

# Attacks on genetic privacy via uploads to genealogical databases

Michael D. Edge & Graham Coop

Center for Population Biology and Department of Evolution and Ecology,  
University of California, Davis

December 17, 2019

## Abstract

Direct-to-consumer (DTC) genetics services are increasingly popular for genetic genealogy, with tens of millions of customers as of 2019. Several DTC genealogy services allow users to upload their own genetic datasets in order to search for genetic relatives. A user and a target person in the database are identified as genetic relatives if the user's uploaded genome shares one or more sufficiently long segments in common with that of the target person—that is, if the two genomes share one or more long regions identical by state (IBS). IBS matches reveal some information about the genotypes of the target person, particularly if the chromosomal locations of IBS matches are shared with the uploader. Here, we describe several methods by which an adversary who wants to learn the genotypes of people in the database can do so by uploading multiple datasets. Depending on the methods used for IBS matching and the information about IBS segments returned to the user, substantial information about users' genotypes can be revealed with a few hundred uploaded datasets. For example, using a method we call IBS tiling, we estimate that an adversary who uploads approximately 900 publicly available genomes could recover at least one allele at SNP sites across up to 82% of the genome of a median person of European ancestries. In databases that detect IBS segments using unphased genotypes, approximately 100 uploads of falsified datasets can reveal enough genetic information to allow accurate genome-wide imputation of every person in the database. Different DTC services use different methods for identifying and reporting IBS segments, leading to differences in vulnerability to the attacks we describe. We provide a proof-of-concept demonstration that the GEDmatch database in particular uses unphased genotypes to detect IBS and is vulnerable to genotypes being revealed by artificial datasets. We suggest simple-to-implement suggestions that will prevent the exploits we describe and discuss our results in light of recent trends in genetic privacy, including the recent use of uploads to DTC genetic genealogy services by law enforcement.

## 31 1 Introduction

32 As genotyping costs have fallen over the last decade, direct-to-consumer (DTC) genetic testing  
33 (Hogarth et al., 2008; Hogarth and Saukko, 2017; Khan and Mittelman, 2018) has become a  
34 major industry, with over 26 million people enrolled in the databases of the five largest companies  
35 (Regalado, 2019). One of the major applications of DTC genetics is genetic genealogy. Customers  
36 of companies such as 23andMe and Ancestry, once they are genotyped, can view a list of other  
37 customers who are likely to be genetic relatives. These putative relatives' full names are often  
38 given, and sometimes contact details are given as well. Such genealogical matching services are  
39 of interest to people who want to find distant genetic relatives to extend their family tree, and  
40 can be particularly useful to people who otherwise may not have information about their genetic  
41 relatives, such as adoptees or the biological children of sperm donors. Several genetic genealogy  
42 services—including GEDmatch, MyHeritage, FamilyTreeDNA, and LivingDNA (Table 1)—allow  
43 users to upload their own genetic data if they have been genotyped by another company. These  
44 entities generally offer some subset of their services at no charge to uploaders, which helps to  
45 grow their databases. Upload services have also been used by law enforcement, with the goal of  
46 identifying relatives of the source of a crime-scene sample (Erich et al., 2018; Edge and Coop,  
47 2019), prompting discussion about genetic privacy (Court, 2018; Ram et al., 2018; Kennett, 2019;  
48 Scudder et al., 2019).

49 The genetic signal used to identify likely genealogical relatives is identity by descent (IBD,  
50 Browning and Browning 2012; Thompson 2013. We use "IBD" to indicate both "identity by  
51 descent" and "identical by descent," depending on context.) Pairs of people who share an ancestor  
52 in the recent past can share segments of genetic material from that ancestor. The distribution  
53 of IBD sharing as a function of genealogical relatedness is well studied (Donnelly, 1983; Huff  
54 et al., 2011; Browning and Browning, 2012; Thompson, 2013; Buffalo et al., 2016; Conomos  
55 et al., 2016; Ramstetter et al., 2018), and DTC genetics entities can use information about  
56 the number and length of inferred IBD segments between a pair of people to estimate their  
57 likely genealogical relationship (Staples et al., 2016; Ramstetter et al., 2017). These shared  
58 segments—IBD segments—result in the sharing of a near-identical stretch of chromosome (a  
59 shared haplotype). Shared haplotypes can most easily be identified looking for long genomic  
60 regions where two people share at least one allele at nearly every locus.

61 For the rest of the main text, we focus on identical-by-state (IBS) segments, which are  
62 genomic runs of (near) identical sequence shared among individuals and can be thought of as a  
63 superset of true IBD segments. Very long IBS segments, say over 7 centiMorgans (cM), are likely  
64 to be IBD, but shorter IBS segments, say  $<4$  cM, may or may not represent true IBD due to  
65 recent sharing—they may instead represent a mosaic of shared ancestry deeper in the past. Many  
66 of the algorithms for IBD detection that scale well to large datasets rely principally on detection  
67 of long IBS segments, at least as their first step (Gusev et al., 2009; Henn et al., 2012; Huang  
68 et al., 2014). We consider the effect on our results of attempting to distinguish IBS and IBD in  
69 supplementary material.

70 Many DTC genetics companies, in addition to sharing a list of putative genealogical relatives,  
71 give customers information about their shared IBS with each putative relative, possibly including  
72 the number, lengths, and locations of shared genetic segments (Table 1). This IBS report  
73 may represent substantial information about one's putative relatives—one already has access to  
74 one's own genotype, and so knowing the locations of IBS sharing with putative relatives reveals

Service	Database Size (millions)	Individuals Shown	IBS/IBD Segments Reported
GEDmatch	1.2	3,000 closest matches shown free; Unlimited w/ \$10/month license; any two kits can be searched against each other	Yes if longer than user-set threshold. Min. threshold 0.1cM, default 7cM
FamilyTreeDNA	1*	All that share at least one 9cM block or one 7.69cM block and 20 total cM	Yes, down to 1cM, for \$19 per kit
MyHeritage	3	All that share at least one 8cM block	Yes, down to 6cM, for \$29 per kit or unlimited for \$129/year. Customers may opt out
LivingDNA	Unknown	Putative relatives out to $\approx$ 4th cousin	Only sum length of matching segments reported
DNA.LAND**	0.159	Top 50 matches shown with minimum 3cM segment	Yes

Table 1: Key parameters for several genetic genealogy services that allow user uploads as of July 26th, 2019. \*Though Regalado (2019) reports that FamilyTreeDNA has two million users, he also suggests that only about half of these are genotyped at genome-wide autosomal SNPs, which is in line with other estimates (Larkin, 2018). \*\*DNA.LAND has discontinued genealogical matching for uploaded samples as of July 26th, 2019.

75 information about those relatives' genotypes in those locations (He et al., 2014). Users of genetic  
76 genealogy services implicitly or explicitly agree to this kind of genetic information sharing, in which  
77 large amounts of genetic information are shared with close biological relatives and small amounts  
78 of information are shared with distant relatives.

79 Here we consider methods by which it may be possible to compromise the genetic privacy of  
80 users of genetic genealogy databases. In particular, we show that for services where genotype data  
81 can be directly uploaded by users, many users may be at risk of sharing a substantial proportion  
82 of their genome-wide genotypes with any party that is able to upload and combine information  
83 about several genotypes. We consider two major tools that might be used by an adversary to  
84 reveal genotypes in a genetic genealogy database. One tool available to the adversary is to  
85 upload real genotype data or segments of real genotype data. When uploading real genotypes,  
86 the information gained comes by virtue of observed sharing between the uploaded genotypes and  
87 genotypes in the database (an issue also raised by Larkin, 2017). Publicly available genotypes from  
88 the 1000Genomes Project (1000 Genomes Project Consortium, 2012), Human Genome Diversity  
89 Project (Cann et al., 2002), OpenSNP project (Greshake et al., 2014), or similar initiatives might  
90 be uploaded.

91 A second tool available to the adversary is to upload artificial genetic datasets (Ney et al.,  
92 2018). In particular, we consider the use of artificial genetic datasets that are tailored to trick  
93 algorithms that use a simple, scalable method for IBS detection, that of identifying long segments  
94 in which a pair of genomes contains no incompatible homozygous sites (Henn et al., 2012; Huang  
95 et al., 2014). Such artificial datasets can be designed to reveal the genotypes of users at single

96 sites of interest or sufficiently widely spaced sites genome-wide. We describe how a set of a  
97 few hundred artificial datasets could be designed to reveal enough genotype information to allow  
98 accurate imputation of common genotypes for every user in the database.

99 Below, we describe these procedures and illustrate them using either publicly available or  
100 artificial data. We show that under some circumstances, many users could be at risk of having  
101 their genotypes revealed, either at key positions or at many sites genome-wide. In particular, we  
102 show that GEDmatch, as of mid-December 2019, was vulnerable to an attack we term IBS baiting  
103 that obtains genotype data via artificial data uploads. Our results are largely complementary to  
104 the independent work of Ney et al. (2020), which was first posted publicly within a week of the  
105 first public posting of this manuscript on bioRxiv. In the discussion, we consider our work in light  
106 of other genetic privacy concerns (Erlich and Narayanan, 2014; Naveed et al., 2015) and the work  
107 of Ney and colleagues (2020), and we give some suggested practices that DTC genetics services  
108 can adopt to prevent privacy breaches by the techniques described here.

## 109 2 Results

110 We describe three general methods for revealing the genotypes of users in genetic genealogy  
111 databases that allow uploads. The first, **IBS tiling**, involves uploading many real genotypes in  
112 order to identify genotype information from many regions in many people. The second, **IBS**  
113 **probing**, involves uploading a dataset containing a long haplotype that includes an allele of  
114 interest, creating matches at this locus. Genotypes at other places in the genome are chosen to  
115 be unlikely to generate IBS with any user in the database, so matches with the uploaded dataset  
116 are likely to be users who carry the allele of interest. The third method, **IBS baiting**, involves  
117 uploading fake datasets with long runs of heterozygosity to induce phase-unaware methods for  
118 IBS calling to reveal genotypes.

### 119 2.1 IBS tiling

120 In IBS tiling, the genotype information shared between a target user in the database and each  
121 member of a set of comparison genomes is aggregated into potentially substantial information  
122 about the target's genotypes. For example, consider a user of European ancestries. She is likely  
123 to have some degree of IBS sharing with a large set of people from across Europe (Ralph and  
124 Coop, 2013) and beyond. If one knows the user's IBS sharing locations with one random person  
125 of European ancestries (and the random person's genotype), then one can learn a little about the  
126 user's genotype. But if one can upload many people's genotypes for comparison, then one can  
127 uncover small proportions of the target user's genotypes from many of the comparison genotypes,  
128 eventually uncovering much of the target user's genome by virtue of a "tiling" of shared IBS with  
129 known genotypes (Figure 1A). Similar ideas have been suggested with application to IBD-based  
130 genotype imputation (Carmi et al., 2014).

131 We consider the amount of IBS tiling possible within a set of publicly available genotypes for  
132 872 people of European origin genotyped at 544,139 sites. We phased the sample using Beagle 5.0  
133 (Browning and Browning, 2007) and used Refined IBD software (Browning and Browning, 2013)  
134 to identify IBS segments (see Methods). In the main text, we include IBS segments that are not  
135 particularly likely to be IBD—these are IBS segments returned by Refined IBD with relatively low

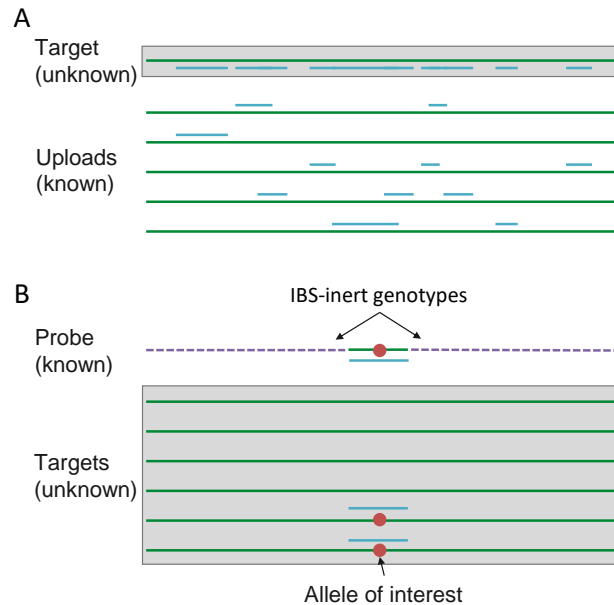


Figure 1: Schematics of the IBS tiling and IBS probing procedures. (A) In IBS tiling, multiple genotypes are uploaded (green lines) and the positions at which they are IBS with the target (represented by blue lines) are recorded. Once enough datasets have been uploaded, the target will eventually have a considerable proportion of their genome "tiled" by IBS with uploads that have known genotypes. (B) In IBS probing, the uploaded probe consists of a haplotype carrying an allele of interest (red dot) surrounded by "IBS-inert" segments (purple dashed lines)—fake genotype data designed to be unlikely to share any IBS regions with anyone in the database. In the event of an IBS match in the database, the matching database entry is likely to carry the allele of interest.

136 LOD scores for IBD, between 1 and 3. We consider the results obtained after filtering segments  
137 likely to be true IBD in Figure S1. True IBD segments reveal more than mere IBS segments about  
138 shared genotypes because untyped variants (including rare variants) within an IBD segment are  
139 likely to be shared. At the same time, mere IBS is sufficient to infer sharing for SNPs that are  
140 genotyped within the segment.

141 Once we identified IBS segments shared among the 872 people in our sample, we set out  
142 to estimate the amount of genotype information that could be identified using IBS tiling. The  
143 amount of genotype information obtainable is strongly influenced by two factors: the size of  
144 the comparison set used (i.e., the number of people used to identify IBS segments with a target  
145 sample), and the restrictiveness of the criteria by which IBS segments are identified. For example,  
146 if only long IBS segments are shown to users, then the proportion of a typical person's genotype  
147 data obtainable will be smaller than if short IBS segments are also shown. The minimum IBS  
148 length reported by several genetic genealogy services as of July 26th, 2019 is shown in Table 1.

149 Figure 2 shows the median amount of coverage obtainable by IBS tiling as a function of  
150 comparison sample size, imposing various restrictions on the minimum segment length in cM.  
151 (For similar results, see Figure 2b of Carmi et al. (2014) and Figure 2 of Panoutsopoulou et al.  
152 (2014).) Approximately 2.8 Giga base-pairs (Gbp) were covered by IBS segments anywhere in

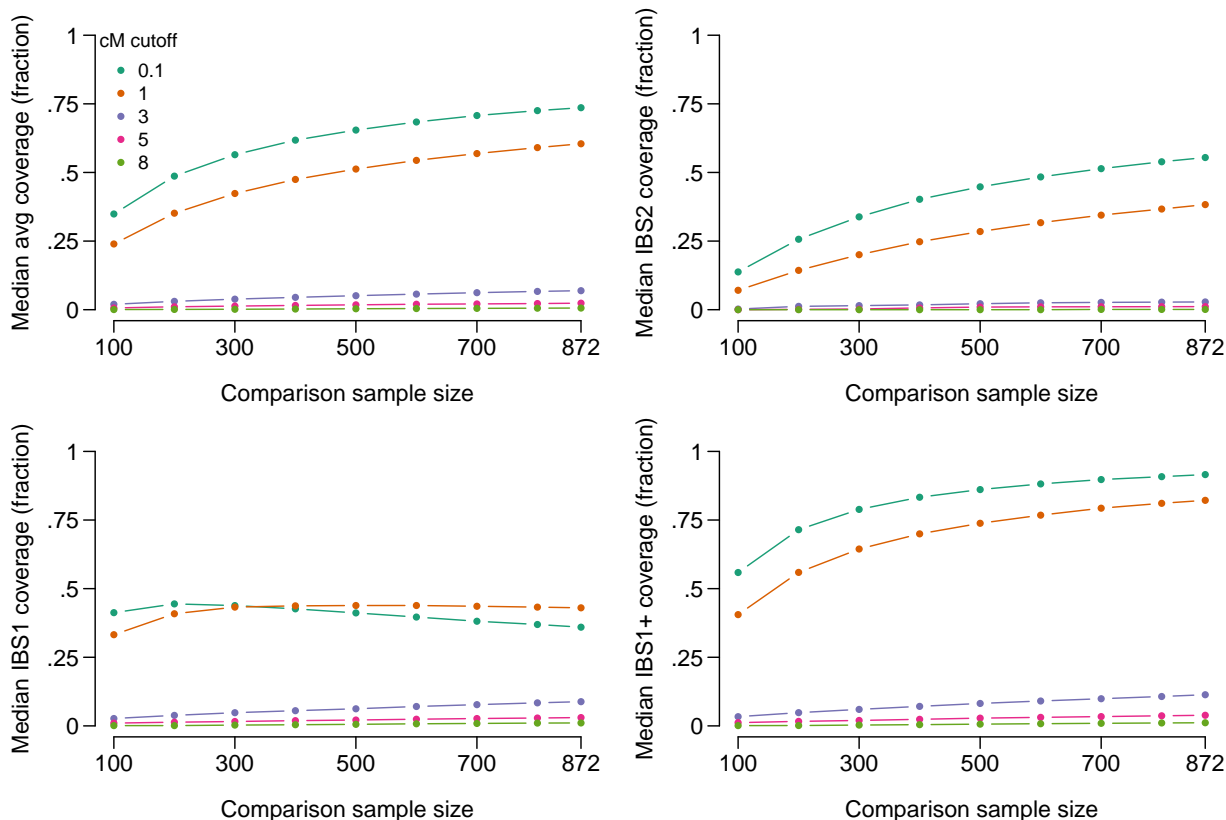


Figure 2: Lengths of genome in Giga base-pairs (Gbp) covered by IBS tiling as a function of minimum required length of IBS segments in centiMorgans (cM) and size of a randomly selected comparison sample for the median person in our dataset. The top-left panel shows the average coverage across each of the person's two haplotypes. The top-right shows IBS2 coverage, the length of genome where both haplotypes are covered by IBS tiles. The bottom-left panel shows IBS1, the length of genome where exactly one haplotype is covered by IBS tiles. (IBS1 coverage can decrease at larger comparison sample sizes because IBS2 coverage increases.) The bottom-right panel shows IBS1+ coverage, the length of genome covered by either IBS1 or IBS2.

153 the genome; we take this to be approximately the maximum possible genomic length recoverable  
 154 by IBS with our SNP set. Using the entire sample (871 genotypes, since the target is left out)  
 155 and including all called IBS segments  $>1$  cM, the median person has an average of 60% of the  
 156 maximum length of 2.8 Gbp covered by IBS segments (with the average taken across their two  
 157 chromosomes), and sites across 82% of this length have at least one of two alleles recoverable by  
 158 IBS tiling. Increasing the cM threshold required for reporting substantially reduces the amount  
 159 of IBS tiling. With a cutoff of 3 cM, approximately 6.9% of the median person's genotype  
 160 information is recoverable, including at least one of two alleles at sites in 11% of the genome.  
 161 When a more stringent cutoff of 8 cM is used, only 1% of the genome has at least one of two  
 162 alleles recoverable for the median person when using a comparison sample of 871. Our reports  
 163 for segments longer than 3 cM may be conservative because Refined IBD sometimes splits long  
 164 IBS segments into multiple shorter segments in the presence of phasing errors (Browning and  
 165 Browning, 2013; Bjelland et al., 2017).



166 For some people, the amount of information obtainable by IBS tiling is even larger. In our  
167 sample, the top 10% of people have genotypes across 76% of their total genome covered by IBS  
168 tiles, including one or more alleles at sites in at least 93% of the 2.8 Gbp covered by IBS tiles  
169 anywhere. If only segments longer than 3 cM are reported, the top 10% of people have one or  
170 both alleles covered at sites in at least 38% of the total, and if only segments longer than 8 cM  
171 are reported, the top 10% have one or both alleles covered at sites in at least 6% of the total.

172 The coverage obtained by IBS tiling and its informativeness about target genotypes depends  
173 on the specific practices used for reporting IBS information (Figures S1-S5). For example, some  
174 DTC genealogy services only report matching segments for pairs of people who share at least one  
175 long IBS segment (Table 1), but then allow users to see shorter IBS segments ( $> 1\text{cM}$ ) for those  
176 pairs of people. Unsurprisingly, we find that this strategy allows a higher level of IBS tiling than  
177 if only long segments are revealed (Figure S2), because people who share a long IBS segment  
178 may also share shorter segments that are hidden if only long segments are reported.

179 In this demonstration of IBS tiling, we used haplotype information provided by the Refined  
180 IBD software to determine which haplotypes were covered by IBS in each person. Most genetic  
181 genealogy services that provide information on the location of IBS matches with putative rela-  
182 tives do not provide haplotype information, making it difficult to distinguish IBS1 (in which one  
183 chromosome is covered by an IBS segment) and IBS2 (in which both chromosomes are covered  
184 by IBS segments). One tool available to an adversary pursuing IBS tiling is to upload genotype  
185 information that is homozygous at all sites using one of two phased haplotypes as a basis, effec-  
186 tively searching for IBS with one chromosome at a time. In the presence of phasing errors, some  
187 IBS segments may be missed, and the assumption that phase is known would render the coverage  
188 rates in Figure 2 overestimates. At the same time, the decrease in tiling performance is small for  
189 short segments, which can be seen by conducting our test of IBS tiling using Germline software  
190 with the haploid flag, which causes putative IBS segments to terminate with a single phasing  
191 error (Figure S3). It may remain difficult to distinguish some cases—such as distinguishing IBS1  
192 from IBS2 with a run of homozygosity on the database genotype—but there will be no question  
193 about which uploaded haplotype is IBS with the database genotype. Thus, at any point where a  
194 homozygous upload and a target are IBS, at least one of the target's alleles is known. Further,  
195 if the target is IBS with any other uploaded datasets at a genetic locus of interest, it will often  
196 be possible to infer the target's full genotype.

197 IBS tiling rates vary somewhat by population, with Finnish samples showing the highest tiling  
198 rates among the 1000Genomes populations included (Figure S4). There also appear to be slight  
199 biases for IBS tiles to appear in regions with low SNP density and lower heterozygosity, meaning  
200 that the proportion of alleles—and particularly the proportion of minor alleles—recovered by tiling  
201 is typically slightly lower than the proportion of the genome length in Mbp covered (Figure S5).

## 202 **2.2 IBS probing**

203 IBS probing is an application of the same idea underlying IBS tiling. By IBS probing, one could  
204 identify people with specific genotypes of interest, such as risk alleles for Alzheimer's disease  
205 (Corder et al., 1993), even if the DTC service does not report chromosomal locations of IBS  
206 matches. To identify people carrying a particular allele at a locus of interest, one could use  
207 haplotypes carrying the allele in publicly available databases. To do so, one would extract a  
208 haplotype that surrounds the allele of interest and place it into a false genetic dataset designed

209 to have no long IBS segments with any real genomes (Figure 1B). Thus, any returned putative  
210 relatives must match at the allele of interest, revealing that they carry the allele. We call this  
211 attack “IBS probing” by analogy with hybridization probes, as the genuine haplotype around the  
212 allele of interest acts as a probe. Whereas IBS tiling recovers genetic information from across the  
213 genome, IBS probing acts only on a single locus of interest. The advantage is that IBS probing  
214 is possible even in databases that do not report the chromosomal locations of IBS segments.

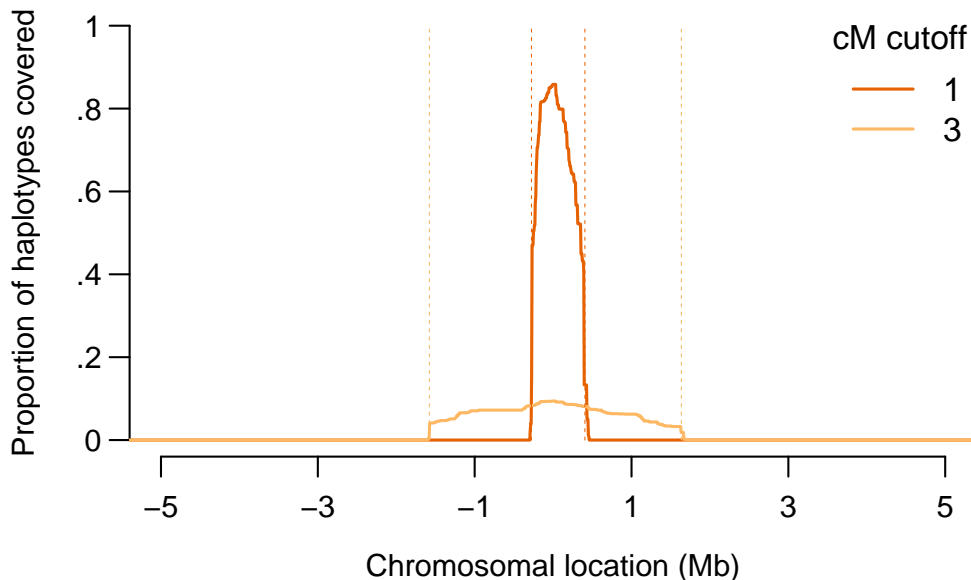


Figure 3: A demonstration of the IBS probing method around position 45411941 on chromosome 19 (GRCh37 coordinates), in the APOE locus. We show the proportion of haplotypes among the 872 Europeans in our sample covered IBS by probes constructed from the sample, as a function of the chromosomal location in a 10-Mb region around the site of interest. In red, we show the coverage using a 1cM threshold for reporting IBS, where the probes are constructed using real data in a 1.9-cM region centered on the site of interest (region boundaries shown in dashed orange). In orange, we show the coverage using a 3cM threshold for reporting IBS, where the probes are constructed using real data in a 5.9-cM region around the site of interest.

215 There are several ways of generating chromosomes unlikely to have long shared segments with  
216 any entries in the database. One simple way is to sample alleles at each locus in proportion to  
217 their frequencies. Chromosomes generated in this way are free of linkage disequilibrium (LD) and  
218 thus unlike genuine chromosomes. If the database distinguishes between IBS and IBD, then these  
219 fake data are unlikely to register as IBD with any genuine haplotypes. However, they may appear  
220 as IBS in segments where genetic diversity is low, depending on the length threshold used by the  
221 database. Near-zero rates of IBS can be obtained by generating more unusual-looking fake data,  
222 such as by sampling alleles from one minus their frequency or by generating a dataset of all minor  
223 alleles.

224 Figure 3 shows a demonstration of IBS probing performance in our set of 872 Europeans in  
225 a window around the APOE locus. For a 1cM threshold for reporting IBS, we generated probes  
226 by retaining 1.9cM of real data around a site of interest in the APOE locus from all 872 people.  
227 Outside that 1.9cM window, we generated data by drawing alleles randomly (see Methods). For a



228 3cM threshold for reporting IBS, we generated probes by retaining 5.9cM of real data around the  
229 site of interest. With 1cM matching, 1497 of 1744 haplotypes (86%) matched one of the probes  
230 at the site of interest. (Target haplotypes were not allowed to match probes constructed from  
231 the same person that carried the target haplotype.) With 3cM matching, 164 of 1744 haplotypes  
232 (9.4%) matched one of the probes at the site of interest. Very few matches occurred outside the  
233 region of interest—none with a 3cM threshold and only 0.1% of matches with a 1cM threshold.  
234 Moreover, we generated different inert genotypes for all 872 probes, and the great majority of  
235 these had no matches with any real sample. An adversary would only need to generate one inert  
236 dataset, which can be tested by uploading to the database and confirming that no matches are  
237 returned. Probes could then be constructed by stitching real haplotypes at the site of interest  
238 into the the same set of inert data. The probes would then be likely to match each other, but  
239 the adversary would know those identities and could ignore those matches.

240 The efficacy of IBS probing will depend on the minimum IBS-match length reported to users,  
241 the specific methods used for identifying IBS segments (Figures S6-S7), and whether the genotype  
242 of interest is included on the SNP chip. These factors vary in terms of whether they affect the  
243 sensitivity of IBS probing—the proportion of people carrying the allele of interest returned by a  
244 probe or set of probes—or the precision of IBS probing—the proportion of people returned by a  
245 probe who in fact carry the genotype of interest. For example, high thresholds for IBS reporting  
246 will mean that uploaded genotypes will need to have long IBS segments with targets at the locus  
247 of interest. Long IBS segments are likely to represent relatively close genealogical relatives (i.e.,  
248 long IBS segments are likely to be IBD segments), and not many targets will be close relatives  
249 of the source of any given haplotype of interest, meaning that the sensitivity of IBS probing is  
250 reduced by reporting thresholds that require long IBS segments. If the locus of interest, or a  
251 highly correlated one, is not included on the chip used to genotype either the uploaded sample or  
252 the target sample, then probing may only expected to work well if the upload and the target are  
253 truly IBD rather than merely IBS, reducing the precision of IBS probing for variants that are not  
254 genotyped. Limiting probing results to likely IBD matches will decrease the number of matches  
255 returned, particularly for short cM thresholds (Figure S6).

256 Another factor that will affect the success of IBS probing is the frequency of the allele of  
257 interest. For example, if the allele of interest is very rare, then it is likely to be only somewhat  
258 enriched on the haplotypes that tend to carry it, and reported matches may not actually carry the  
259 allele, even if they are IBD with an uploaded haplotype that carries it. IBS probing will perhaps  
260 be most sensitive and precise when the allele of interest is both common and relatively young,  
261 as is the case for founder mutations. In this case, most carriers of the allele will share the same  
262 long haplotype around the site of interest, meaning that fewer probes would need to be uploaded  
263 in order to learn the identities of the majority of the carriers in the database.

## 264 **2.3 IBS baiting**

265 IBS tiling and IBS probing take advantage of publicly available genotype data. The idea of both  
266 is that an adversary uploads genuine genetic datasets—or, in the case of IBS probing, datasets  
267 with genuine segments—to learn about entries in the database that share segments with the  
268 uploaded genomes.

269 In this section, we describe an exploit called IBS baiting. The specific strategy for IBS baiting  
270 that we describe is possible if the database identifies putative IBS segments by searching for

271 long regions where a pair of people has no incompatible homozygous sites. An incompatible  
272 homozygous site is a site at which one person in the pair is homozygous for one allele, and the  
273 other person is homozygous for the other allele. Identifying IBS segments in this way does not  
274 require phased genotypes and scales relatively easily to large datasets—we refer to methods in  
275 this class as "phase-unaware" and contrast them with phase-aware methods for IBS detection.  
276 Phase-unaware methods are robust to phasing errors, which are an issue for long IBD segments  
277 (Durand et al., 2014). Major DTC genetics companies have used phase-unaware methods in  
278 the past for IBS detection (Henn et al., 2012; Hon et al., 2013), and some state-of-the-art IBD  
279 detection and phasing pipelines feature an initial phase-unaware step (Huang et al., 2014; Loh  
280 et al., 2016).

281 The main tool used in IBS baiting is the construction of apparently IBS segments by assigning  
282 every uploaded site in the region to be heterozygous. (SNPs with missing data may also included  
283 in these regions). These runs of heterozygosity, which are unlikely to occur naturally (unlike  
284 runs of homozygosity, McQuillan et al., 2008; Pemberton et al., 2012), will be identified as  
285 IBS with every genome in the database using phase-unaware methods: because they contain no  
286 homozygous sites at all, they cannot contain homozygous sites incompatible with any person in  
287 the database.

288 Here, we consider a database in which an apparent IBS segment is halted exactly at the places  
289 at which the first incompatible homozygous site occurs on each side of the segment. We also  
290 assume that the database detects all segments without incompatible homozygous sites that pass  
291 the required length threshold. Ney et al. (2020) independently proposed a similar approach in their  
292 section VII 'Genetic Marker Extraction Using Matching Segments,' showing that GEDmatch was  
293 vulnerable to it. Similarly, we demonstrate below that IBS baiting can be implemented against  
294 GEDmatch.

### 295 **2.3.1 Single-site IBS Baiting**

296 The simplest application of IBS baiting is to use it to reveal genotypes at a single site. If IBS is  
297 identified by looking for single incompatible homozygous sites and missing data can be ignored,  
298 then users' genotypes at any single biallelic site of interest can be determined by examining their  
299 putative IBS with each of two artificial datasets (Figure 4A). In each artificial dataset, the site  
300 of interest is flanked by a run of heterozygosity. The combined length of these two runs of  
301 heterozygosity must exceed the minimum length of IBS segment reported by the database. The  
302 adversary uploads two datasets with these runs of heterozygosity in place. In one dataset, the site  
303 of interest is homozygous for the major allele, and in the other, the site of interest is homozygous  
304 for the minor allele. If the target user is homozygous at the site of interest, then one of these two  
305 uploads will not show a single, uninterrupted IBS segment—IBS will be interrupted at the site of  
306 interest (or may not be called at all). If the IBS segment with the dataset homozygous for the  
307 major allele is interrupted, then the target user is homozygous for the minor allele. Similarly, if  
308 the IBS segment with the dataset homozygous for the minor allele is interrupted, then the target  
309 user is homozygous for the major allele. If both uploads show uninterrupted IBS segments with  
310 the target, then the target user is heterozygous at the site of interest. Thus, for any genotyped  
311 biallelic site of interest, the genotypes of every user shown as a match can be revealed after  
312 uploading two artificial datasets. Depending on how possible matches are made accessible to  
313 the adversary, the genotypes of every user could be returned. Genotypes of medical interest that

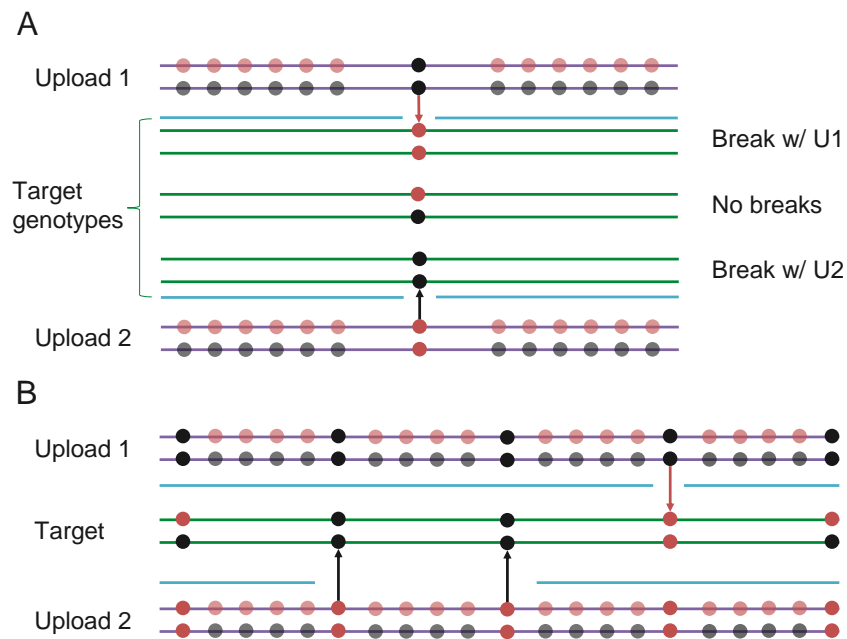


Figure 4: Schematics of the IBS baiting procedure. (A) To perform IBS baiting at a single site, two uploads are required, each with runs of heterozygous genotypes flanking the key site. At the key site, the two uploaded datasets are homozygous for different alleles. The three possible target genotypes at the key site can each be determined by examining their IBS coverage with the uploads. If there is a break in IBS with either upload, then the target is homozygous for the allele not carried by the upload that shows the break in IBS (with the broken IBS segment shown as a cyan line). If there is no break in IBS with either upload, then the target is heterozygous at the key site. (B) Target genotypes at many key sites across the genome can be learned by comparison with two uploaded datasets, as long as key sites are spaced widely enough.

314 are often included in SNP chips, such as those in the APOE locus (Corder et al., 1993), are  
 315 potentially vulnerable to single-site IBS baiting.

316 Here, we have considered a database using the simplest possible version of a phase-unaware  
 317 method for detecting IBS, that in which an apparent IBS segment is halted exactly at the places  
 318 at which the first incompatible homozygous site occurs on each side of the segment. In principle,  
 319 phase-unaware IBS-detection algorithms can be altered to allow for occasional incompatible ho-  
 320 mozygous sites before halting as an allowance for genotyping error, or the extent of the reported  
 321 region might be modified to be less than the full range between incompatible homozygous sites.  
 322 Versions of IBS baiting might be developed to work within such modifications. The key insight  
 323 is that if two artificial kits differ at exactly one site in a region and they produce two different  
 324 patterns of called IBS with a target, then the target's genotype is revealed at that site. For  
 325 example, if a database uses a phase-unaware method for IBS calling that requires two incompat-  
 326 ible homozygous sites before a putative IBS segment is halted, then an attacker might modify  
 327 our scheme by putting in a rare homozygote at a site near the key site. For most target users,  
 328 the rare homozygote in the uploaded files would be an incompatible homozygous site, implying  
 329 that a mismatch at the key site will cause a break in a putative IBS region. By using different

330 homozygote genotypes nearby, an attacker might still identify the genotypes of everyone in the  
331 database at the key site. As discussed below, such measures do not appear to be necessary to  
332 perform IBS baiting in GEDmatch. Further, in GEDmatch, uploading a third bait dataset with a  
333 missing genotype at the key site can distinguish targets with missing genotypes from heterozygous  
334 targets.

335 Single-site IBS baiting could also be used if chromosomal locations of matches are not re-  
336 ported. To do so, one would use the the scheme we describe in a large region surrounding the  
337 locus of interest and use fake IBS-inert segments to fill in the rest of the dataset.

### 338 2.3.2 Parallel IBS Baiting

339 The second method we consider applies the IBS baiting technique to many sites in parallel (Figure  
340 4B). By parallel application of IBS baiting, users' genotypes at hundreds or thousands of sites  
341 across the genome can be identified by comparison with each pair of artificial genotypes. By re-  
342 peated parallel IBS baiting, eventually enough genotypes can be learned that genotype imputation  
343 becomes accurate, and genome-wide genotypes could in principle be imputed for every user in the  
344 database. If IBS segments as short as 1cM are reported to the user, then accurate imputation  
345 (97-98% accuracy) becomes possible after comparison with only  $\approx 100$  uploaded datasets. The  
346 procedure starts by designing a single pair of uploaded files as follows:

- 347 1. Identify a set of key sites to be revealed by the IBS baiting procedure. For every key site,  
348 the sum of the distances in cM to the nearest neighboring key site on each side (or the  
349 end of the chromosome, if there is no flanking key site on one side) must be at least the  
350 minimum IBS length reported by the database.
- 351 2. Produce two artificial genetic datasets. In each, every non-key site is heterozygous. In one,  
352 each key site is homozygous for the major allele, in the other, each key site is homozygous  
353 for the minor allele.
- 354 3. Upload each artificial dataset and compare them to a target user. Key sites that are covered  
355 by putative IBS segments between the target and both artificial datasets are heterozygous  
356 in the target. The target is homozygous for the major allele at key sites that are covered by  
357 putative IBS segments between the target and the major-allele-homozygous dataset only.  
358 Similarly, the target is homozygous for the minor allele at key sites that are covered by  
359 putative IBS segments between the target and the minor-allele-homozygous dataset only.

360 Carrying out this procedure reveals the target's genotype at every key site. If IBS segments of  
361 length at least  $t$  cM are reported, and a chromosome is  $c$  cM long, then up to  $2c/t - 1$  key sites  
362 can be revealed with each pair of uploaded files. (To see this, consider the case where  $c = tk$ ,  
363 with  $k$  a positive integer, and place key sites at  $t/2, t, 3t/2, \dots, c - t/2$ . This calculation ignores  
364 the possibility of missing data at key sites in the target.) This means that with a minimum  
365 reported IBS threshold of 1cM, 100 uploaded datasets could reveal approximately 100 genotypes  
366 per cM, which is enough to impute genome-wide genotypes at 97 – 98% accuracy (Shi et al.,  
367 2018). In principle, the key sites could also be chosen to ensure good LD coverage and higher  
368 imputation accuracy. Of course, higher accuracy imputation can be obtained by recovering exact  
369 genotypes for more sites, and with several thousand uploads, the genotypes at every genotyped  
370 site could be revealed by IBS baiting without the need to impute.

### 371 2.3.3 IBS baiting in GEDmatch

372 We hypothesized that IBS baiting would work in the GEDmatch DTC database. GEDmatch  
373 provides no public documentation of the IBS algorithm they use, but IBS segments identified by  
374 GEDMatch seem to terminate only on incompatible homozygous sites, as would be expected if  
375 they use phase-unaware IBS detection. Specifically, the GEDmatch 1-to-1 match tool identifies  
376 the locations of IBS segments between pairs of genetic datasets ("kits" in GEDMatch terminology)  
377 and allows the user to specify the minimum genetic length and minimum number of matching  
378 SNPs to include in a segment. The 1-to-1 tool also returns a 'full resolution' picture of the  
379 chromosome that appears to be a SNP-by-SNP picture of the match between the kits along  
380 each chromosome. (These pictures are themselves a major security risk. We alerted GEDmatch  
381 to the risk in a July 24th email (posted here: [https://github.com/mdedge/IBS\\_privacy/  
382 blob/master/IBS\\_baiting\\_demo/GEDmatch\\_emails.pdf](https://github.com/mdedge/IBS_privacy/blob/master/IBS_baiting_demo/GEDmatch_emails.pdf)) but did not analyze them further.  
383 Ney et al. (2020) showed in detail that the images provided by GEDmatch allow an adversary to  
384 learn the full genotype of a target person.)

385 To demonstrate IBS baiting in GEDmatch, we uploaded a small number of artificial genotypes  
386 to their database beginning in late November 2019. These kits were designed in accordance with  
387 the algorithm discussed above, but with some slight alterations to bypass counter-measures that  
388 GEDMatch has put in place since we (and, independently, Ney and colleagues) informed them  
389 of the risk of IBS baiting in summer 2019. Before uploading any data to GEDmatch, we first  
390 confirmed our planned procedure with the UC Davis IRB and with GEDmatch representatives.  
391 We uploaded our kits into the GEDmatch 'research' and not 'public' category to prevent matches  
392 to the public database, and only used the 1-to-1 IBS match tool among our own uploaded test  
393 kits. In this way, we avoided interacting with any genotype data of real GEDmatch users and did  
394 not violate GEDmatch's terms and conditions.

395 We targeted four random SNPs along chromosome 22 for IBS baiting. We uploaded two bait  
396 genotype kits (B1 & B2) that had opposite-homozygote genotypes at each of these key SNPs.  
397 Each key SNP was in turn surrounded by a  $\sim 1\text{cM}$  stretch of SNPs containing genotypes that  
398 were either heterozygous or coded missing. The rest of the genome was specified to be IBS-inert.  
399 We then uploaded three target genotype datasets whose genotypes we wanted to determine at  
400 the key sites. Two of these target kits (T1 and T3) had opposite-homozygous genotypes at each  
401 of the key SNPs, while the third (T2) was heterozygous at each key SNP. (See section 4.4.1  
402 for more details on the kit design.) We then used the GEDmatch 1-to-1 match tool, choosing  
403 the parameters so a single opposite-homozygous genotype between a bait and target kit would  
404 interrupt a putative IBS segment.

405 In each case, our two bait kits had the correct IBS patterns with the target kits, allowing  
406 correct determination of the target genotypes by IBS baiting. On the left of Figure 5, we show a  
407 zoomed-in view of the three targets' matches around one of the key SNP sites. The homozygous  
408 targets have IBS matches with only one of the bait kits, whereas the heterozygous target has  
409 IBS matches with both bait kits. This pattern is seen across all four target regions (right side of  
410 Figure 5, see Section 4.4.2 for more detailed results). The target and bait kits displayed in Figure  
411 5 were uploaded and analyzed on December 15, 2019, showing that GEDmatch has remained  
412 vulnerable to IBS-baiting attacks even after its acquisition by Verogen, which was announced on  
413 December 9, 2019.

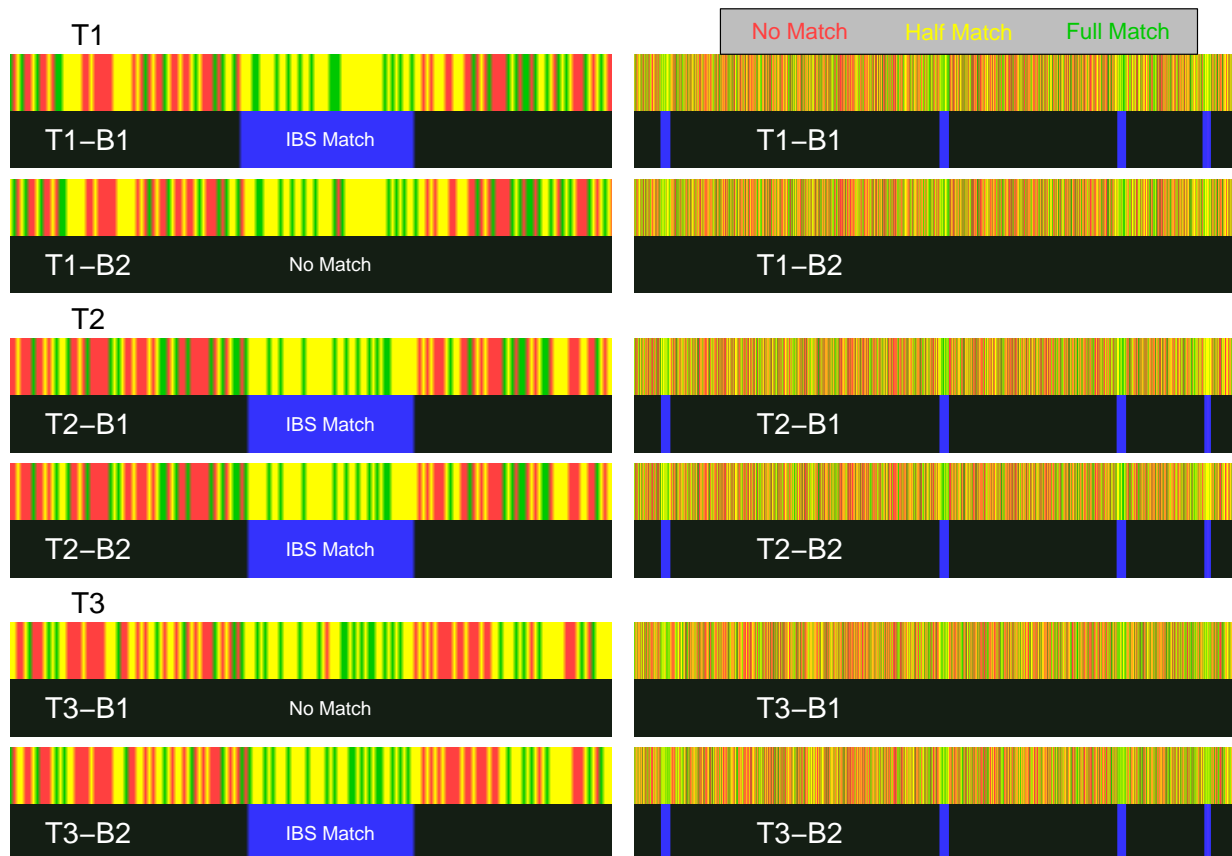


Figure 5: Visualization of IBS baiting using GEDmatch's 1-to-1 chromosome browser. **Left:** Zoomed-in view of the region containing key SNP 1, showing the three target kits (T1-T3) matched to the two bait kits (B1 & B2). **Right:** Zoomed-out views of regions containing all four key SNPs on chromosome 22. For each pair of bait & target kits, the top rectangle (red, yellow, or green) shows the GEDmatch SNP-level pairwise genotype-match image (colored to show no match, half match, or full match) returned by the 1-to-1 GEDmatch tool. The bottom rectangle (black & blue) shows the GEDmatch IBD-track image, black for no putative IBD match, blue regions showing putative IBD segments. The white text on the IBS track is not provided by GEDmatch and was added as a guide to the eye. Opposite-homozygote calls at the key SNP are seen in the left panel as a red line in an otherwise matching region (yellow & green). The spatial positions of SNPs in the match panel appears to have been jittered; e.g. the location of the red line varies slightly in the different plots that should have the same coordinate system (perhaps as a countermeasure against a Ney et al. (2020)-style attack).

### 414 3 Discussion

415 We have suggested several methods by which an adversary might learn the genotypes of people  
416 included in a genetic genealogy database that allows uploads. Our methods take advantage of



417 both the population-genetic distributions of IBS segments and of methods used for calling IBS.  
418 In particular, IBS tiling works simply because there are background levels of IBS (and IBD) even  
419 among distantly related members of a population (e.g. Ralph and Coop, 2013). In our dataset,  
420 the median person had the majority of their genetic information susceptible to IBS tiling on the  
421 basis of other members of the dataset, depending on the procedures used for reporting IBS.  
422 IBS tiling performance will also depend on the ancestries of the target and comparison samples  
423 because IBD rates differ within and among populations (Palamara et al., 2012; Carmi et al.,  
424 2013; Ralph and Coop, 2013), as well as on the prevalence of close biological relatives in the  
425 dataset. IBS tiling performance improves as the size of the comparison sample increases. Thus,  
426 if enough genomes are compared with a target for IBS, eventually a substantial amount of the  
427 target genome is covered by IBS with one or more of the comparison genomes.

428 IBS probing combines the principles behind IBS tiling with the idea of "IBS-inert" artificial  
429 segments. If the majority of the genome—everywhere except a locus of interest—can be replaced  
430 with artificial segments that will not have IBS with any genome in the database, then the adversary  
431 knows that any matches identified are in a locus of interest. As such, IBS probing could be used  
432 to reveal sensitive genetic information about database participants even if chromosomal locations  
433 of matches are not reported to users.

434 Finally, IBS baiting exploits phase-unaware IBS calling algorithms that use incompatible ho-  
435 mozygous sites to delimit putative IBS regions. Although such methods can be useful in genetic  
436 genealogy because they scale well to large data, they are vulnerable to fake datasets that include  
437 runs of heterozygous sites, which will be identified as IBS with everyone in the database. By  
438 inserting homozygous genotypes at key sites and heterozygotes everywhere else, we estimate that  
439 approximately 100 well-designed uploads could reveal enough genotypes to impute genome-wide  
440 information for any user in a database, provided that the threshold for reporting a matching  
441 segment is  $\approx 1$  cM. Similarly, two uploads could reveal any genotype at a single site of interest,  
442 such as rs429358, which reveals whether the user carries an APOE- $\epsilon 4$  variant and is associated  
443 with risk of late-onset Alzheimer's disease.

444 There are millions of people enrolled in genetic genealogy databases that allow uploads (Table  
445 1). Genetic genealogy has many applications, and uploads are popular with users who want to  
446 find relatives who may be scattered across different databases. Though allowing uploads brings  
447 several benefits for both customers and DTC services, it also entails additional privacy risks.  
448 Users of DTC genetic genealogy services that allow uploads could be at risk of having their  
449 genetic information extracted by the procedures we describe here, depending on the methods  
450 that these services use to identify and report IBS. Concerns arising from the methods we report  
451 are in addition to standard digital security concerns. The attacks we describe require little special  
452 expertise in computing; the adversary only needs to be able to procure or create the appropriate  
453 data files and to process and aggregate the information returned from the database.

454 We have not set out to determine precisely how vulnerable users of each specific DTC service  
455 are. We do not know the full details of methods used by each service for matching, nor have  
456 we attempted to deanonymize any real users' genotypes. We contacted representatives of each  
457 of the organizations listed in Table 1 90 days (July 24th, 2019) before posting this manuscript  
458 publicly in order to give them time to repair any security vulnerabilities related to the methods  
459 we describe. We have posted our emails to GEDmatch representatives here: [https://github.com/mdedge/IBS\\_privacy/blob/master/IBS\\_baiting\\_demo/GEDmatch\\_emails.pdf](https://github.com/mdedge/IBS_privacy/blob/master/IBS_baiting_demo/GEDmatch_emails.pdf).

461 On the basis of our results, we do have serious concerns about the privacy of GEDmatch

462 users. As of this writing, GEDmatch uses length thresholds for displaying matching segments  
463 that are too short, allowing for effective IBS tiling attacks, and GEDmatch also appears to use  
464 phase-unaware IBD detection methods, allowing for IBS baiting attacks. Additionally, as detailed  
465 by Ney et al. (2020), whose work was independent of ours, GEDmatch provides users with high-  
466 resolution images comparing the chromosomes of any two users at SNP-level resolution, allowing  
467 for reconstruction of a target's genotype using these images. GEDmatch was recently purchased  
468 by Verogen, a forensic genetics company, but as of December 15, 2019, GEDmatch has not as  
469 yet prevented the attacks we describe. Since our and Ney et al.'s initial communications with  
470 GEDmatch in July, GEDmatch has placed a recaptcha on its upload and 1-to-1 tool forms, which  
471 may deter bulk bot attacks to harvest large numbers of kit genotypes. However, as we outline  
472 below, there are simple steps that could be taken to make IBS attacks much less of a risk.

473 In our estimation, the other active services listed in Table 1 (MyHeritage, FamilyTreeDNA, and  
474 LivingDNA) are likely substantially less vulnerable than GEDmatch to the attacks we describe here.  
475 LivingDNA does not provide a chromosome browser, precluding IBS tiling attacks. MyHeritage  
476 and FamilyTreeDNA use thresholds for revealing matching segment locations that make IBS tiling  
477 much less efficient. (However, FamilyTreeDNA's practice of showing matches as short as 1cM  
478 given that two people share at least one long match is still somewhat permissive, see Figure S2.)  
479 Representatives of MyHeritage, FamilyTreeDNA, and LivingDNA have confirmed to us that their  
480 IBD-calling algorithms rely on phased data, which should preclude IBS baiting. (We have not  
481 tested this ourselves.) DTC genetic genealogy is a growing field, and any new entities that begin  
482 offering upload services may also face threats of the kind we describe.

483 Genetic genealogy databases that allow uploads have been in the public eye recently because  
484 of their role in long-range familial search strategies recently adopted by law enforcement. In  
485 long-range familial search, investigators seek to identify the source of a crime-scene sample by  
486 identifying relatives of the sample in a genetic genealogy database that allows uploads. Searching  
487 in SNP-based genealogy databases allows the detection of much more distant relationships than  
488 does familial searching in traditional forensic microsatellite datasets (Rohlf et al., 2012), vastly  
489 increasing the number of people detectable by familial search (Erlich et al., 2018; Edge and  
490 Coop, 2019). At this writing, both GEDmatch and FamilyTreeDNA have been searched in this  
491 way. Long-range familial search raises a range of privacy concerns (Court, 2018; Ram et al., 2018;  
492 Kennett, 2019; Scudder et al., 2019). One response from advocates of long-range search has been  
493 to note that "raw genetic data are not disclosed to law enforcement...Search results display only  
494 the length and chromosomal location of shared DNA blocks" (Greytak et al., 2018). However,  
495 the methods we describe here illustrate that there are several ways to reveal users' raw genetic  
496 data on the basis of the locations of shared DNA blocks. Because companies that work with law  
497 enforcement on long-range familial searching—including Parabon Nanolabs and Bode Technology  
498 (Kennett, 2019)—now routinely upload tens of datasets to genetic genealogy databases, they may  
499 be accidentally accumulating information that would allow them to reconstruct many people's  
500 genotypes.

501 Data breaches via IBS tiling, IBS probing, and IBS baiting are preventable. We have identified  
502 a set of strategies that genetic genealogy services could adopt to protect their genotype data from  
503 IBS-based attacks. We give a detailed list of these strategies in Appendix A (also summarized in  
504 Table 2). Broadly, the suggestions consist of restrictions on they types of datasets that can be  
505 uploaded, restrictions on the kinds of information shared with users, and restrictions on classes  
506 of methods used for identifying putative IBD segments. For example, to prevent IBS tiling,

Strategy	Prevents IBS tiling	Prevents IBS probing	Prevents IBS baiting
Require cryptographic signature from genotyping service	Yes	Yes	Yes
Restrict reporting of IBS to long segments (e.g. >8 cM)	Partially	Partially	Weakly
Report number and lengths of IBS segments but not locations	Yes	No	Partially
Block homozygous uploads	Partially	No	No
Report small number of matching individuals per kit	Partially	Partially	Partially
Disallow matching between arbitrary kits	Partially	Partially	Partially
Block uploads of publicly available genomes	Partially	No	No
Block uploads with evidence of IBS-inert segments	No	Yes	No
Block uploads with long runs of heterozygosity	No	No	Partially
Use phase-aware methods for IBS detection	No	No	Yes

Table 2: Potential countermeasures against the methods of attack outlined here, with their likely effectiveness against IBS tiling, IBS probing, and IBS baiting.

507 the simplest measures are either to forgo the use of a chromosome browser feature or only to  
 508 show users the positions of long IBS segments, such as segments of at least 8cM. To prevent  
 509 IBS baiting, the most robust countermeasure is to phase data before identifying IBS segments,  
 510 allowing only relatively few phase switches in any putative segment. Phasing the data and only  
 511 reporting long segments both decrease the uncertainty of IBD calls and so may improve user  
 512 experience as well. Finally, we also support the strategy of requiring encrypted signatures on  
 513 uploaded files, proposed by Erlich et al. (2018), which would allow DTC databases to block any  
 514 files that do not originate from trusted sources. Some of our suggestions limit the potential uses  
 515 of genetic genealogy data, and users will vary in the degree to which they value these potential  
 516 uses and in the degree to which they want to protect their genetic information.

517 All of these suggestions assume that genealogy services will maintain raw genetic data for  
 518 people in their database. Another possibility would be for individual people instead to upload  
 519 an encrypted version of their genetic data, with relative matching performed on the encrypted  
 520 datasets, as has been suggested previously (He et al., 2014).

521 Our IBS tiling and IBS probing results focus on users of European ancestries, in part because  
 522 most users of DTC genetic genealogy services appear to have substantial European ancestries.  
 523 (DTC genetics companies generally do not release this kind of information on their users, but  
 524 their research papers suggest that they have access to especially large samples with European  
 525 ancestries—for example, a 23andMe paper on demography in the United States included almost  
 526 150,000 self-described European Americans and less than 10,000 each of self-described African  
 527 Americans and Latino Americans (Bryc et al., 2015). For a qualitatively similar sample compo-  
 528 sition in a study from Ancestry, see Han et al. (2017).) One question is how these results would

529 generalize to other populations. Because IBD sharing is generally greater within populations than  
530 between populations (e.g., Ralph and Coop, 2013), potential users are more vulnerable if there are  
531 more publicly available genomes from people with similar ancestries. If IBD-detection algorithms  
532 are not well calibrated to differences in heterozygosity across populations, then spurious IBD calls  
533 will be more common in populations with lower heterozygosity, leading to greater risk of IBD  
534 tiling. Finally, we show in Figure S4 that in our sample, Finnish samples are more vulnerable to  
535 IBS tiling than other populations, which is likely due to Finns tracing substantial ancestry to a  
536 founder population that experienced a bottleneck  $\sim 100$  generations ago (Kere, 2001). Members  
537 of other groups with similar demographic histories are likely to be at elevated risk of IBS tiling  
538 and IBS probing as well.

539 We have focused on genetic genealogy databases that allow uploads because at this writing,  
540 it is straightforward to download publicly available genetic datasets and to produce fake genetic  
541 datasets for upload. In principle, however, another way to perform attacks like the ones we de-  
542 scribe would be to use biological samples. For example, a group of people willing to share their  
543 genetic data with each other could collaborate to perform IBS tiling by sending actual biological  
544 samples for genotyping. Even IBS probing and IBS baiting could be performed with biological  
545 samples by adversaries who could synthesize the samples. Though synthesizing such samples is  
546 technically challenging now, it may become easier in the future. Such methods could present  
547 opportunities to attack databases that do not allow uploads, such as the large databases main-  
548 tained by Ancestry ( $>14$  million) and 23andMe ( $>9$  million) (Regalado, 2019). They would also  
549 thwart the countermeasure of requiring uploaded datasets to include a cryptographic signature  
550 indicating their source.

551 The IBS-based privacy attacks we describe here add to a growing set of threats to genetic  
552 privacy (Homer et al., 2008; Nyholt et al., 2009; Im et al., 2012; Gymrek et al., 2013; Humbert  
553 et al., 2015; Shringarpure and Bustamante, 2015; Edge et al., 2017; Ayday and Humbert, 2017;  
554 Kim et al., 2018; Erlich et al., 2018). A person's genotype includes sensitive health information  
555 that might be used for discrimination, and people whose genetic information is compromised  
556 may be vulnerable to scams involving falsified relatives (Ney et al., 2020). Though there are  
557 many emerging threats to privacy, some of the more unsettling of which have nothing to do  
558 with genetics, genetic data do have special features that might require special considerations. In  
559 particular, genetic privacy concerns not only the person whose genotypes are directly revealed but  
560 also their relatives whose genotypes may be revealed indirectly (Humbert et al., 2013), a point  
561 highlighted by the use of genetic genealogy for long-range forensic searches (Erlich et al., 2018;  
562 Edge and Coop, 2019).

563 Though many forms of genetic discrimination are prohibited legally, rules vary between coun-  
564 tries and states. For example, in the United States, the Genetic Information Nondiscrimination  
565 Act (GINA) protects against genetic discrimination in the provision of health insurance but does  
566 not explicitly disallow genetic discrimination in the provision of life insurance, disability insurance,  
567 or long-term care insurance (Bélisle-Pipon et al., 2019). In addition to measures for protecting  
568 genetic privacy in the short term, there is a need for more complete frameworks governing the  
569 circumstances under which genetic data can be used (Clayton et al., 2019).

## 4 Methods

### 4.1 Data assembly

We performed IBS tiling with publicly available genotypes from 872 people of European ancestries. Of these 872 genotypes, 503 came from the EUR subset of the 20130502 release of phase 3 of the 1000 Genomes project (1000 Genomes Project Consortium, 2012), downloaded from <ftp://ftp.1000genomes.ebi.ac.uk/vol11/ftp/release/20130502/>. This release set has been pruned to remove close biological relatives. The EUR subset includes the following population codes and numbers of people: CEU (Utah residents with Northern and Western European Ancestry, 99 people), FIN (Finnish in Finland, 99 people), GBR (British in England and Scotland, 91 people), IBS (Iberian Population in Spain, 107 people), TSI (Toscani in Italia, 107 people).

The remaining 369 were selected from samples typed on the Human Origins SNP array (Patterson et al., 2012), including 142 genotypes from the Human Genome Diversity Project (Cann et al., 2002). Specifically, we downloaded the Human Origins data from <https://reich.hms.harvard.edu/downloadable-genotypes-present-day-and-ancient-dna-data-compiled-published-papers>, using the 1240K+HO dataset, version 37.2. The 372 selected people were all contemporary samples chosen according to population labels. We also excluded people from the Human Origins dataset if they appeared in the 1000 Genomes dataset. The populations used for selecting data, along with the number of participants included after excluding 1000 Genomes samples, were as follows: "Adygei" (16), "Albanian" (6), "Basque" (29), "Belarusian" (10), "Bulgarian" (10), "Croatian" (10), "Czech" (10), "English" (0), "Estonian" (10), "Finnish" (0), "French" (61), "Greek" (20), "Hungarian" (20), "Icelandic" (12), "Italian\_North" (20), "Italian\_South" (4), "Lithuanian" (10), "Maltese" (8), "Mordovian" (10), "Norwegian" (11), "Orcadian" (13), "Romanian" (10), "Russian" (22), "Sardinian" (27), "Scottish" (0), "Sicilian" (11), "Spanish" (0), "Spanish\_North" (0), and "Ukrainian" (9). The populations with 0 people included are those for which all the samples in the Human Origins dataset are included in the 1000 Genomes phase 3 panel. Samples with group labels marked "ignore" were excluded, including samples marked as close relatives.

We down-sampled the sequence data from the 1000 Genomes project to include only sites typed by the Human Origins chip. Of the 597,573 SNPs included in the Human Origins dataset, 558,257 sites appeared at the same position in the 1000 Genomes dataset, 557,999 of which appear as biallelic SNPs. For 546,530 of these, both the SNP identifier and position match in 1000 Genomes, and for 544,139 of them, the alleles agreed as well. We merged the dataset at the set of 544,139 SNPs at which SNP identifiers, positions, and alleles matched between the Human Origins and 1000 Genomes datasets.

We used vcftools (Danecek et al., 2011), bcftools (Li, 2011), PLINK (Purcell et al., 2007), and EIGENSOFT Price et al. (2006) to create the merged file. The script used to create it is available at [github.com/mdedge/IBS\\_privacy/](https://github.com/mdedge/IBS_privacy/), and the merged data file is available at <https://doi.org/10.25338/B8X619>.

### 4.2 Phasing, IBS calling, and IBS tiling

We phased the combined dataset using Beagle 5.0 Browning and Browning (2007) using the default settings and genetic maps for each chromosome. We used linear interpolation to ob-



tain the genetic map position of each SNP on the build GRCh37 LDhat genetic map (International HapMap Consortium, 2007) downloaded from the Beagle website ([http://bochet.gcc.biostat.washington.edu/beagle/genetic\\_maps/](http://bochet.gcc.biostat.washington.edu/beagle/genetic_maps/)). We used Refined IBD software (Browning and Browning, 2013) to identify IBS segments, retaining segments of at least .1 centiMorgans (cM) with LOD scores  $>1$ . We also used Germline (Gusev et al., 2009) to identify IBS segments under alternative parameters, shown in the supplement. The resulting IBS segments were analyzed using the GenomicRanges package (Lawrence et al., 2013) in R (R Core Team, 2013). Scripts used for phasing, IBS calling, and IBS tiling are available at [github.com/mdedge/IBS\\_privacy/](https://github.com/mdedge/IBS_privacy/) website.

### 620 4.3 IBS probing

621 To generate IBS-inert genotypes for IBS probing in Figure 3, we computed allele frequencies within  
622 the set of 872 Europeans for chromosome 19. Allele frequencies less than 10% were changed to  
623 10%, and then alleles were sampled at one minus their frequency. This strategy generates genetic  
624 data that look quite unlike real data, with the advantage (for the purposes of IBS probing) of  
625 being unlikely to return IBS matches anywhere. An adversary attempting IBS probing in a real  
626 database would need to tailor the approach to the quality control and IBS calling methods used  
627 by the database.

628 After inert genotypes were produced, we stitched them with real phased genotypes from  
629 windows around GRCh position 45411941 on chromosome 19, the site of SNP rs429358. SNP  
630 rs429358 is in the APOE locus; if a haplotype has a C at rs429358 and a C at nearby SNP rs7412,  
631 then that haplotype is said to harbor the APO- $\epsilon 4$  allele, which confers risk for Alzheimer's disease  
632 Corder et al. (1993). rs429358 is not genotyped on the Human Origins chip, but it is included on  
633 recent chips used by both Ancestry and 23andMe. To simulate probing with a 1cM threshold for  
634 matching, we pulled real data from a region of 1.9cM around the site, and to simulate probing  
635 with a 3cM threshold, we pulled real data from a region of 5.9cM around the site. Distances in  
636 cM were computed by linear interpolation from a genetic map in GRCh37 coordinates. Scripts  
637 used to generate Figure 3 are available at [github.com/mdedge/IBS\\_privacy/](https://github.com/mdedge/IBS_privacy/).

### 638 4.4 GEDmatch demonstration

639 On November 21st, 2019, we first uploaded artificial genetic datasets to GEDmatch's research  
640 mode in order to demonstrate the possibility of IBS baiting. GEDMatch has not published details  
641 of its IBS detection procedures. However, the options available to users in the 1-to-1 match  
642 tool and the description of how those options can be used to ignore single-site matches led us to  
643 hypothesize that GEDmatch uses phase-unaware IBS detection and that the 1-to-1 match tool  
644 might be vulnerable to IBS baiting.

645 **Description of GEDmatch 1-to-1 tool** GEDmatch's 1-to-1 match tool allows the user to  
646 compare the IBS matches of any two genetic datasets (or, in GEDmatch parlance, "kits"), as  
647 long as the kit numbers are known to the user. Thus, to identify the genotypes of many users an  
648 adversary would need access to the kit numbers of many users. The 1-to-many tool in default  
649 GEDmatch reports 3000 of the closest genetic relatives of any kit whose number is known to the



650 user, and reports the kit numbers of those match kits (along with names and email addresses).  
651 Thus an adversary can iteratively search for all the kit numbers matching a known kit, and so  
652 obtain many kit numbers to use in 1-to-1 searches. We alerted GEDmatch to this issue with the  
653 1-to-many tool, as nearly the entire GEDmatch database of kit numbers and genetic relationships  
654 could be scraped.

655 The 1-to-1 match tool allows the user to specify parameters that govern IBS calling. In  
656 particular, the user can specify the minimum cM length of the blocks (down to 0.1cM) and  
657 the minimum number of SNPs in a block (down to 25 SNPs). GEDmatch also allows the  
658 user to specify the 'mis-match bunch limit,' which appears to be the minimum number of IBS-  
659 compatible SNPs after an opposite-homozygous site that are required in order for a second  
660 opposite-homozygous site not to break the IBS segment.

661 **Ethics.** In order to comply with GEDmatch's terms and conditions, we used artificial datasets  
662 designed not to match any genuine genetic data uploaded to GEDmatch. The kits were uploaded  
663 in "Research" mode, where they are not visible to other users via 1-to-many search. We did not  
664 interact with any other users' data; we ran GEDmatch's 1-to-1 comparison tools only comparing  
665 among our artificial kits. We exercised care not to interact with any other tools and to avoid  
666 accidental discoveries. Prior to uploading the artificial datasets, we also consulted with the UC  
667 Davis Institutional Review Board to ensure that these uploads do not constitute human subjects  
668 research. Upon receiving confirmation from the IRB that our uploads do not constitute human  
669 subjects research and before uploading the datasets, we alerted GEDmatch that we would be  
670 making the uploads, and we also shared the kit numbers with them after we had completed our  
671 analyses.

#### 672 **4.4.1 Construction of artificial datasets**

673 We constructed artificial "target" and "bait" kits using the SNPs included in the 23&Me v4 chip.  
674 (The "target" kits are the targets of inference, and the "bait" kits are designed to reveal their  
675 genotypes.) We identified the alleles at these SNP positions in the 1000Genomes dataset, along  
676 with their frequencies in the EUR subset of 1000Genomes. We assigned as missing ('- -') any SNP  
677 that we could not match by position in 1000Genomes. We chose four target SNPs at random  
678 on chromosome 22. These SNPs were chosen at random from the set of strand-unambiguous  
679 polymorphisms, i.e. not A/T and G/C SNPs. These strand-unambiguous sites include the  
680 majority of SNPs on the chip, e.g. 89% of the SNPs on the 23&Me chip on chromosome 22.

681 **Target genomes** We uploaded three artificial target genome kits (T1-T3). These vary in their  
682 genotypes at the target SNPs. T1 and T3 are homozygous for different alleles; T2 is heterozygous.  
683 At the rest of the loci, we constructed genotypes by randomly sampling alleles according to their  
684 frequencies at each SNP. Thus, there is no LD among loci.

685 **Bait genomes.** We uploaded two artificial bait genome kits. These two kits have opposite-  
686 homozygote genotypes at each of the target SNPs. The two bait uploads were then set to have  
687 identical genotypes in the rest of their autosomes, with their genotypes specified as below.

688 To create a region around the target that would bait a phase-unaware method into calling  
689 IBD, we took SNPs in the 0.6cM on either side of the target SNP, selected at random 22 on each  
690 side, and set them to be heterozygous in both bait genomes. The rest of the SNPs within this  
691 bait region were set to be missing. We used only 22 heterozygous SNPs on each side and filled  
692 in the rest with missing data (rather than making all sites heterozygous) because large numbers  
693 of heterozygous sites generated an error on upload, “HTZ string too long” and would not be  
694 processed further. Blocking uploads with long runs of heterozygous sites is a countermeasure put  
695 in place by GEDmatch after we and Ney et al. initially alerted GEDmatch to the risks of upload-  
696 based privacy attacks. However, we found that the countermeasure was not triggered by runs of  
697 heterozygous sites with missing sites interspersed, and these runs of heterozygosity interspersed  
698 with missingness also effectively baited GEDmatch into calling IBD segments. Additionally, we  
699 confirmed with Peter Ney (personal communication) that his previously-uploaded kits including  
700 long runs of heterozygosity remain active even though re-uploads of those same kits are blocked  
701 as of December 3rd, 2019, suggesting that the block applies only to newly uploaded kits and not  
702 to existing data on GEDmatch.

703 The alleles in target kits at all other autosomal SNPs in the genome were drawn at random  
704 with frequency  $1 - p$ , where  $p$  is the frequency in the 1000Genomes EUR subsample. This scheme  
705 was chosen to ensure that the bait genomes were unlikely to have spurious IBS matches anywhere  
706 with any target genome, so that the only potential IBS was in the target regions.

#### 707 4.4.2 Detailed results of baiting

708 We compared each target to both bait genomes using the 1-to-1 GEDmatch tool. We set the  
709 minimum block to a length of  $> 0.7cM$  and 25 SNPs, with a mismatch cutoff of 25 SNPs. This  
710 ensured that we could detect IBS in the key regions, but that a single opposite-homozygous  
711 mismatch would be sufficient to prevent the identification of a putative IBS segment in the key  
712 region.

713 The baiting attempt was successful; we observed IBS only where we expected it between bait  
714 and target kits (Figure 5). We observed no putative IBD segments on any chromosome except  
715 22, as expected on the basis of our procedure for filling in artificial genotypes in both sets of  
716 kits. The details of the matches on chromosome 22 are reported in Table 3. We observed 4  
717 putative IBD segments overlapping our target bait regions in the comparisons with matching  
718 homozygote genotypes at the bait site, i.e. in the T1-B1 and T3-B2 comparisons, as well in  
719 both heterozygote-homozygote comparisons, i.e. T2-B1 and T2-B2. We observed no putative  
720 IBD segments in the pairs with opposite-homozygous mismatches, T1-B1 and T3-B2. Thus the  
721 genotypes of the targets are readily discernable from from the putative IBD segments output  
722 by GEDmatch. The full results returned by GEDmatch are available as images here ([https://github.com/mdedge/IBS\\_privacy/tree/master/IBS\\_baiting\\_demo](https://github.com/mdedge/IBS_privacy/tree/master/IBS_baiting_demo); the kit numbers are  
723 redacted to prevent reuse).

724  
725 Some of the IBS blocks have fewer SNPs than we expect. We believe this to be due to the  
726 removal of SNPs during the tokenization stage, during which rare SNPs and SNPs with stand-  
727 ambiguous alleles seem to be removed Ney et al. (2020). We did not investigate this further, but  
728 multiple uploads could be used to determine the approximate criteria for SNPs to be included,  
729 and hence determine where an adversary should set cutoffs.

730 Our two bait kits could both generate IBS matches to the target because the target genotype is

Matching pairs		Target 1	Target 2	Target 3	Target 4
	target bp	27613130	34024097	37673781	42008068
T1-(B1 B?)	IBS L bp	27427698	33771672	37519864	40054428
	IBS R bp	27680780	34328741	37827711	43112674
	IBS cM	1.3	0.8	1.1	1.2
	# SNPs	47	45	42	40
	# SNPs B?	46	44	41	39
T2-(B1 B2 B?)	IBS L bp	27433179	33771672	37508507	40357667
	IBS R bp	27680780	34328741	37827711	43112674
	IBS cM	1.3	0.8	1.2	0.9
	# SNPs	45	45	45	32
	# SNPs B?	44	44	44	31
T3-(B3 B?)	IBS L bp	27433179	33771672	37519864	40357667
	IBS R bp	27680780	34328741	37827711	43112674
	IBS cM	1.3	0.8	1.1	0.9
	# SNPs	45	45	45	32
	# SNPs B?	44	44	41	31
T?-(All Baits)	IBS L bp	27433179	33771672	37519864	40357667
	IBS R bp	27680780	34328741	37827711	43112674
	IBS cM	1.3	0.8	1.1	0.9
	# SNPs	44	44	44	31
	# SNPs B?	44	44	44	31

Table 3: Summary of the SNPs targeted by baiting and the IBS returned by GEDmatch. For each region, we give the position of the key SNP (target bp). Because by design our bait kits are genetically identical outside of the target SNPs, the IBS regions returned by GEDmatch's 1-to-1 tool are identical across bait kits generating a match. For each pairwise comparison, we report the IBS information returned: Left-Right bp of the IBS region, the cM length, the number (#) of SNPs in the IBS region with a non-missing target. We also report the number (#) of SNPs spanned by the region IBS when matched to the missing target Bmiss.

731 missing rather than heterozygous. To determine whether a genotype was missing, we implemented  
732 a trick borrowed from Ney et al. (2020), and uploaded a third bait kit (Bmiss) with the target  
733 SNP set to missing (i.e. '- -') and then looked at the number of SNPs an IBS match across  
734 the target site spans. In each case, the non-missing baits (B1 and B2) generated an IBS block  
735 match with with T1-T3 that was one SNP longer than the IBS block generated by the Bmiss bait  
736 (Table 3). Comparing these baits to a new target with a missing genotype at each target site  
737 (T?), we see that in each pairwise comparison the IBS blocks are the same number of SNPs long  
738 regardless of whether the target SNP bait genotype was missing (Table 3). Therefore, we can  
739 distinguish the target being heterozygote or missing by the use of a third bait kit and inspection  
740 of the number of SNPs included in a IBS match.

741 The possibility of IBS-baiting-like procedures also interacts with the vulnerabilities arising from  
742 the presentation of SNP-level visualizations explored by Ney et al.. Even if short IBS blocks were  
743 not reported to the user explicitly, it is clear from the zoomed-in view that we can see the target  
744 mismatches in question (see Figure 5). One measure that GEDmatch appears to have taken  
745 against a Ney et al.-style attack is to jitter the positions of SNPs in their visualization slightly.  
746 However, an attacker could counter such jittering by embedding key sites in runs of heterozygosity,  
747 making it easier to identify them in visualizations after jittering. Thus, the images displayed by  
748 GEDmatch still pose additional security risks.

## 749 Acknowledgments

750 We thank Matt Bishop, Elizabeth Joh, Peter Ney, and Mike Sweeney for useful conversations,  
751 and we thank Shai Carmi, Yaniv Erlich, Debbie Kennett, Leah Larkin, Magnus Nordborg, Rori  
752 Rohlf, Noah Rosenberg, Ann Turner, Amy Williams, and an anonymous reviewer for helpful  
753 comments on the manuscript. Swapan Mallick and David Reich answered questions about the  
754 Human Origins dataset, Brian Browning answered questions about Refined IBD, and Alexander  
755 Gusev answered questions about Germline software. We acknowledge support from the National  
756 Institutes of Health (R01-GM108779 and F32-GM130050).

## 757 References

- 758 1000 Genomes Project Consortium (2012). An integrated map of genetic variation from 1,092  
759 human genomes. *Nature* 491(7422), 56.
- 760 Ayday, E. and M. Humbert (2017). Inference attacks against kin genomic privacy. *IEEE Security*  
761 *& Privacy* 15(5), 29–37.
- 762 Bélisle-Pipon, J.-C., E. Vayena, R. C. Green, and I. G. Cohen (2019). Genetic testing, insurance  
763 discrimination and medical research: what the united states can learn from peer countries.  
764 *Nature Medicine* 25(8), 1198–1204.
- 765 Bjelland, D. W., U. Lingala, P. S. Patel, M. Jones, and M. C. Keller (2017). A fast and accurate  
766 method for detection of ibd shared haplotypes in genome-wide snp data. *European Journal of*  
767 *Human Genetics* 25(5), 617.

- 768 Browning, B. L. and S. R. Browning (2013). Improving the accuracy and efficiency of identity-  
769 by-descent detection in population data. *Genetics* 194(2), 459–471.
- 770 Browning, S. R. and B. L. Browning (2007). Rapid and accurate haplotype phasing and missing-  
771 data inference for whole-genome association studies by use of localized haplotype clustering.  
772 *The American Journal of Human Genetics* 81(5), 1084–1097.
- 773 Browning, S. R. and B. L. Browning (2012). Identity by descent between distant relatives:  
774 detection and applications. *Annual Review of Genetics* 46, 617–633.
- 775 Bryc, K., E. Durand, J. Macpherson, D. Reich, and J. Mountain (2015). The genetic ancestry  
776 of african americans, latinos, and european americans across the united states. *The American*  
777 *Journal of Human Genetics* 96(1), 37 – 53.
- 778 Buffalo, V., S. M. Mount, and G. Coop (2016). A genealogical look at shared ancestry on the x  
779 chromosome. *Genetics* 204(1), 57–75.
- 780 Cann, H. M., C. De Toma, L. Cazes, M.-F. Legrand, V. Morel, L. Piouffre, J. Bodmer, W. F.  
781 Bodmer, B. Bonne-Tamir, A. Cambon-Thomsen, et al. (2002). A human genome diversity cell  
782 line panel. *Science* 296(5566), 261–262.
- 783 Carmi, S., K. Y. Hui, E. Kochav, X. Liu, J. Xue, F. Grady, S. Guha, K. Upadhyay, D. Ben-  
784 Avraham, S. Mukherjee, et al. (2014). Sequencing an ashkenazi reference panel supports  
785 population-targeted personal genomics and illuminates jewish and european origins. *Nature*  
786 *Communications* 5, 4835.
- 787 Carmi, S., P. F. Palamara, V. Vacic, T. Lencz, A. Darvasi, and I. Pe'er (2013). The variance of  
788 identity-by-descent sharing in the wright–fisher model. *Genetics* 193(3), 911–928.
- 789 Clayton, E. W., B. J. Evans, J. W. Hazel, and M. A. Rothstein (2019). The law of genetic  
790 privacy: applications, implications, and limitations. *Journal of Law and the Biosciences*, 1–36.
- 791 Conomos, M., A. Reiner, B. Weir, and T. Thornton (2016). Model-free estimation of recent  
792 genetic relatedness. *The American Journal of Human Genetics* 98(1), 127 – 148.
- 793 Corder, E. H., A. M. Saunders, W. J. Strittmatter, D. E. Schmechel, P. C. Gaskell, G. Small,  
794 A. D. Roses, J. Haines, and M. A. Pericak-Vance (1993). Gene dose of apolipoprotein e type  
795 4 allele and the risk of alzheimer's disease in late onset families. *Science* 261(5123), 921–923.
- 796 Court, D. S. (2018). Forensic genealogy: Some serious concerns. *Forensic Science International:*  
797 *Genetics* 36, 203–204.
- 798 Danecek, P., A. Auton, G. Abecasis, C. A. Albers, E. Banks, et al. (2011). The variant call  
799 format and vcftools. *Bioinformatics* 27(15), 2156–2158.
- 800 Donnelly, K. P. (1983). The probability that related individuals share some section of genome  
801 identical by descent. *Theoretical Population Biology* 23(1), 34–63.

- 802 Durand, E. Y., N. Eriksson, and C. Y. McLean (2014, 04). Reducing Pervasive False-Positive  
803 Identical-by-Descent Segments Detected by Large-Scale Pedigree Analysis. *Molecular Biology*  
804 *and Evolution* 31(8), 2212–2222.
- 805 Edge, M. D., B. F. B. Algee-Hewitt, T. J. Pemberton, J. Z. Li, and N. A. Rosenberg (2017).  
806 Linkage disequilibrium matches forensic genetic records to disjoint genomic marker sets. *Pro-*  
807 *ceedings of the National Academy of Sciences* 114(22), 5671–5676.
- 808 Edge, M. D. and G. Coop (2019). How lucky was the genetic investigation in the golden state  
809 killer case? *bioRxiv*.
- 810 Erlich, Y. and A. Narayanan (2014). Routes for breaching and protecting genetic privacy. *Nature*  
811 *Reviews Genetics* 15(6), 409.
- 812 Erlich, Y., T. Shor, I. Pe’er, and S. Carmi (2018). Identity inference of genomic data using  
813 long-range familial searches. *Science* 362(6415), 690–694.
- 814 Greshake, B., P. E. Bayer, H. Rausch, and J. Reda (2014). Opensnp—a crowdsourced web resource  
815 for personal genomics. *PLoS One* 9(3), e89204.
- 816 Greytak, E., D. H. Kaye, B. Budowle, C. Moore, and S. Armentrout (2018). Privacy and genetic  
817 genealogy data. *Science* 361, 857.
- 818 Gusev, A., J. K. Lowe, M. Stoffel, M. J. Daly, D. Altshuler, J. L. Breslow, J. M. Friedman,  
819 and I. Pe’er (2009). Whole population, genome-wide mapping of hidden relatedness. *Genome*  
820 *Research* 19(2), 318–326.
- 821 Gymrek, M., A. L. McGuire, D. Golan, E. Halperin, and Y. Erlich (2013). Identifying personal  
822 genomes by surname inference. *Science* 339(6117), 321–324.
- 823 Han, E., P. Carbonetto, R. E. Curtis, Y. Wang, J. M. Granka, J. Byrnes, K. Noto, A. R. Kermany,  
824 N. M. Myres, M. J. Barber, et al. (2017). Clustering of 770,000 genomes reveals post-colonial  
825 population structure of north america. *Nature Communications* 8, 14238.
- 826 He, D., N. A. Furlotte, F. Hormozdiari, J. W. J. Joo, A. Wadia, R. Ostrovsky, A. Sahai, and  
827 E. Eskin (2014). Identifying genetic relatives without compromising privacy. *Genome Re-*  
828 *search* 24(4), 664–672.
- 829 Henn, B. M., L. Hon, J. M. Macpherson, N. Eriksson, S. Saxonov, I. Pe’er, and J. L. Mountain  
830 (2012, 04). Cryptic distant relatives are common in both isolated and cosmopolitan genetic  
831 samples. *PLOS ONE* 7(4), 1–13.
- 832 Hogarth, S., G. Javitt, and D. Melzer (2008). The current landscape for direct-to-consumer  
833 genetic testing: Legal, ethical, and policy issues. *Annual Review of Genomics and Human*  
834 *Genetics* 9(1), 161–182. PMID: 18767961.
- 835 Hogarth, S. and P. Saukko (2017). A market in the making: the past, present and future of  
836 direct-to-consumer genomics. *New Genetics and Society* 36(3), 197–208.



- 837 Homer, N., S. Szelinger, M. Redman, D. Duggan, W. Tembe, J. Muehling, J. V. Pearson, D. A.  
838 Stephan, S. F. Nelson, and D. W. Craig (2008). Resolving individuals contributing trace  
839 amounts of dna to highly complex mixtures using high-density snp genotyping microarrays.  
840 *PLoS Genetics* 4(8), e1000167.
- 841 Hon, L., S. Saxonov, B. T. Naughton, J. L. Mountain, A. Wojcicki, and L. Avey (2013, June 11).  
842 Finding relatives in a database. US Patent 8,463,554.
- 843 Huang, L., S. Bercovici, J. M. Rodriguez, and S. Batzoglou (2014). An effective filter for ibd  
844 detection in large data sets. *PLoS ONE* 9(3), e92713.
- 845 Huff, C. D., D. J. Witherspoon, T. S. Simonson, J. Xing, W. S. Watkins, Y. Zhang, T. M. Tuohy,  
846 D. W. Neklason, R. W. Burt, S. L. Guthery, et al. (2011). Maximum-likelihood estimation of  
847 recent shared ancestry (ersa). *Genome Research* 21(5), 768–774.
- 848 Humbert, M., E. Ayday, J.-P. Hubaux, and A. Telenti (2013). Addressing the concerns of the  
849 lacks family: quantification of kin genomic privacy. In *Proceedings of the 2013 ACM SIGSAC*  
850 *conference on Computer & Communications Security*, pp. 1141–1152. ACM.
- 851 Humbert, M., K. Huguenin, J. Hugonot, E. Ayday, and J.-P. Hubaux (2015). De-anonymizing  
852 genomic databases using phenotypic traits. *Proceedings on Privacy Enhancing Technolo-*  
853 *gies 2015*(2), 99–114.
- 854 Im, H. K., E. R. Gamazon, D. L. Nicolae, and N. J. Cox (2012). On sharing quantitative trait  
855 gwas results in an era of multiple-omics data and the limits of genomic privacy. *The American*  
856 *Journal of Human Genetics* 90(4), 591–598.
- 857 International HapMap Consortium (2007). A second generation human haplotype map of over  
858 3.1 million snps. *Nature* 449(7164), 851.
- 859 Kennett, D. (2019). Using genetic genealogy databases in missing persons cases and to develop  
860 suspect leads in violent crimes. *Forensic Science International* 301, 107 – 117.
- 861 Kere, J. (2001). Human population genetics: Lessons from finland. *Annual Review of Genomics*  
862 *and Human Genetics* 2(1), 103–128. PMID: 11701645.
- 863 Khan, R. and D. Mittelman (2018). Consumer genomics will change your life, whether you get  
864 tested or not. *Genome Biology* 19(1), 120.
- 865 Kim, J., M. D. Edge, B. F. Algee-Hewitt, J. Z. Li, and N. A. Rosenberg (2018). Statistical  
866 detection of relatives typed with disjoint forensic and biomedical loci. *Cell* 175(3), 848 –  
867 858.e6.
- 868 Larkin, L. (2017, Mar). Cystic fibrosis: A case study in genetic privacy.
- 869 Larkin, L. (2018, Sept). Database sizes—september 2018 update.
- 870 Lawrence, M., W. Huber, H. Pages, P. Aboyoun, M. Carlson, R. Gentleman, M. T. Morgan,  
871 and V. J. Carey (2013). Software for computing and annotating genomic ranges. *PLoS*  
872 *Computational Biology* 9(8), e1003118.

- 873 Li, H. (2011, 09). A statistical framework for SNP calling, mutation discovery, association  
874 mapping and population genetical parameter estimation from sequencing data. *Bioinformat-*  
875 *ics* 27(21), 2987–2993.
- 876 Loh, P.-R., P. F. Palamara, and A. L. Price (2016). Fast and accurate long-range phasing in a  
877 uk biobank cohort. *Nature Genetics* 48(7), 811.
- 878 McQuillan, R., A.-L. Leutenegger, R. Abdel-Rahman, C. S. Franklin, M. Pericic, L. Barac-Lauc,  
879 N. Smolej-Narancic, B. Janicijevic, O. Polasek, A. Tenesa, et al. (2008). Runs of homozygosity  
880 in european populations. *The American Journal of Human Genetics* 83(3), 359–372.
- 881 Naveed, M., E. Ayday, E. W. Clayton, J. Fellay, C. A. Gunter, J.-P. Hubaux, B. A. Malin,  
882 and X. Wang (2015, August). Privacy in the genomic era. *ACM Computing Surveys* 48(1),  
883 6:1–6:44.
- 884 Ney, P., L. Ceze, and T. Kohno (2020). Genotype extraction and false relative attacks: Se-  
885 curity risks to third-party genetic genealogy services beyond identity inference. Network  
886 and Distributed System Security Symposium (NDSS). Preprint Posted 10/29/19, [https://dnasec.cs.washington.edu/genetic-genealogy/ney\\_ndss.pdf](https://dnasec.cs.washington.edu/genetic-genealogy/ney_ndss.pdf).  
887
- 888 Ney, P. M., L. Ceze, and T. Kohno (2018). Computer security risks of distant relative matching  
889 in consumer genetic databases. *CoRR abs/1810.02895*.
- 890 Nyholt, D. R., C.-E. Yu, and P. M. Visscher (2009). On jim watson's apoe status: genetic  
891 information is hard to hide. *European Journal of Human Genetics* 17(2), 147.
- 892 Palamara, P. F., T. Lencz, A. Darvasi, and I. Pe'er (2012). Length distributions of identity by  
893 descent reveal fine-scale demographic history. *The American Journal of Human Genetics* 91(5),  
894 809–822.
- 895 Panoutsopoulou, K., K. Hatzikotoulas, D. K. Xifara, V. Colonna, A.-E. Farmaki, G. R. Ritchie,  
896 L. Southam, A. Gilly, I. Tachmazidou, S. Fatumo, et al. (2014). Genetic characterization of  
897 greek population isolates reveals strong genetic drift at missense and trait-associated variants.  
898 *Nature Communications* 5, 5345.
- 899 Patterson, N., P. Moorjani, Y. Luo, S. Mallick, N. Rohland, Y. Zhan, T. Genschoreck, T. Webster,  
900 and D. Reich (2012). Ancient admixture in human history. *Genetics* 192(3), 1065–1093.
- 901 Pemberton, T. J., D. Absher, M. W. Feldman, R. M. Myers, N. A. Rosenberg, and J. Z. Li  
902 (2012). Genomic patterns of homozygosity in worldwide human populations. *The American*  
903 *Journal of Human Genetics* 91(2), 275–292.
- 904 Price, A. L., N. J. Patterson, R. M. Plenge, M. E. Weinblatt, N. A. Shadick, and D. Reich  
905 (2006). Principal components analysis corrects for stratification in genome-wide association  
906 studies. *Nature Genetics* 38(8), 904.
- 907 Purcell, S., B. Neale, K. Todd-Brown, L. Thomas, M. A. Ferreira, D. Bender, J. Maller, P. Sklar,  
908 P. I. de Bakker, M. J. Daly, and P. C. Sham (2007). Plink: A tool set for whole-genome associ-  
909 ation and population-based linkage analyses. *The American Journal of Human Genetics* 81(3),  
910 559 – 575.

- 911 R Core Team (2013). *R: A Language and Environment for Statistical Computing*. Vienna, Austria:  
912 R Foundation for Statistical Computing.
- 913 Ralph, P. and G. Coop (2013). The geography of recent genetic ancestry across europe. *PLoS*  
914 *Biology* 11(5), e1001555.
- 915 Ram, N., C. J. Guerrini, and A. L. McGuire (2018). Genealogy databases and the future of  
916 criminal investigation. *Science* 360(6393), 1078–1079.
- 917 Ramstetter, M. D., T. D. Dyer, D. M. Lehman, J. E. Curran, R. Duggirala, J. Blangero, J. G.  
918 Mezey, and A. L. Williams (2017). Benchmarking relatedness inference methods with genome-  
919 wide data from thousands of relatives. *Genetics* 207(1), 75–82.
- 920 Ramstetter, M. D., S. A. Shenoy, T. D. Dyer, D. M. Lehman, J. E. Curran, R. Duggirala,  
921 J. Blangero, J. G. Mezey, and A. L. Williams (2018). Inferring identical-by-descent sharing of  
922 sample ancestors promotes high-resolution relative detection. *The American Journal of Human*  
923 *Genetics* 103(1), 30–44.
- 924 Regalado, A. (2019). More than 26 million people have taken an at-home ancestry test. *MIT*  
925 *Technology Review*.
- 926 Rohlf, R. V., S. M. Fullerton, and B. S. Weir (2012, 02). Familial identification: Population  
927 structure and relationship distinguishability. *PLOS Genetics* 8(2), 1–13.
- 928 Scudder, N., D. McNevin, S. F. Kelty, C. Funk, S. J. Walsh, and J. Robertson (2019). Policy and  
929 regulatory implications of the new frontier of forensic genomics: direct-to-consumer genetic  
930 data and genealogy records. *Current Issues in Criminal Justice* 31(2), 194–216.
- 931 Shi, S., N. Yuan, M. Yang, Z. Du, J. Wang, X. Sheng, J. Wu, and J. Xiao (2018). Comprehensive  
932 assessment of genotype imputation performance. *Human Heredity* 83(3), 107–116.
- 933 Shringarpure, S. S. and C. D. Bustamante (2015). Privacy risks from genomic data-sharing  
934 beacons. *The American Journal of Human Genetics* 97(5), 631–646.
- 935 Staples, J., D. J. Witherspoon, L. B. Jorde, D. A. Nickerson, J. E. Below, C. D. Huff,  
936 U. of Washington Center for Mendelian Genomics, et al. (2016). Padre: pedigree-aware distant-  
937 relationship estimation. *The American Journal of Human Genetics* 99(1), 154–162.
- 938 Thompson, E. A. (2013). Identity by descent: variation in meiosis, across genomes, and in  
939 populations. *Genetics* 194(2), 301–326.

## Appendix A Detailed rationale for proposed countermeasures

Here, we detail the rationale and possible advantages and disadvantages of the countermeasures listed in Table 2.

- 1. Require uploaded files to include cryptographic signatures identifying their source.** This recommendation was initially made by Erlich et al. (2018). Under this suggestion, DTC genetics services would cryptographically sign the genetic data files they provide to users. Upload services might then check for a signature from an approved DTC service on each uploaded dataset, blocking datasets from upload otherwise. An alternative procedure that would accomplish the same goal would be for the DTC entities to exchange data directly at the user's request (Ney et al., 2018). Such a procedure would allow upload services to know the source of the files they analyze and to disallow uploaded datasets produced by non-approved entities and user-modified datasets. All the methods we describe require the upload of multiple genetic datasets. As such, requiring cryptographic signatures would force the adversary to have multiple biological samples analyzed by a DTC service in order to implement any of our procedures, and IBS probing and IBS baiting would require synthetic samples, which are much harder to produce than fake datasets. Another benefit of this approach is that it would protect research participants against being reidentified using DTC genetic genealogy services (Erlich et al., 2018). A disadvantage of this strategy is that it requires the cooperation of several distinct DTC services.
- 2. Restrict reporting of IBS to long segments.** Reporting short IBS segments increases the typical coverage of IBS tiling (Figure 2) and IBS probing (3), as well as the efficiency of IBS baiting. Very short blocks may be of little practical utility for genetic genealogy (Huff et al., 2011). Reporting only segments longer than 8 cM would substantially limit IBS tiling attacks. A partially effective variant of this strategy is to report short segments only for pairs of people who share at least one long segment (Figure S2). One disadvantage is that short segments, though less reliably inferred than longer segments, may still be of interest to genealogists.
- 3. Do not report locations of IBS segments.** Another tactic for preventing IBS tiling is not to report chromosomal locations at all. If chromosomal locations are not reported, IBS tiling as we have described it becomes impossible.
- 4. Block uploads of genomes with excessive homozygosity.** IBS tiling is especially informative if genotypes that are homozygous for phased haplotypes are uploaded, so blocking genomes with excessive homozygosity presents a barrier to IBS tiling attacks. However, runs of homozygosity occur naturally (Pemberton et al., 2012), and allowing for naturally occurring patterns of homozygosity would leave a loophole for an adversary who could upload many genotypes, using including homozygous regions and using only those for tiling.
- 5. Report only a small number of putative relatives per uploaded kit.** Reporting only the closest relatives (say the  $\approx 50$  closest relatives) of an uploaded kit would decrease the efficiency of all the methods we describe here. Only a small number of people could

- 980 have their privacy compromised by each upload. This countermeasure comes with costs to  
981 genealogists, who may want to explore as many matches as possible in order to build family  
982 trees.
- 983 6. **Disallow arbitrary matching between kits.** Some services allow searches for IBS be-  
984 tween any pair of individuals in the database. Allowing such searches makes all potential  
985 IBS attacks easier. This countermeasure might hamper the investigations of genealogists  
986 exploring complex hypotheses about relatedness.
- 987 7. **Block uploads of publicly available genomes.** There are now thousands of genomes  
988 available for public download, and these publicly available genomes can be used for IBS  
989 tiling. Genetic genealogy databases could include publicly available genomes (potentially  
990 without allowing them to be returned as IBS matches for typical users) and flag accounts  
991 that upload them. This strategy would go some distance toward blocking IBS tiling, but it  
992 could be thwarted in several ways, for example by uploading genetic datasets produced by  
993 splicing together haplotypes from publicly available genomes.
- 994 8. **Block uploads with evidence of IBS-inert segments.** IBS-inert segments—i.e. false  
995 genetic segments designed to be unlikely to be IBS with anyone in the database—are key  
996 to IBS probing. Some methods for constructing IBS-inert segments are easy to identify,  
997 but others may not be. If a database is large enough, genomes with IBS-inert segments  
998 could be identified by looking for genomes that have much less apparent IBS with other  
999 database members than might be expected.
- 1000 9. **Block uploads with long runs of heterozygosity.** Long runs of heterozygosity do  
1001 not arise naturally but are key to the IBS baiting approaches we describe here. Blocking  
1002 genomes with long runs of heterozygosity—or alternatively, blocking genomes that have  
1003 much more apparent IBS with a range of other database members than expected—would  
1004 hamper IBS baiting. However, this countermeasure might be hard to apply to a small-scale  
1005 IBS baiting attack, where only one or a few short runs of heterozygosity might be necessary.  
1006 In our sample, the longest run of heterozygosity (in terms of number of SNPs) consisted of  
1007 38 SNPs and spanned .06 cM. This suggests that filtering out long runs of heterozygosity  
1008 might be a promising strategy, though identifying a specific procedure would require more  
1009 careful consideration of variation in non-European populations and of the composition of  
1010 commercial SNP chips (including SNP density and allele frequencies).
- 1011 10. **Use phase-aware methods for IBS detection.** Although calling IBS by looking for long  
1012 segments without incompatible homozygous genotypes scales well to large datasets, such  
1013 methods are easy to trick, allowing IBS baiting approaches. In addition to allowing IBS  
1014 estimation methods that are harder to trick, faked samples may stand out as unusual during  
1015 the process of phasing, raising more opportunities for quality-control checks.

## 1016 Supplementary Figures

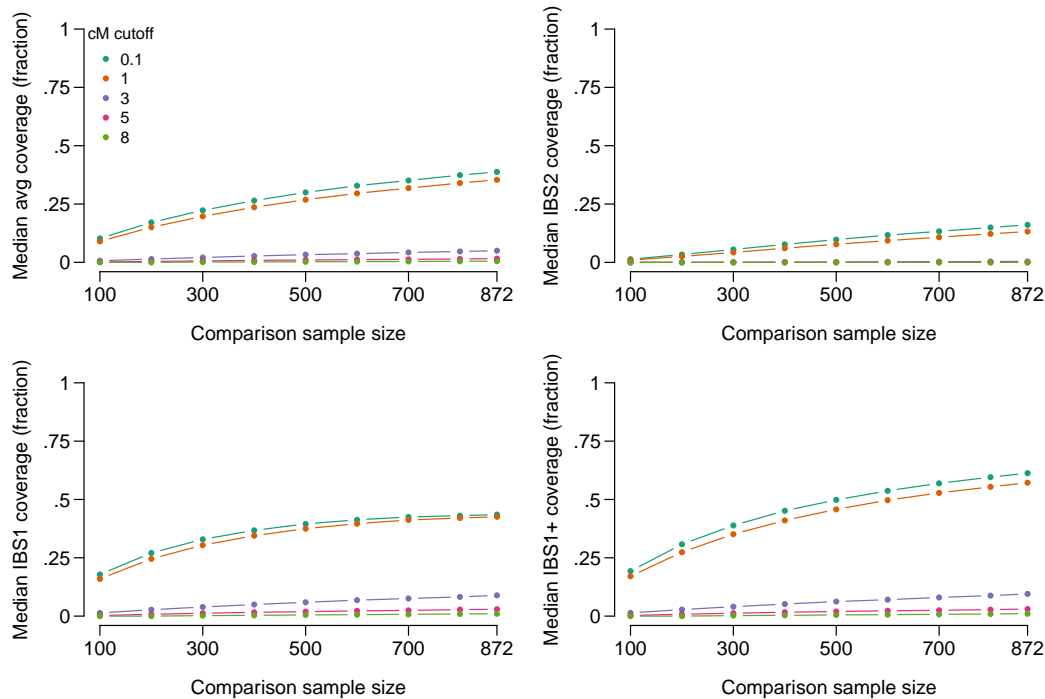


Figure S1: Tiling performance with IBS segments that are unlikely to be IBD filtered out. Conventions are the same as in Figure 2; the difference is that now only IBS segments that represent likely IBD (LOD score  $> 3$ ) are included. As expected, the amount of tiling possible is reduced when the LOD score threshold is increased, particularly when segments as short as 1 cM are allowed. However, tiling still reveals a substantial amount of information about target genotypes. Using a comparison sample of 871, and including all called IBS segments  $>1$  cM, the median person has an average of 35% of the maximum length of 2.8 Gbp covered by IBD segments with LOD  $>3$ , and has at least one chromosome covered for approximately 57% of the genome. If only segments  $>3$  cM are included, then averaging across the two chromosomes, median coverage is 5.0%, and the median proportion for which at least one chromosome is covered is 9.5%. As before, the percentage of the genome recoverable by tiling varies among people, and some people still have large proportions of their genetic data recoverable by tiling. With a LOD score threshold of 3, the top 10% of people have at least 58% of their total genotype information covered by IBD tiles, including one or more alleles at sites in at least 81% of the genome covered by IBD tiles.



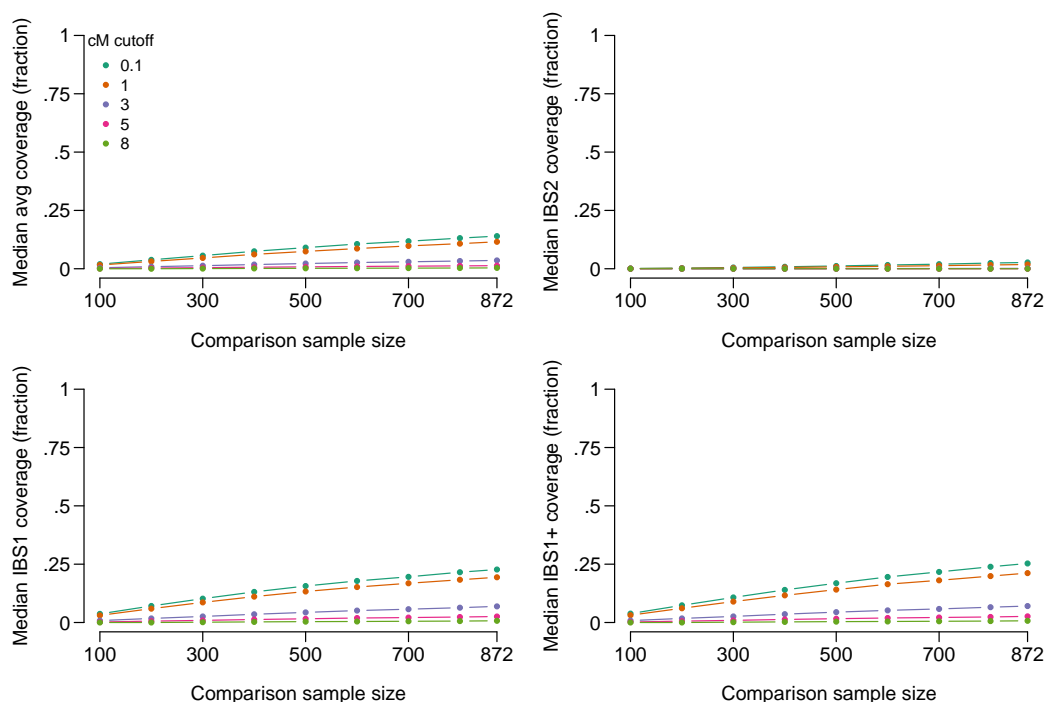


Figure S2: IBS tiling performance, limiting to comparison samples who share at least 1 IBS segment of 8 cM or more with the target. Conventions are the same as in Figure 2. Some DTC genetics companies use a two-step approach for reporting IBS information to users. For example, at this writing, MyHeritage identifies people who are likely matches of a given user as all those who share an apparent IBD segment of at least 8 cM with the user. However, once matches are identified, inferred IBD segments down to a minimum length of 6 cM are reported to the user (see Table 1). Similarly, FamilyTreeDNA only reports matching segments for pairs of people who pass a sharing threshold, and for those pairs of individuals they report all matches down to 1cM. As expected, reporting only IBS segments for pairs of people who share at least one long IBS segment ( $>8$  cM) substantially reduces but does not eliminate the effectiveness of IBS tiling. With 872 comparison samples, the median person has approximately 12% of their genome covered by IBS tiles of 1 cM or more (averaged across both chromosomes) and at least one chromosome covered for 21% of the genome. People in the top 10% of IBS tiling coverage have 44% of their genome length recoverable by tiling (averaging across both chromosomes), with at least one chromosome tiled over more than 67% of the genome. Importantly, the practice of requiring at least one long IBS match in order to report any IBS segments will not reduce the effectiveness of IBS tiling if phase-unaware methods are used for calling IBS. In that case, the attacker could simply insert a long run of heterozygous sites in each of the genomic datasets uploaded, causing an apparent long run of IBS with every user in the database (see section 2.3). After getting "in the door" with a long run of heterozygous sites, the attacker could then use tiling to find out about the rest of the genome.

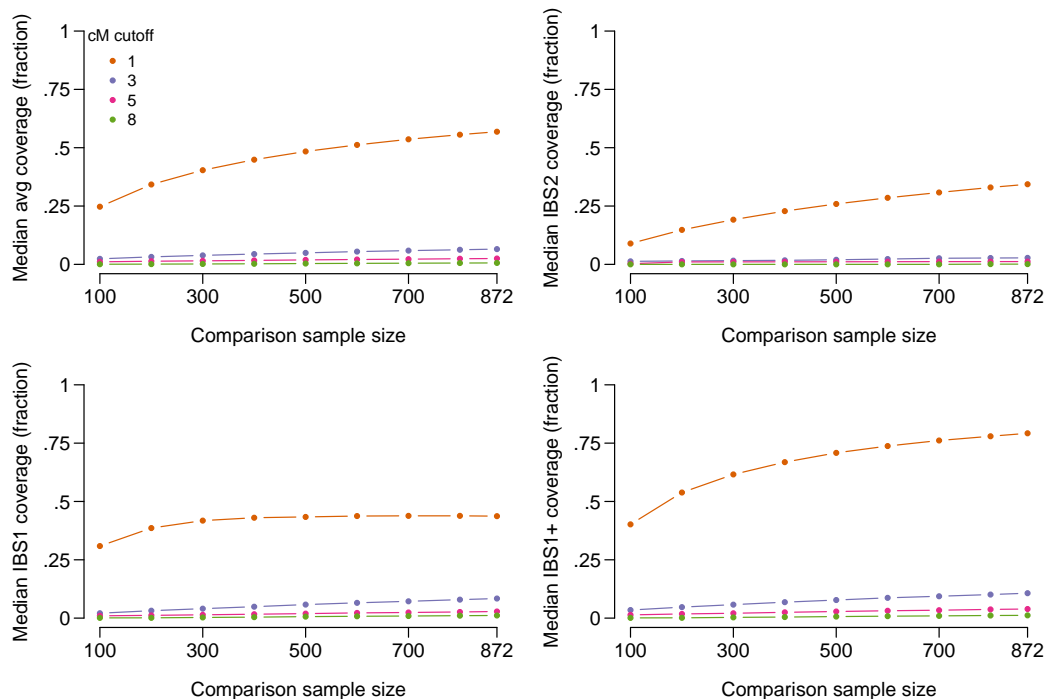


Figure S3: IBS tiling performance when genotype phasing switches are disallowed. Conventions are the same as in the Figure 2. We called IBS segments using Germline (Gusev et al., 2009), using the haploid flag to find IBS segments within the phased chromosomes produced by Beagle. We also set the `err_hom` argument to zero, set the `bits` argument to 32 to increase sensitivity for short segments, used the `w_extend` flag to extend segments beyond the slices produced by Germline, and set the minimum IBS segment length to 1cM. (Setting the minimum segment length shorter than 1cM does not appear to be possible in Germline.) The amount of tiling possible is reduced somewhat when phase switches are disallowed. However, tiling still reveals substantial information about target genotypes. Using a comparison sample of 871, and including all called IBS segments  $>1$  cM, the median person has an average of 57% of the maximum length of 2.8 Gbp covered by IBS segments, and has at least one chromosome covered for approximately 79% of the genome. If only segments  $>3$  cM are included, then averaging across the two chromosomes, median coverage is 6.5%, and the median proportion for which at least one chromosome is covered is 11%. The top 10% of people have at least 73% of their genomes covered by IBS tiles of 1 cM or more, including one or more alleles at sites in at least 91% of the genome covered by IBS tiles.

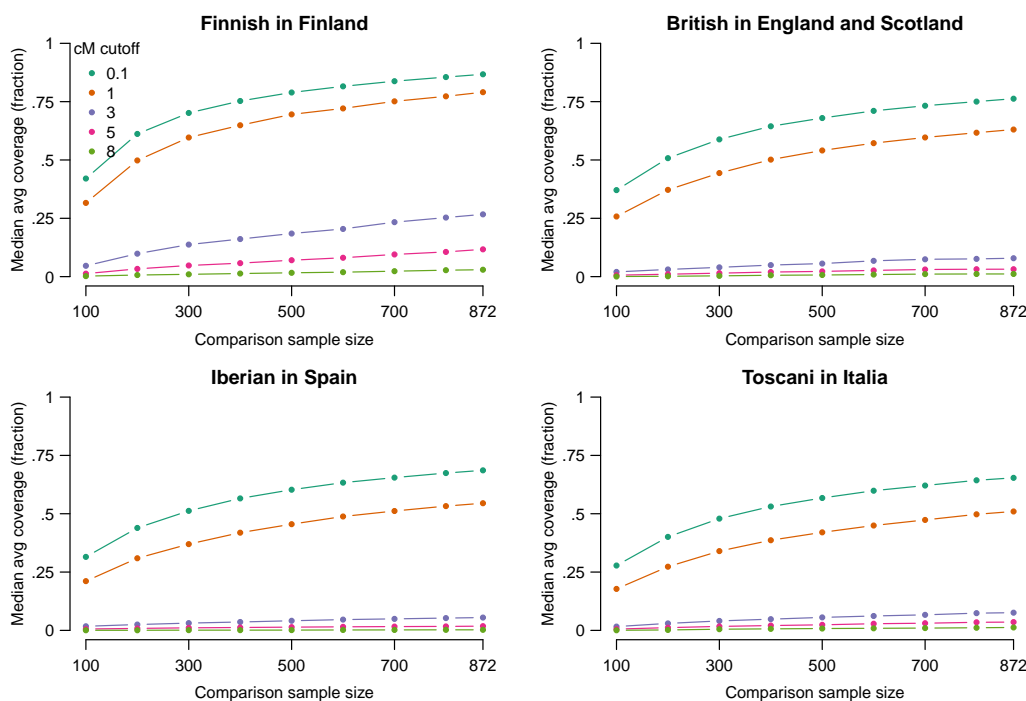


Figure S4: IBS tiling performance in selected populations. We examined IBS tiling performance in four European populations from the 1000Genomes data—Finnish in Finland (FIN, top left, 99 people), British in England and Scotland (GBR, top right, 91 people), Iberian in Spain (IBS, bottom left, 107 people), and Toscani in Italia (TSI, bottom left, 107 people). In all panels, median average IBS tiling coverage is shown on the vertical axis using refinedIBD with  $\text{LOD} > 1$ , as in the top-left panel of Figure 2. Median tiling accuracy varies among populations. For example, using IBS tiles  $> 1\text{cM}$  and with all 871 other individuals used in the comparison sample, the median coverage percentages by population were 79% (FIN), 63% (GBR), 55% (IBS), and 51% (TSI). The most striking population difference is the higher IBS tiling rates attained among Finns, who have long been of interest as a founder population, having experienced a bottleneck approximately 100 generations ago (Kere, 2001). Another factor that likely influences these results is the genetic similarity of members of each population to members of other populations included—for example, the CEU population of 1000Genomes is closely related to the GBR population, and the inclusion of CEU may partially account for the higher tiling rates in GBR than in IBS or TSI.

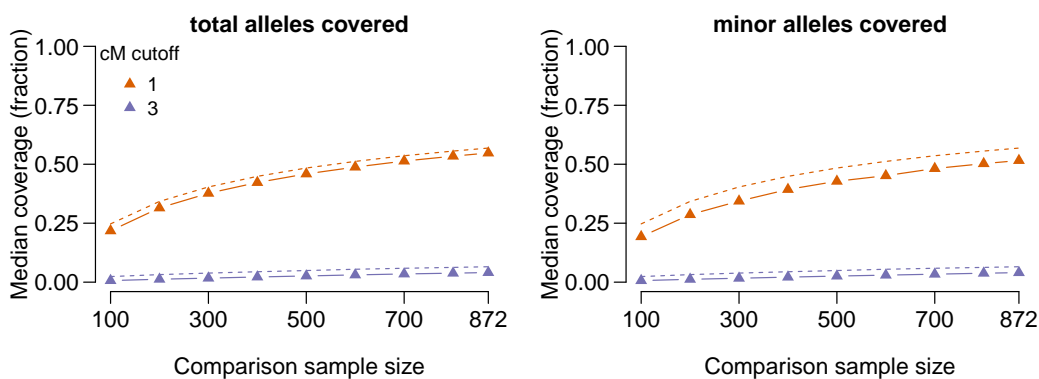


Figure S5: IBS tiling performance in terms of number of total alleles covered (left panel) and number of minor alleles covered (right panel, 18.6% of total alleles were minor alleles). We used Germline in haploid mode (as in Figure S3), as it allows easier identification of which allele is covered by a given IBS segment. Dashed lines show the results in terms of fraction of base pairs covered, whereas the solid lines show results in terms of alleles covered. Results for cM cutoffs  $<1$  are not shown because they cannot be run in Germline, and results for cM cutoffs  $>3$  are not shown because it is difficult to distinguish the dashed and solid lines. As would be expected, there is a slight bias for IBS tiles to fall in regions with lower SNP density, leading to slightly fewer alleles on the chip being covered than would be expected on the basis of total base pairs covered. For example, with a 1cM cutoff and all samples included, the median is 57% of the genome length in base pairs covered by IBS tiles, whereas the median proportion of total alleles covered by IBS tiles is 55%. It also appears that IBS tiles may be more likely to appear in regions that are less genetically diverse, as the proportion of minor alleles tiled is slightly lower than the proportion of total alleles covered. For example, with a 1cM cutoff and all samples included, the proportion of minor alleles covered by IBS tiles is 52%.

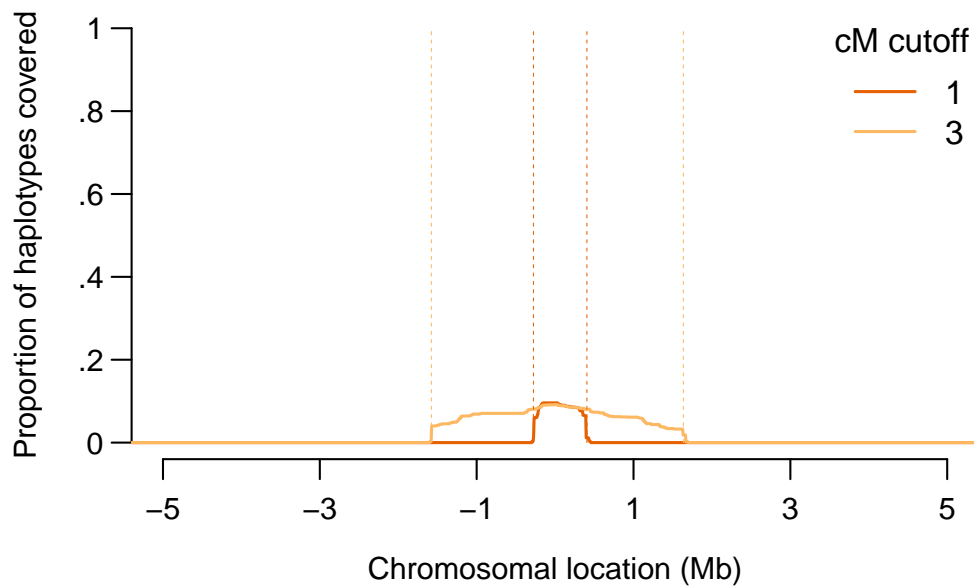


Figure S6: A demonstration of the IBD probing method around position 45411941 on chromosome 19 (GRCh37 coordinates), in the APOE locus. Conventions are the same as in Figure 3; the difference is that only IBS segments with a LOD score  $>3$  for IBD are included. When IBD probing is performed with a 1-cM threshold, 9.6% of haplotypes had a match among the probes constructed from the other 871 people in the dataset. With a 3-cM threshold, 9.2% of haplotypes had a match.

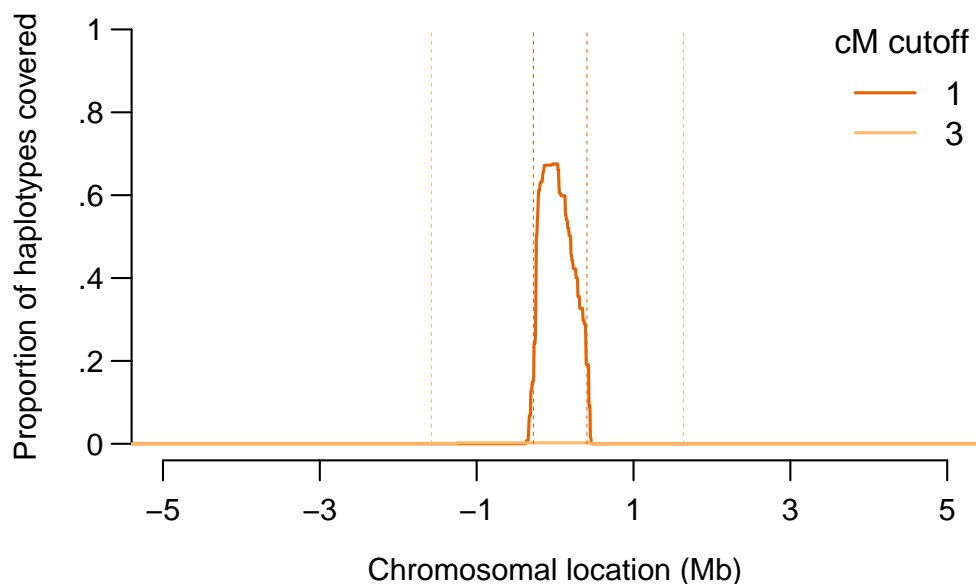


Figure S7: A demonstration of the IBS probing method around position 45411941 on chromosome 19 (GRCh37 coordinates), in the APOE locus. Conventions are the same as in Figure 3; the difference is that IBS calling was performed by Germline (Gusev et al., 2009) in haploid mode, meaning that phasing switches are disallowed. We set the `err_hom` argument to zero, we used the `w_extend` flag to extend segments beyond the slices produced by Germline, and we set the minimum IBS segment length to 1cM. All other arguments were kept at their default values. When IBS probing is performed with a 1-cM threshold, 67.5% of haplotypes had a match among the probes constructed from the other 871 people in the dataset. With a 3-cM threshold, 0.2% of haplotypes had a match.