Jeroen F J Laros, Lennart C Karssen, Alexandros Kanterakis, Najaf
763 Amin, Jouke Jan Hottenga, Eric-Wubbo Lameijer, Mathijs Kattenberg, Martijn Dijkstra,
764 Heorhiy Byelas, Jessica van Setten, Barbera D C van Schaik, Jan Bot, Isac J Nijman, Ivo
765 Renkens, Tobias Marschall, Alexander Schnhuth, Jayne Y Hehir-Kwa, Robert E Hand-

766 saker, Paz Polak, Mashaal Sohail, Dana Vuzman, Fereydoun Hormozdiari, David van
767 Enckevort, Hailiang Mei, Vyacheslav Koval, Matthijs H Moed, K Joeri van der Velde,
768 Fernando Rivadeneira, Karol Estrada, Carolina Medina-Gomez, Aaron Isaacs, Steven A
769 McCarroll, Marian Beekman, Anton J M de Craen, H Eka D Suchiman, Albert Hofman,
770 Ben Oostra, Andr G Uitterlinden, Gonneke Willemsen, LifeLines Cohort Study, Math-
771 ieu Platteel, Jan H Veldink, Leonard H van den Berg, Steven J Pitts, Shobha Potluri,
772 Purnima Sundar, David R Cox, Shamil R Sunyaev, Johan T den Dunnen, Mark Stonek-
773 ing, Peter de Knijff, Manfred Kayser, Qibin Li, Yingrui Li, Yuanping Du, Ruoyan Chen,
774 Hongzhi Cao, Ning Li, Sujie Cao, Jun Wang, Jasper A Bovenberg, Itsik Pe'er, P Eline
775 Slagboom, Cornelia M van Duijn, Dorret I Boomsma, Gert-Jan B van Ommen, Paul I W
776 de Bakker, Morris A Swertz & Cisca Wijmenga

777 **References**

- 778 [1] Crow, J. F. (2000). The origins, patterns and implications of human spontaneous
779 mutation. *Nature Reviews Genetics* *1*, 40–47.
- 780 [2] Arnheim, N. and Calabrese, P. (2009). Understanding what determines the frequency
781 and pattern of human germline mutations. *Nature Reviews Genetics* *10*, 478–488.
- 782 [3] Michaelson, J. J., Shi, Y., Gujral, M., Zheng, H., Malhotra, D., Jin, X., Jian, M.,
783 Liu, G., Greer, D., Bhandari, A., et al. (2012). Whole-genome sequencing in autism
784 identifies hot spots for de novo germline mutation. *Cell* *151*, 1431–1442.
- 785 [4] Samocha, K. E., Robinson, E. B., Sanders, S. J., Stevens, C., Sabo, A., McGrath,
786 L. M., Kosmicki, J. A., Rehnström, K., Mallick, S., Kirby, A., et al. (2014). A
787 framework for the interpretation of de novo mutation in human disease. *Nature*
788 *Genetics* *46*, 944–950.
- 789 [5] McVicker, G., Gordon, D., Davis, C., and Green, P. (2009). Widespread genomic
790 signatures of natural selection in hominid evolution. *PLoS Genetics* *5*, e1000471.
- 791 [6] Gronau, I., Arbiza, L., Mohammed, J., and Siepel, A. (2013). Inference of natural
792 selection from interspersed genomic elements based on polymorphism and divergence.
793 *Molecular Biology and Evolution* *30*, 1159–1171.

- 794 [7] Schaibley, V. M., Zawistowski, M., Wegmann, D., Ehm, M. G., Nelson, M. R., Jean,
795 P. L. S., Abecasis, G. R., Novembre, J., Zöllner, S., and Li, J. Z. (2013). The
796 influence of genomic context on mutation patterns in the human genome inferred
797 from rare variants. *Genome Research* *23*, 1974–1984.
- 798 [8] Presgraves, D. C. and Soojin, V. Y. (2009). Doubts about complex speciation
799 between humans and chimpanzees. *Trends in ecology & evolution* *24*, 533–540.
- 800 [9] Li, H. and Durbin, R. (2011). Inference of human population history from individual
801 whole-genome sequences. *Nature* *475*, 493–496.
- 802 [10] Scally, A. and Durbin, R. (2012). Revising the human mutation rate: implications
803 for understanding human evolution. *Nature Reviews Genetics* *13*, 745–753.
- 804 [11] Schiffels, S. and Durbin, R. (2014). Inferring human population size and separation
805 history from multiple genome sequences. *Nature Genetics* *46*, 919–925.
- 806 [12] Francioli, L. C., Polak, P. P., Koren, A., Menelaou, A., Chun, S., Renkens, I., van
807 Duijn, C. M., Swertz, M., Wijmenga, C., van Ommen, G., et al. (2015). Genome-wide
808 patterns and properties of de novo mutations in humans. *Nature Genetics*.
- 809 [13] Ségurel, L., Wyman, M. J., and Przeworski, M. (2014). Determinants of mutation
810 rate variation in the human germline. *Annual Review of Genomics and Human*
811 *Genetics* *15*, 47–70.

- 812 [14] Campbell, C. D. and Eichler, E. E. (2013). Properties and rates of germline mutations
813 in humans. *Trends in Genetics* *29*, 575–584.
- 814 [15] Nachman, M. W. and Crowell, S. L. (2000). Estimate of the mutation rate per
815 nucleotide in humans. *Genetics* *156*, 297–304.
- 816 [16] Sequencing, T. C., Consortium, A., et al. (2005). Initial sequence of the chimpanzee
817 genome and comparison with the human genome. *Nature* *437*, 69–87.
- 818 [17] Roach, J. C., Glusman, G., Smit, A. F., Huff, C. D., Hubley, R., Shannon, P. T.,
819 Rowen, L., Pant, K. P., Goodman, N., Bamshad, M., et al. (2010). Analysis of
820 genetic inheritance in a family quartet by whole-genome sequencing. *Science* *328*,
821 636–639.
- 822 [18] Conrad, D. F., Keebler, J. E., DePristo, M. A., Lindsay, S. J., Zhang, Y., Casals, F.,
823 Idaghdour, Y., Hartl, C. L., Torroja, C., Garimella, K. V., et al. (2011). Variation
824 in genome-wide mutation rates within and between human families. *Nature* *43*,
825 712–714.
- 826 [19] Campbell, C. D., Chong, J. X., Malig, M., Ko, A., Dumont, B. L., Han, L., Vives,
827 L., O’Roak, B. J., Sudmant, P. H., Shendure, J., et al. (2012). Estimating the
828 human mutation rate using autozygosity in a founder population. *Nature Genetics*
829 *44*, 1277–1281.

- 830 [20] Kong, A., Frigge, M. L., Masson, G., Besenbacher, S., Sulem, P., Magnusson, G.,
831 Gudjonsson, S. A., Sigurdsson, A., Jonasdottir, A., Jonasdottir, A., et al. (2012).
832 Rate of de novo mutations and the importance of fathers age to disease risk. *Nature*
833 *488*, 471–475.
- 834 [21] Besenbacher, S., Liu, S., Izarzugaza, J., Grove, J., Belling, K., Bork-Jensen, J.,
835 Huang, S., Als, T., Li, S., Yadav, R., et al. (2015). Novel variation and de novo
836 mutation rates in population-wide de novo assembled danish trios. *Nature Commu-*
837 *nications 6*, 5969.
- 838 [22] Fu, Q., Li, H., Moorjani, P., Jay, F., Slepchenko, S. M., Bondarev, A. A., Johnson,
839 P. L., Aximu-Petri, A., Prüfer, K., de Filippo, C., et al. (2014). Genome sequence
840 of a 45,000-year-old modern human from western siberia. *Nature 514*, 445–449.
- 841 [23] Lipson, M., Loh, P.-R., Sankararaman, S., Patterson, N., Berger, B., and Reich, D.
842 (2015). Calibrating the human mutation rate via ancestral recombination density in
843 diploid genomes. bioRxiv pp. 015560.
- 844 [24] Kong, A., Masson, G., Frigge, M. L., Gylfason, A., Zusmanovich, P., Thorleifsson,
845 G., Olason, P. I., Ingason, A., Steinberg, S., Rafnar, T., et al. (2008). Detection of
846 sharing by descent, long-range phasing and haplotype imputation. *Nature Genetics*
847 *40*, 1068–1075.

- 848 [25] Gusev, A., Lowe, J. K., Stoffel, M., Daly, M. J., Altshuler, D., Breslow, J. L.,
849 Friedman, J. M., and Pe'er, I. (2009). Whole population, genome-wide mapping of
850 hidden relatedness. *Genome Research* *19*, 318–326.
- 851 [26] Powell, J. E., Visscher, P. M., and Goddard, M. E. (2010). Reconciling the analysis
852 of ibd and ibs in complex trait studies. *Nature Reviews Genetics* *11*, 800–805.
- 853 [27] Browning, S. R. and Browning, B. L. (2012). Identity by descent between distant
854 relatives: detection and applications. *Annual review of genetics* *46*, 617–633.
- 855 [28] Browning, S. R. and Thompson, E. A. (2012). Detecting rare variant associations
856 by identity-by-descent mapping in case-control studies. *Genetics* *190*, 1521–1531.
- 857 [29] Palamara, P. F., Lencz, T., Darvasi, A., and Peer, I. (2012). Length distributions of
858 identity by descent reveal fine-scale demographic history. *The American Journal of*
859 *Human Genetics* *91*, 809–822.
- 860 [30] Ralph, P. and Coop, G. (2013). The geography of recent genetic ancestry across
861 europe. *PLoS Biology* *11*, e1001555.
- 862 [31] Browning, B. L. and Browning, S. R. (2013). Detecting identity by descent and
863 estimating genotype error rates in sequence data. *The American Journal of Human*
864 *Genetics* *93*, 840–851.

- 865 [32] Gudbjartsson, D. F., Sulem, P., Helgason, H., Gylfason, A., Gudjonsson, S. A.,
866 Zink, F., Oddson, A., Magnusson, G., Halldorsson, B. V., Hjartarson, E., et al.
867 (2015). Sequence variants from whole genome sequencing a large group of icelanders.
868 *Scientific Data* 2, EP.
- 869 [33] Hellmann, I., Ebersberger, I., Ptak, S. E., Pääbo, S., and Przeworski, M. (2003).
870 A neutral explanation for the correlation of diversity with recombination rates in
871 humans. *The American Journal of Human Genetics* 72, 1527–1535.
- 872 [34] Purcell, S., Neale, B., Todd-Brown, K., Thomas, L., Ferreira, M. A., Bender, D.,
873 Maller, J., Sklar, P., De Bakker, P. I., Daly, M. J., et al. (2007). Plink: a tool set
874 for whole-genome association and population-based linkage analyses. *The American*
875 *Journal of Human Genetics* 81, 559–575.
- 876 [35] Browning, B. L. and Browning, S. R. (2013). Improving the accuracy and efficiency
877 of identity-by-descent detection in population data. *Genetics* 194, 459–471.
- 878 [36] Rodriguez, J. M., Bercovici, S., Huang, L., Frostig, R., and Batzoglou, S. (2015).
879 Parente2: a fast and accurate method for detecting identity by descent. *Genome*
880 *research* 25, 280–289.
- 881 [37] Gusev, A., Palamara, P. F., Aponte, G., Zhuang, Z., Darvasi, A., Gregersen, P., and
882 Pe'er, I. (2012). The architecture of long-range haplotypes shared within and across
883 populations. *Molecular Biology and Evolution* 29, 473–486.

- 884 [38] Palamara, P. F. and Pe'er, I. (2013). Inference of historical migration rates via
885 haplotype sharing. *Bioinformatics* *29*, i180–i188.
- 886 [39] Williams, A., Genevrese, G., Dyer, T., Truax, K., Jun, G., Patterson, N., Curran,
887 J. E., Duggirala, R., Blangero, J., Reich, D., et al. (2014). Non-crossover gene
888 conversions show strong gc bias and unexpected clustering in humans. bioRxiv pp.
889 009175.
- 890 [40] Wiuf, C. and Hein, J. (2000). The coalescent with gene conversion. *Genetics* *155*,
891 451–462.
- 892 [41] Odenthal-Hesse, L., Berg, I. L., Veselis, A., Jeffreys, A. J., and May, C. A. (2014).
893 Transmission distortion affecting human noncrossover but not crossover recombina-
894 tion: a hidden source of meiotic drive. *PLoS Genetics* *10*, e1004106.
- 895 [42] Jeffreys, A. J. and May, C. A. (2004). Intense and highly localized gene conversion
896 activity in human meiotic crossover hot spots. *Nature Genetics* *36*, 151–156.
- 897 [43] Busing, F. M., Meijer, E., and Van Der Leeden, R. (1999). Delete-m jackknife for
898 unequal m. *Statistics and Computing* *9*, 3–8.
- 899 [44] Shlyakhter, I., Sabeti, P. C., and Schaffner, S. F. (2014). Cosi2: An efficient simulator
900 of exact and approximate coalescent with selection. *Bioinformatics* *30*, 3427–3429.

- 901 [45] Liang, L., Zöllner, S., and Abecasis, G. R. (2007). Genome: a rapid coalescent-based
902 whole genome simulator. *Bioinformatics* *23*, 1565–1567.
- 903 [46] He, Z., Li, X., Ling, S., Fu, Y.-X., Hungate, E., Shi, S., and Wu, C.-I. (2013).
904 Estimating dna polymorphism from next generation sequencing data with high error
905 rate by dual sequencing applications. *BMC Genomics* *14*, 535.
- 906 [47] The Genome of the Netherlands Consortium. (2014). Whole-genome sequence varia-
907 tion, population structure and demographic history of the dutch population. *Nature*
908 *Genetics* *46*, 818–825.
- 909 [48] McKenna, A., Hanna, M., Banks, E., Sivachenko, A., Cibulskis, K., Kernytzky, A.,
910 Garimella, K., Altshuler, D., Gabriel, S., Daly, M., et al. (2010). The genome analysis
911 toolkit: a mapreduce framework for analyzing next-generation dna sequencing data.
912 *Genome Research* *20*, 1297–1303.
- 913 [49] Menelaou, A. and Marchini, J. (2013). Genotype calling and phasing using next-
914 generation sequencing reads and a haplotype scaffold. *Bioinformatics* *29*, 84–91.
- 915 [50] Genovese, G., Kähler, A. K., Handsaker, R. E., Lindberg, J., Rose, S. A., Bakhoun,
916 S. F., Chambert, K., Mick, E., Neale, B. M., Fromer, M., et al. (2014). Clonal
917 hematopoiesis and blood-cancer risk inferred from blood dna sequence. *New England*
918 *Journal of Medicine* *371*, 2477–2487.

- 919 [51] Wang, K., Li, M., and Hakonarson, H. (2010). Annovar: functional annotation of
920 genetic variants from high-throughput sequencing data. *Nucleic acids research* *38*,
921 e164–e164.
- 922 [52] Adzhubei, I. A., Schmidt, S., Peshkin, L., Ramensky, V. E., Gerasimova, A., Bork,
923 P., Kondrashov, A. S., and Sunyaev, S. R. (2010). A method and server for predicting
924 damaging missense mutations. *Nature methods* *7*, 248–249.
- 925 [53] Davydov, E. V., Goode, D. L., Sirota, M., Cooper, G. M., Sidow, A., and Batzoglou,
926 S. (2010). Identifying a high fraction of the human genome to be under selective
927 constraint using *gerp++*. *PLoS computational biology* *6*, e1001025.
- 928 [54] Maurano, M. T., Humbert, R., Rynes, E., Thurman, R. E., Haugen, E., Wang,
929 H., Reynolds, A. P., Sandstrom, R., Qu, H., Brody, J., et al. (2012). Systematic
930 localization of common disease-associated variation in regulatory dna. *Science* *337*,
931 1190–1195.
- 932 [55] Andersson, R., Gebhard, C., Miguel-Escalada, I., Hoof, I., Bornholdt, J., Boyd, M.,
933 Chen, Y., Zhao, X., Schmidl, C., Suzuki, T., et al. (2014). An atlas of active
934 enhancers across human cell types and tissues. *Nature* *507*, 455–461.
- 935 [56] Trynka, G., Westra, H.-J., Slowikowski, K., Hu, X., Xu, H., Stranger, B. E., Han,
936 B., and Raychaudhuri, S. (2014). Disentangling effects of colocalizing genomic

937 annotations to functionally prioritize non-coding variants within complex trait loci.

938 bioRxiv pp. 009258.

939 [57] Hnisz, D., Abraham, B. J., Lee, T. I., Lau, A., Saint-André, V., Sigova, A. A., Hoke,
940 H. A., and Young, R. A. (2013). Super-enhancers in the control of cell identity and
941 disease. *Cell* *155*, 934–947.

942 [58] of the Psychiatric Genomics Consortium, S. W. G. et al. (2014). Biological insights
943 from 108 schizophrenia-associated genetic loci. *Nature* *511*, 421–427.

944 [59] Gusev, A., Lee, S. H., Trynka, G., Finucane, H., Vilhjálmsón, B. J., Xu, H., Zang,
945 C., Ripke, S., Bulik-Sullivan, B., Stahl, E., et al. (2014). Partitioning heritability of
946 regulatory and cell-type-specific variants across 11 common diseases. *The American*
947 *Journal of Human Genetics* *95*, 535–552.

948 [60] Ward, L. D. and Kellis, M. (2012). Evidence of abundant purifying selection in
949 humans for recently acquired regulatory functions. *Science* *337*, 1675–1678.

950 [61] Consortium, E. P. et al. (2012). An integrated encyclopedia of dna elements in the
951 human genome. *Nature* *489*, 57–74.

952 [62] Hoffman, M. M., Ernst, J., Wilder, S. P., Kundaje, A., Harris, R. S., Libbrecht, M.,
953 Giardine, B., Ellenbogen, P. M., Bilmes, J. A., Birney, E., et al. (2012). Integrative

- 954 annotation of chromatin elements from encode data. *Nucleic acids research* pp.
955 gks1284.
- 956 [63] Sankararaman, S., Mallick, S., Dannemann, M., Prüfer, K., Kelso, J., Pääbo, S.,
957 Patterson, N., and Reich, D. (2014). The genomic landscape of neanderthal ancestry
958 in present-day humans. *Nature* *507*, 354–357.
- 959 [64] Cabili, M. N., Trapnell, C., Goff, L., Koziol, M., Tazon-Vega, B., Regev, A., and
960 Rinn, J. L. (2011). Integrative annotation of human large intergenic noncoding rnas
961 reveals global properties and specific subclasses. *Genes & development* *25*, 1915–
962 1927.
- 963 [65] Krawczak, M., Ball, E. V., and Cooper, D. N. (1998). Neighboring-nucleotide effects
964 on the rates of germ-line single-base-pair substitution in human genes. *The American*
965 *Journal of Human Genetics* *63*, 474–488.
- 966 [66] Kryukov, G. V., Pennacchio, L. A., and Sunyaev, S. R. (2007). Most rare missense
967 alleles are deleterious in humans: implications for complex disease and association
968 studies. *The American Journal of Human Genetics* *80*, 727–739.
- 969 [67] Consortium, . G. P. et al. (2012). An integrated map of genetic variation from 1,092
970 human genomes. *Nature* *491*, 56–65.

- 971 [68] Kloosterman, W. P., Francioli, L. C., Hormozdiari, F., Marschall, T., Hehir-Kwa,
972 J. Y., Abdellaoui, A., Lameijer, E.-W., Moed, M. H., Koval, V., Renkens, I., et al.
973 (2015). Characteristics of de novo structural changes in the human genome. *Genome*
974 *Research*.
- 975 [69] Albrechtsen, A., Moltke, I., and Nielsen, R. (2010). Natural selection and the
976 distribution of identity-by-descent in the human genome. *Genetics* *186*, 295–308.
- 977 [70] Tennessen, J. A., Bigham, A. W., OConnor, T. D., Fu, W., Kenny, E. E., Gravel, S.,
978 McGee, S., Do, R., Liu, X., Jun, G., et al. (2012). Evolution and functional impact
979 of rare coding variation from deep sequencing of human exomes. *Science* *337*, 64–69.
- 980 [71] Nelson, M. R., Wegmann, D., Ehm, M. G., Kessner, D., Jean, P. S., Verzilli, C.,
981 Shen, J., Tang, Z., Bacanu, S.-A., Fraser, D., et al. (2012). An abundance of rare
982 functional variants in 202 drug target genes sequenced in 14,002 people. *Science* *337*,
983 100–104.
- 984 [72] Fu, W., OConnor, T. D., Jun, G., Kang, H. M., Abecasis, G., Leal, S. M., Gabriel,
985 S., Rieder, M. J., Altshuler, D., Shendure, J., et al. (2013). Analysis of 6,515
986 exomes reveals the recent origin of most human protein-coding variants. *Nature* *493*,
987 216–220.
- 988 [73] Kiezun, A., Pulit, S. L., Francioli, L. C., van Dijk, F., Swertz, M., Boomsma, D. I.,
989 van Duijn, C. M., Slagboom, P. E., van Ommen, G., Wijmenga, C., et al. (2013).

- 990 Deleterious alleles in the human genome are on average younger than neutral alleles
991 of the same frequency. *PLoS Genetics* *9*, e1003301.
- 992 [74] Gudbjartsson, D. F., Helgason, H., Gudjonsson, S. A., Zink, F., Oddson, A., Gylfa-
993 son, A., Besenbacher, S., Magnusson, G., Halldorsson, B. V., Hjartarson, E., et al.
994 (2015). Large-scale whole-genome sequencing of the icelandic population. *Nature*
995 *Genetics* *47*, 435–444.
- 996 [75] Polak, P., Karlić, R., Koren, A., Thurman, R., Sandstrom, R., Lawrence, M. S.,
997 Reynolds, A., Rynes, E., Vlahoviček, K., Stamatoyannopoulos, J. A., et al. (2015).
998 Cell-of-origin chromatin organization shapes the mutational landscape of cancer.
999 *Nature* *518*, 360–364.
- 1000 [76] Fenner, J. N. (2005). Cross-cultural estimation of the human generation interval for
1001 use in genetics-based population divergence studies. *American journal of physical*
1002 *anthropology* *128*, 415–423.
- 1003 [77] Harris, K. (2015). Evidence for recent, population-specific evolution of the human
1004 mutation rate. *Proceedings of the National Academy of Sciences* *112*, 3439–3444.
- 1005 [78] Sun, J. X., Helgason, A., Masson, G., Ebenesersdóttir, S. S., Li, H., Mallick, S.,
1006 Gnerre, S., Patterson, N., Kong, A., Reich, D., et al. (2012). A direct characterization
1007 of human mutation based on microsatellites. *Nature Genetics* *44*, 1161–1165.

- 1008 [79] Arbeithuber, B., Betancourt, A. J., Ebner, T., and Tiemann-Boege, I. (2015).
1009 Crossovers are associated with mutation and biased gene conversion at recombination
1010 hotspots. *Proceedings of the National Academy of Sciences* *112*, 2109–2114.
- 1011 [80] Hussin, J. G., Hodgkinson, A., Idaghdour, Y., Grenier, J.-C., Goulet, J.-P., Gbeha,
1012 E., Hip-Ki, E., and Awadalla, P. (2015). Recombination affects accumulation of
1013 damaging and disease-associated mutations in human populations. *Nature Genetics*
1014 *47*, 400–404.
- 1015 [81] Ramu, A., Noordam, M. J., Schwartz, R. S., Wuster, A., Hurles, M. E., Cartwright,
1016 R. A., and Conrad, D. F. (2013). Denovogear: de novo indel and point mutation
1017 discovery and phasing. *Nature methods* *10*, 985–987.
- 1018 [82] Palamara, P. F. (2014). Population genetics of identity by descent. PhD thesis (New
1019 York City: Columbia University).
- 1020 [83] Marjoram, P. and Wall, J. D. (2006). Fast "coalescent" simulation. *BMC Genetics*
1021 *7*, 16.
- 1022 [84] Harris, K. and Nielsen, R. (2013). Inferring demographic history from a spectrum
1023 of shared haplotype lengths. *PLoS Genetics* *9*, e1003521.

- 1024 [85] Hobolth, A. and Jensen, J. L. (2014). Markovian approximation to the finite loci
1025 coalescent with recombination along multiple sequences. *Theoretical population bi-*
1026 *ology*.
- 1027 [86] Carmi, S., Wilton, P. R., Wakeley, J., and Peer, I. (2014). A renewal theory approach
1028 to ibd sharing. *Theoretical population biology* *97*, 35–48.
- 1029 [87] Nei, M. (1978). Estimation of average heterozygosity and genetic distance from a
1030 small number of individuals. *Genetics* *89*, 583–590.
- 1031 [88] Watterson, G. (1975). On the number of segregating sites in genetical models without
1032 recombination. *Theoretical Population Biology* *7*, 256–276.
- 1033 [89] Tajima, F. (1989). Statistical method for testing the neutral mutation hypothesis
1034 by DNA polymorphism. *Genetics* *123*, 585–595.
- 1035 [90] Fu, Y.-X. (1995). Statistical properties of segregating sites. *Theoretical Population*
1036 *Biology* *48*, 172–197.
- 1037 [91] Kruglyak, L., Nickerson, D. A., et al. (2001). Variation is the spice of life. *Nature*
1038 *Genetics* *27*, 234–235.

List of Figures

- 1 tMRCA regression. We simulated a chromosome of 50 cM for 250 diploid samples, using $\mu = 2 \times 10^{-8}$ for the mutation rate and no genotyping error. We matched the allele frequency spectrum of the simulated samples to the spectrum found in real data for IBD detection with GERMLINE, and used the IBD detection parameters used in real data. The slope of this regression captures the simulated mutation rate; the intercept is proportional to genotyping error rate. 4
- 2 MaAF-threshold regression. We simulated 250 diploid samples as described in Figure 1, and a probability of 6×10^{-6} for a basepair to be involved in a non-crossover gene conversion event. We performed the MaAF-threshold regression to correct for the occurrence of gene conversion. The regression intercept is used to estimate the corrected mutation rate, while the difference between corrected and uncorrected mutation rates captures the effects of gene conversion, whose magnitude can be estimated using the observed population heterozygosity. 5
- 3 Inferred mutation rates under several values of simulated genotyping error rate, for three types of genotyping errors. The simulated true underlying mutation rate was $\mu = 2 \times 10^{-8}$. All simulations involved a single chromosome of 250 cM for 200 haploid individuals from a GoNL-like population using $Beta(\alpha = 0.5, \beta = 1)$ as a prior for allele frequency of erroneous variants. True IBD segments were extracted from the simulated ancestral recombination graph. Additional simulation results are shown in Figure S12. 6
- 4 Comparison of the estimate standard error for trios and tMRCA under different demographic models and minimum IBD segment length cutoffs. We report the estimated standard deviation from the analysis of several simulations of a single 100 Mb chromosome. For illustrative purposes, we show results of analyses using IBD length cutoffs of 1.0 and 1.5 cM. Analysis of the GoNL data was performed using a length cutoff of 1.6 cM. 7

5	Inference of gene conversion-corrected mutation rate in simulated data. We simulated a chromosome of 50 cM for 250 diploid samples, using $\mu = 2 \times 10^{-8}$ for the mutation rate and a probability of 6×10^{-6} for a basepair to be involved in a non-crossover gene conversion event. We matched the allele frequency spectrum of the simulated samples to the spectrum found in real data for IBD detection with GERMLINE. We used several values of the GERMLINE allowed mismatching sites (“-het”) to assess the impact of this parameter in the results. Negligible biases are observed for the recovered mutation rate.	8
6	tMRCA regression for segments of length ≥ 1.6 cM in the GoNL data set. The obtained slope is used to estimate mutation rate per generation per base pair, before the effects of gene conversion are accounted for.	9
7	MaAF-threshold regression for segments of length ≥ 1.6 cM in the GoNL data set. We compute mutation rates for several allowed maximum allele frequency thresholds between 0.125 and 0.5 (green dots), and regress the observed heterozygosity on the maximum allele frequency. The intercept of the resulting linear model reflects the corrected mutation rate estimate.	10
8	Association between recombination rate and mutation rate. We annotated the genome based on uniform bins of recombination rate, and estimated mutation rates for each obtained annotation. We observed a strong association between mutation and recombination rate before correcting for the occurrence of gene conversion events. After applying the correction we detect no significant association, which suggests the linear relationship observed for the uncorrected estimates is induced by gene conversion (See Figure S18).	11
9	Relationship between region-specific values of the B statistic and the average length of IBD segments > 1.6 cM spanning the regions. Equally-spaced bins of the B statistic were used. Reduced local effective population size has similar effects on the B statistic and the length of IBD haplotypes, which are longer in regions of strong background selection due to earlier average coalescent times between pairs of individuals.	12

List of Tables

- 1 Analysys of Polyphen 2 and Gerp++ annotated variants: genome-wide vs. mismatching within IBD segments. 13

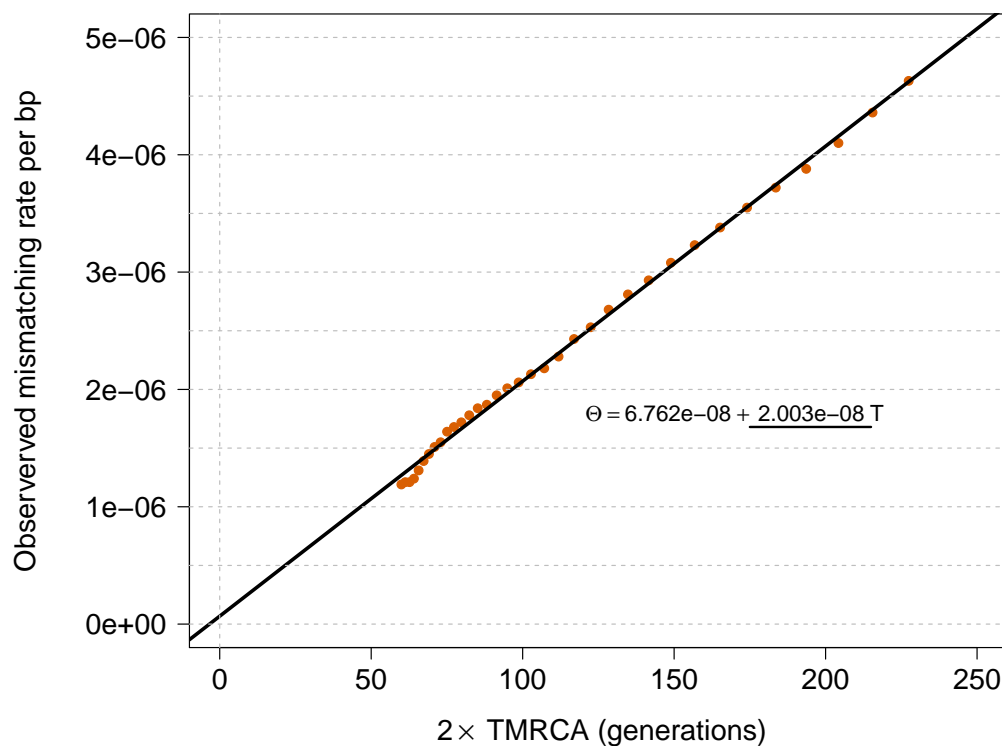


Figure 1: tMRCA regression. We simulated a chromosome of 50 cM for 250 diploid samples, using $\mu = 2 \times 10^{-8}$ for the mutation rate and no genotyping error. We matched the allele frequency spectrum of the simulated samples to the spectrum found in real data for IBD detection with GERMLINE, and used the IBD detection parameters used in real data. The slope of this regression captures the simulated mutation rate; the intercept is proportional to genotyping error rate.

—=

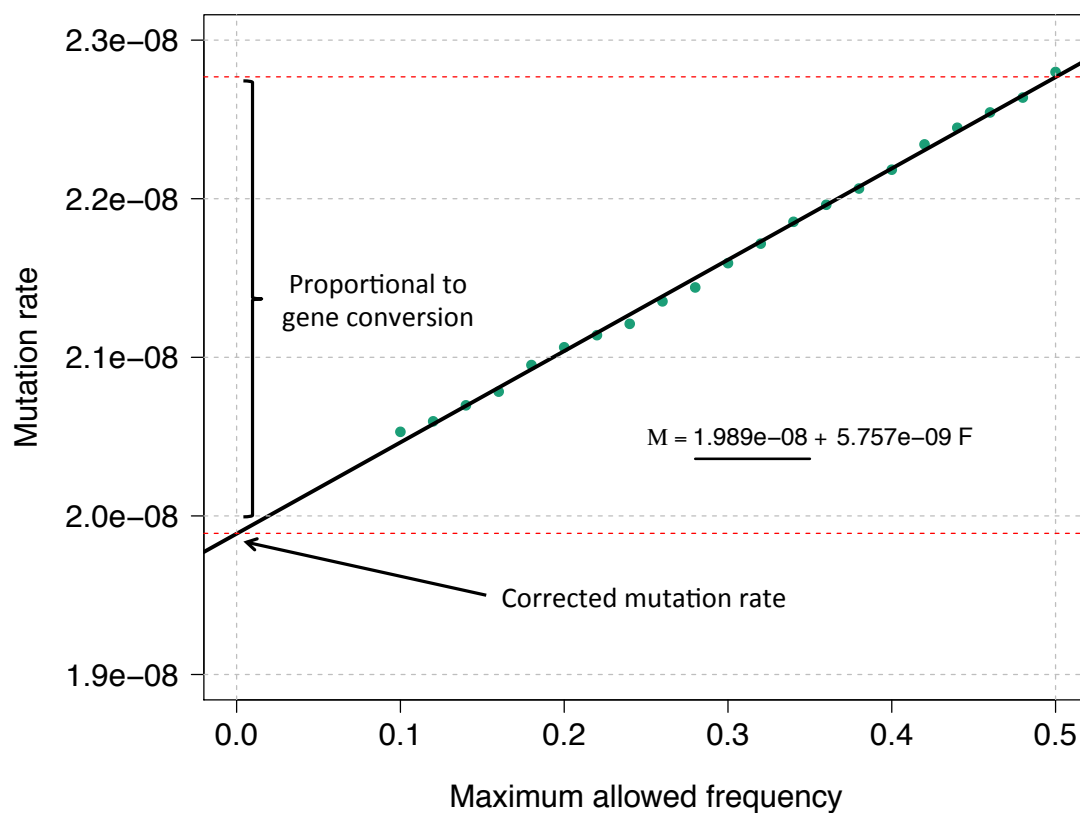


Figure 2: MaAF-threshold regression. We simulated 250 diploid samples as described in Figure 1, and a probability of 6×10^{-6} for a basepair to be involved in a non-crossover gene conversion event. We performed the MaAF-threshold regression to correct for the occurrence of gene conversion. The regression intercept is used to estimate the corrected mutation rate, while the difference between corrected and uncorrected mutation rates captures the effects of gene conversion, whose magnitude can be estimated using the observed population heterozygosity.

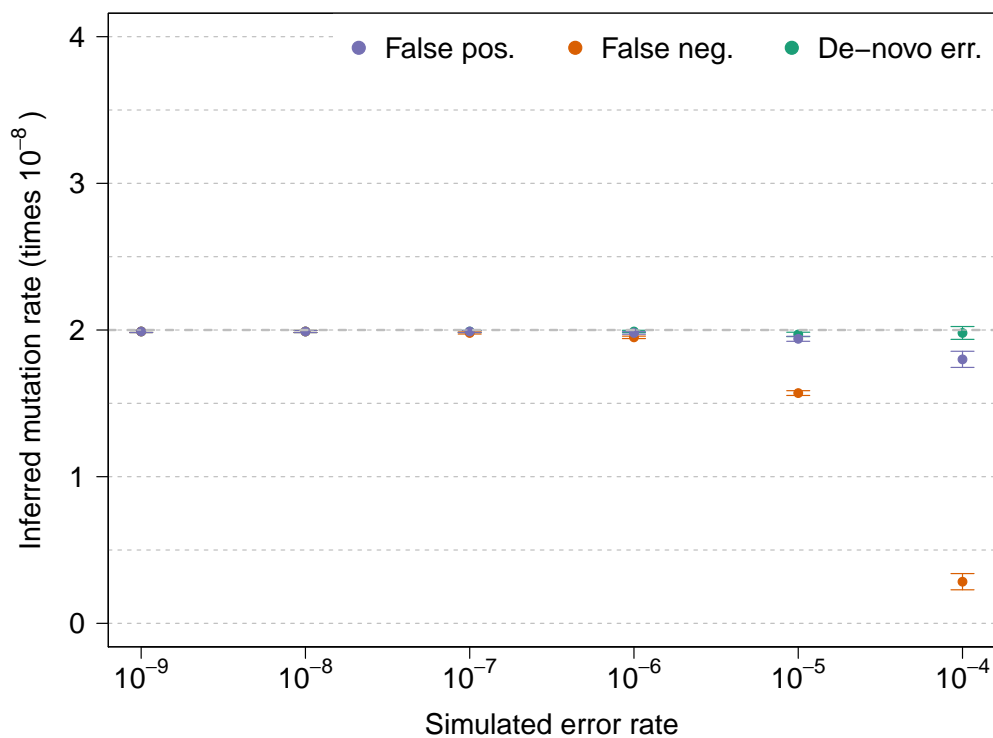


Figure 3: Inferred mutation rates under several values of simulated genotyping error rate, for three types of genotyping errors. The simulated true underlying mutation rate was $\mu = 2 \times 10^{-8}$. All simulations involved a single chromosome of 250 cM for 200 haploid individuals from a GoNL-like population using $Beta(\alpha = 0.5, \beta = 1)$ as a prior for allele frequency of erroneous variants. True IBD segments were extracted from the simulated ancestral recombination graph. Additional simulation results are shown in Figure S12.

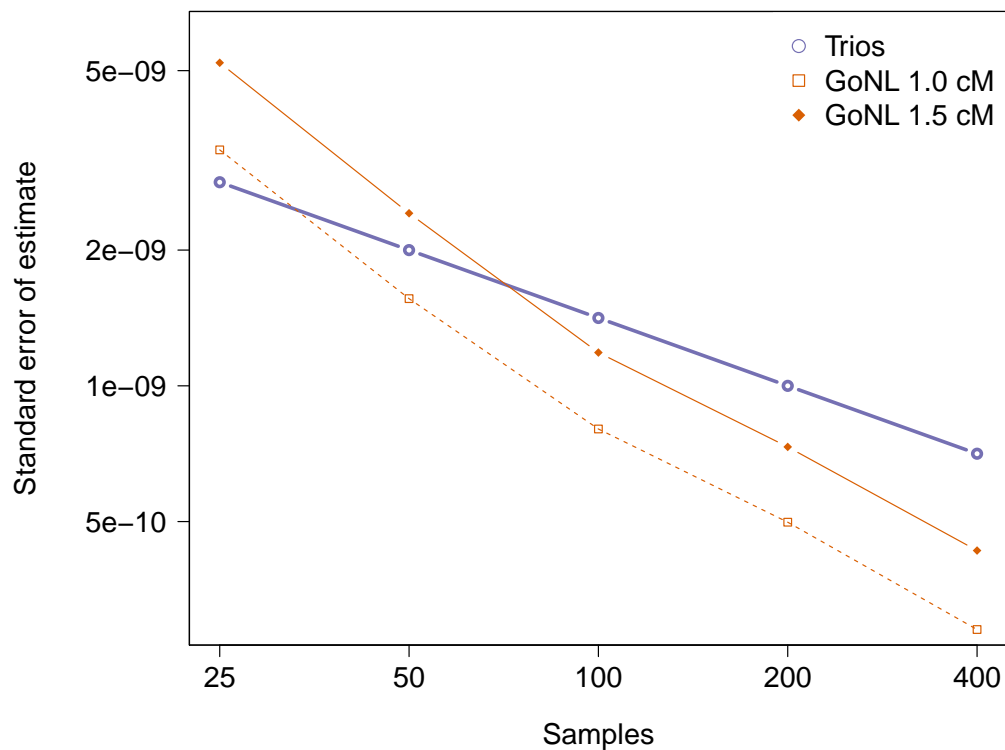


Figure 4: Comparison of the estimate standard error for trios and tMRCA under different demographic models and minimum IBD segment length cut-offs. We report the estimated standard deviation from the analysis of several simulations of a single 100 Mb chromosome. For illustrative purposes, we show results of analyses using IBD length cutoffs of 1.0 and 1.5 cM. Analysis of the GoNL data was performed using a length cutoff of 1.6 cM.

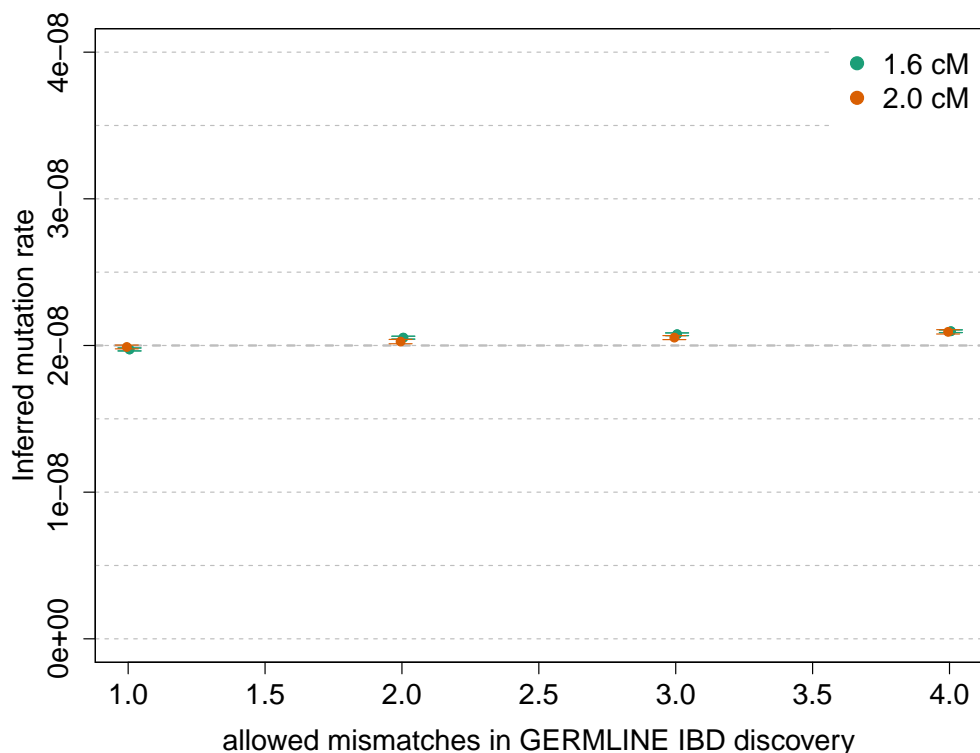


Figure 5: Inference of gene conversion-corrected mutation rate in simulated data. We simulated a chromosome of 50 cM for 250 diploid samples, using $\mu = 2 \times 10^{-8}$ for the mutation rate and a probability of 6×10^{-6} for a basepair to be involved in a non-crossover gene conversion event. We matched the allele frequency spectrum of the simulated samples to the spectrum found in real data for IBD detection with GERMLINE. We used several values of the GERMLINE allowed mismatching sites (“-het”) to assess the impact of this parameter in the results. Negligible biases are observed for the recovered mutation rate.

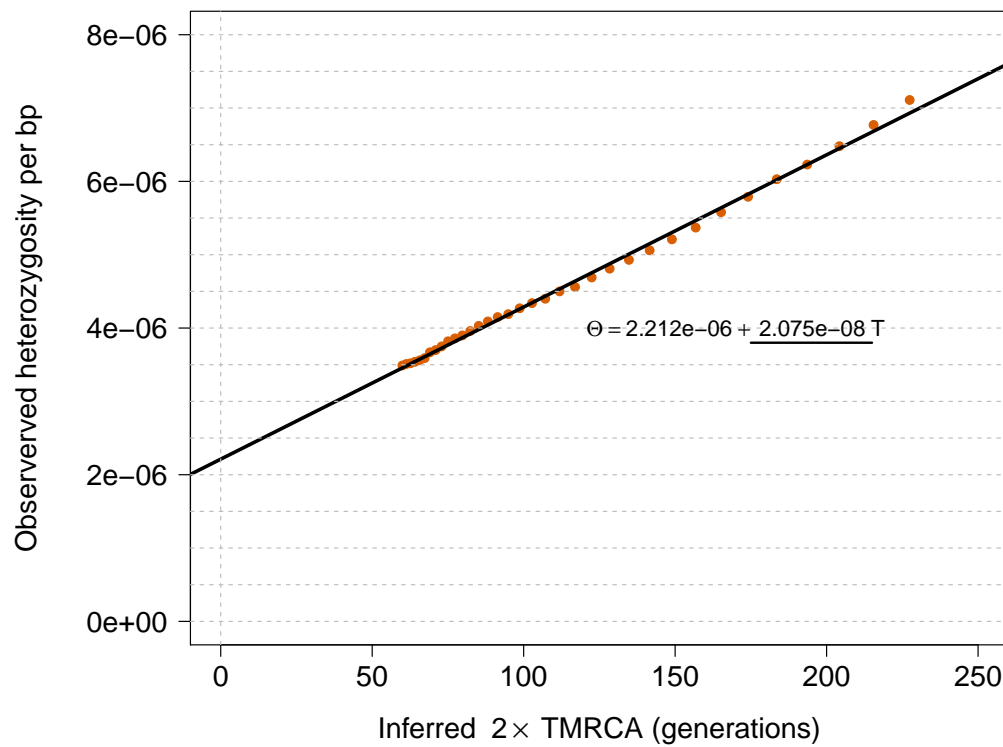


Figure 6: tMRCA regression for segments of length ≥ 1.6 cM in the GoNL data set. The obtained slope is used to estimate mutation rate per generation per base pair, before the effects of gene conversion are accounted for.

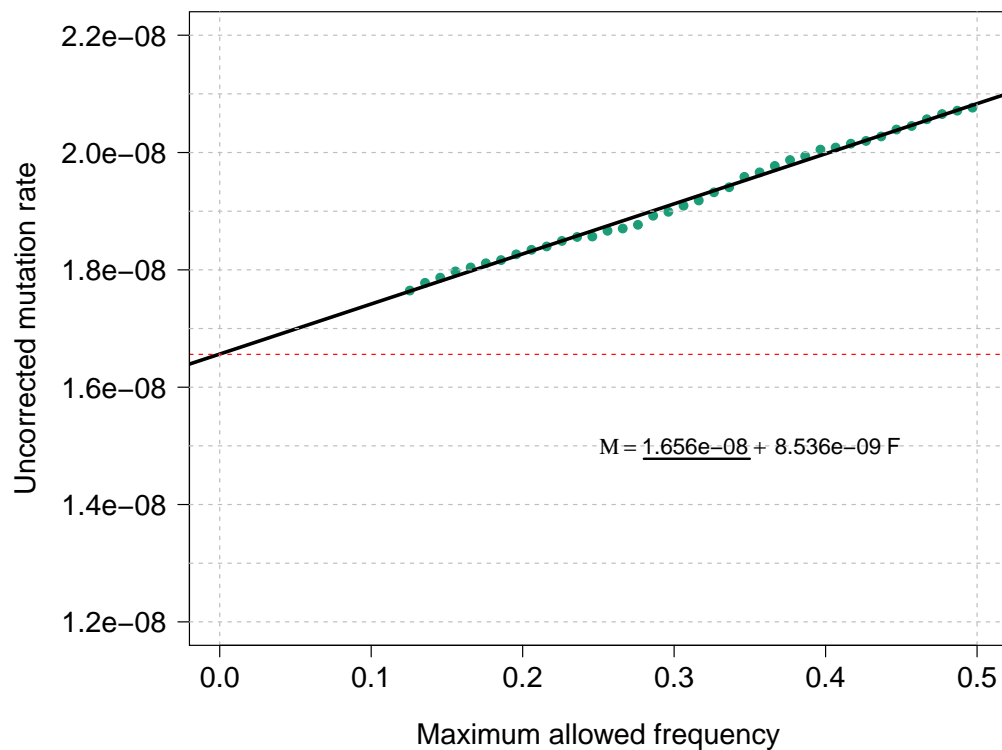


Figure 7: MaAF-threshold regression for segments of length ≥ 1.6 cM in the GoNL data set. We compute mutation rates for several allowed maximum allele frequency thresholds between 0.125 and 0.5 (green dots), and regress the observed heterozygosity on the maximum allele frequency. The intercept of the resulting linear model reflects the corrected mutation rate estimate.

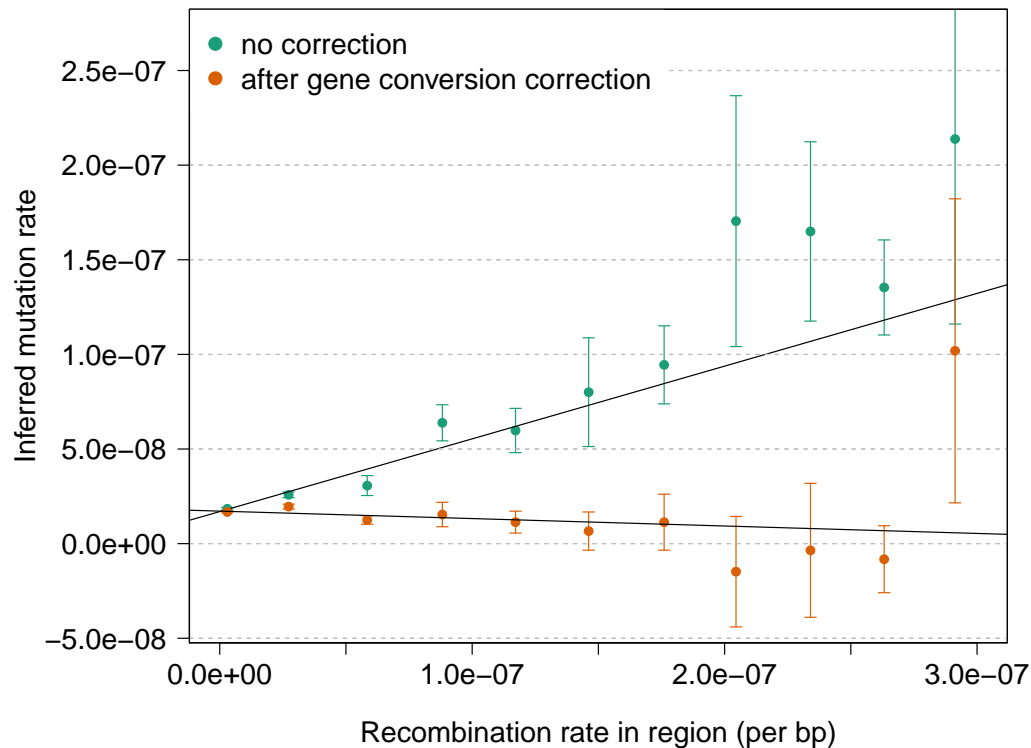


Figure 8: Association between recombination rate and mutation rate. We annotated the genome based on uniform bins of recombination rate, and estimated mutation rates for each obtained annotation. We observed a strong association between mutation and recombination rate before correcting for the occurrence of gene conversion events. After applying the correction we detect no significant association, which suggests the linear relationship observed for the uncorrected estimates is induced by gene conversion (See Figure S18).

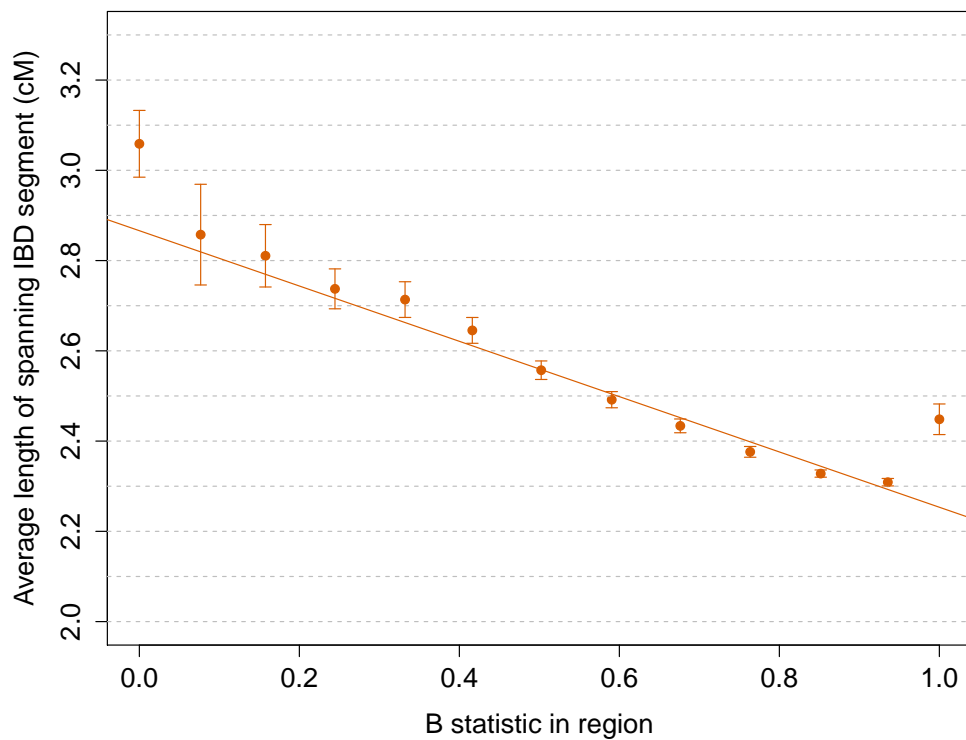


Figure 9: Relationship between region-specific values of the B statistic and the average length of IBD segments > 1.6 cM spanning the regions. Equally-spaced bins of the B statistic were used. Reduced local effective population size has similar effects on the B statistic and the length of IBD haplotypes, which are longer in regions of strong background selection due to earlier average coalescent times between pairs of individuals.

	Genome-wide	Mismatching in IBD
Annotated variants	54,960	1,843
Score mean	0.41 ± 0.0018	0.45 ± 0.0099

(a) Polyphen 2 results.

	Genome-wide	Mismatching in IBD
Annotated variants	948,782	27,900
Score mean	3.08 ± 0.00098	3.11 ± 0.0059

(b) Gerp++ (> 2) results.

Table 1: Analysis of Polyphen 2 and Gerp++ annotated variants: genome-wide vs. mismatching within IBD segments.