

## A scalable approach for genome-wide inference of ancestral recombination graphs

Árni Freyr Gunnarsson<sup>1,2,3</sup>, Jiazheng Zhu<sup>2</sup>, Brian C. Zhang<sup>2</sup>, Zoi Tsangalidou<sup>2</sup>,  
Alex Allmont<sup>4</sup>, Pier Francesco Palamara<sup>1,2</sup>

<sup>1</sup> Centre for Human Genetics, University of Oxford, UK

<sup>2</sup> Department of Statistics, University of Oxford, UK

<sup>3</sup> deCODE genetics/Amgen, Reykjavík, Iceland

<sup>4</sup> Doctoral Training Centre, University of Oxford, Oxford, UK

Correspondence: palamara@stats.ox.ac.uk

### Abstract

1 The ancestral recombination graph (ARG) is a graph-like structure that encodes a detailed  
2 genealogical history of a set of individuals along the genome. ARGs that are accurately  
3 reconstructed from genomic data have several downstream applications, but inference from  
4 data sets comprising millions of samples and variants remains computationally challenging.  
5 We introduce Threads, a threading-based method that significantly reduces the computational  
6 costs of ARG inference while retaining high accuracy. We apply Threads to infer the ARG of  
7 487,409 genomes from the UK Biobank using ~10 million high-quality imputed variants,  
8 reconstructing a detailed genealogical history of the samples while compressing the input  
9 genotype data. Additionally, we develop ARG-based imputation strategies that increase  
10 genotype imputation accuracy for ultra-rare variants ( $MAC \leq 10$ ) from UK Biobank exome  
11 sequencing data by 5-10%. We leverage ARGs inferred by Threads to detect associations with  
12 52 quantitative traits in non-European UK Biobank samples, identifying 22.5% more signals  
13 than ARG-Needle. These analyses underscore the value of using computationally efficient  
14 genealogical modeling to improve and complement genotype imputation in large-scale  
15 genomic studies.

### 16 Introduction

17 The ancestral recombination graph (ARG) is a graph in which nodes represent the genomes of  
18 a set of individuals or their ancestors, and edges represent genealogical connections between  
19 them. At any genomic position, the genealogical connections encoded in the ARG form a single

20 genealogical tree, which changes along the genome due to recombination events that occurred  
21 in the transmission of genetic material from ancestors to descendants. The ARG efficiently  
22 integrates all these marginal trees into a single object and has been widely studied<sup>1-3</sup>. ARGs  
23 have also been leveraged in a wide range of genomic analyses, including generating synthetic  
24 data<sup>4-7</sup>, analyzing demographic history<sup>8-11</sup> and natural selection<sup>8,11</sup>, detecting complex trait  
25 associations<sup>12-14</sup>, and facilitating polygenic prediction<sup>15</sup>.

26 Inference of ARGs from genotype data is computationally challenging due to the vast search  
27 space of graph topologies and coalescence times that could give rise to the observed  
28 genotypes<sup>16</sup>. For this reason, most methods rely on a combination of probabilistic inference  
29 and computational heuristics<sup>8,10-13,17-20</sup>. Recent approaches have improved the inference and  
30 analysis of genome-wide ARGs from sparse genotyping array data in large biobank data sets<sup>13</sup>.  
31 However, current methods struggle to scale to biobank-sized collections comprising millions  
32 of samples and genomic variants. Furthermore, recent work has shown that genealogical  
33 modeling may be used to improve genotype imputation accuracy<sup>13,21</sup>, but these potential  
34 benefits have not yet been attained in modern genomic data sets.

35 We introduce a scalable method for inferring ARGs, called Threads, which can be applied to  
36 biobank-scale datasets of genotyped, imputed, or sequenced individuals. We use extensive  
37 simulations to show that Threads requires significantly less computation and memory  
38 compared to other ARG inference methods and remains accurate despite relying on several  
39 modeling simplifications. We apply Threads to infer a genome-wide genealogy for 487,409  
40 genomes from the UK Biobank using ~10 million high-quality imputed variants. The resulting  
41 ARG integrates both genotypic and genealogical information about the analyzed samples, and  
42 an encoding derived from threading operations allows for more compact storage of input  
43 genotype data compared to commonly used genotype formats. Finally, we develop strategies  
44 that use genealogies inferred using Threads to increase genotype imputation accuracy of ultra-  
45 rare variants and use the inferred ARG to complement genotype imputation in association  
46 analyses of non-European UK Biobank samples.

## 47 **Results**

### 48 **Overview of Threads**

49 Threads infers ARGs through a process called *threading*<sup>11,13</sup>, whereby new haploid genomes  
50 are sequentially grafted onto a partial ARG by computing a set of “threading instructions”<sup>13</sup>.

51 At each position along the genome, and for each sequence, these indicate a closest genetic  
52 relative (or cousin) from among samples already in the ARG, as well as a coalescence time.  
53 Once inferred for the whole sample, these threading instructions uniquely identify an ARG<sup>13</sup>.  
54 The threading instructions output by Threads also optionally specify whether the sequence and  
55 the closest genetic relative share the same allele. When allele sharing information is included,  
56 threading instructions are also sufficient to recover the input genotypes without assembling the  
57 ARG, as shown in Supplementary Figure 1.

58 Threads iteratively builds an ARG by considering all individuals in a fixed order and inferring  
59 the threading instruction for each of these individuals, as depicted in Figure 1. For each target  
60 individual, the threading instruction is computed with respect to the samples that have been  
61 added to the ARG in a previous iteration. To achieve high scalability, Threads breaks down the  
62 inference of threading instructions into three steps, which may be followed by an additional  
63 step that uses these instructions to assemble the ARG. First, Threads uses a haplotype matching  
64 algorithm based on the positional Burrows-Wheeler transform<sup>22</sup> (PBWT) to select, for each  
65 haplotype, a set of candidate matches with an index lower than that of the target sample. Next,  
66 once all such candidate sets have been built, one among the set of most closely related  
67 haplotypes is selected at each site from among the candidate haplotypes by running the Li-  
68 Stephens algorithm<sup>23</sup> in parallel across multiple samples. Finally, coalescence times are  
69 inferred using a likelihood-based approach based on the length of matching segments, the  
70 number of mismatching alleles, and the demographic history of the samples. A more detailed  
71 description of the Threads algorithm can be found in the Supplementary Note.

## 72 **Accuracy and scalability of Threads in simulations**

73 We performed extensive simulations to test the computational scalability and accuracy of  
74 Threads. We included three other ARG inference methods in these benchmarks: Relate<sup>8</sup>,  
75 tsinfer<sup>19</sup> combined with tsdate<sup>10</sup> (tsinfer+tsdate), and ARG-Needle<sup>13</sup>, which is primarily  
76 optimized for inference from genotyping arrays and was only evaluated in that setting. We  
77 simulated sequencing and genotyping array data for up to 16,000 diploid genomes over a 15  
78 megabase (Mb) region. The results of these analyses are summarized in Figure 2; additional  
79 results may be found in Supplementary Figures 2-10.

80 We first measured runtime, memory usage, and disk space used to store results in all simulation  
81 settings. All methods were provided with 8 computational cores for parallel computation. At

82 its slowest, Threads ran ~3 times faster than tsinfer+tsdate, the second fastest method. In many  
83 settings, Threads achieved speed-ups over other methods of an order of magnitude or more,  
84 consuming less memory and disk space (Figure 2, Supplementary Figures 2-4). We observed  
85 similar improvements when only one CPU was made available to all methods (Supplementary  
86 Figure 5).

87 We used several metrics to evaluate the accuracy of Threads and the other tested approaches.  
88 Due to computational costs, only three of these metrics were applicable to these large sample  
89 sizes, namely the Robinson-Foulds metric<sup>24</sup> (RF), a tree-based measure of topological  
90 accuracy, the tree-total variation<sup>13</sup> (TV), which probabilistically measures accuracy of both  
91 topology and branch lengths, and the max-r<sup>2</sup> score, a stochastic accuracy score we devised to  
92 capture the extent to which variants implied by edges in the ARG tag other underlying genomic  
93 variants (see Methods). Because the RF metric is affected by the presence of polytomies (i.e.,  
94 ARG nodes with more than two descendants) these were randomly resolved when present.

95 We performed several other secondary analyses using four additional metrics. We assessed the  
96 number of mutations that are needed to heuristically map a held-out genomic variant to an  
97 inferred marginal tree (mapping score, MS; see Methods). For completeness, we also employed  
98 three other metrics that measure the Euclidean distance between vectors summarizing marginal  
99 trees in various ways. These included the root-mean-square error<sup>13</sup> (RMSE), which evaluates  
100 the accuracy of inferring pairwise genealogical distances, the topological Kendall-Colijn  
101 distance<sup>25</sup> (KC), which measures distances to the root from internal nodes and was evaluated  
102 using randomly resolved polytomies, and the split-size vector metric<sup>26</sup> (SV), which assesses  
103 clade size similarity. Due to their higher computational requirements, these vector-based  
104 metrics were only tested in smaller experiments. We note that the KC and the SV metric have  
105 been observed to be particularly sensitive to features such as tree shape and balance<sup>13,26</sup>, which  
106 may influence their interpretability in downstream analyses.

107 For sequencing data, Threads outperformed other methods under both the TV and max-r2  
108 scores, while Relate achieved the highest accuracy under the RF and MS metrics (Figure 2,  
109 Supplementary Figures 6, 10). In simulated genotyping array data, either Threads or ARG-  
110 Needle was the most accurate or tied for highest accuracy across all simulation settings (Figure  
111 2, Supplementary Figures 8, 10). Threads proved robust to genotyping errors (Supplementary  
112 Figure 7) and to an artificially low mutation rate (Supplementary Figure 9). In smaller

113 experiments involving vector-based metrics, Relate outperformed Threads and tsinfer+tsdate  
114 under RMSE and SV, whereas tsinfer+tsdate achieved the highest accuracy under KC.

115 Genealogical encodings of genotype data have been shown to enable compressing genotype  
116 information in simulated data<sup>7</sup>, although real data sets yielded lower compression rates<sup>19</sup>. The  
117 threading instructions output by Threads provide an alternative efficient encoding of both  
118 genealogical and genotypic data. In addition to uniquely determining the inferred ARG<sup>13</sup>, these  
119 instructions can be used to reconstruct the genotypic data used to infer it, as described in  
120 Supplementary Figure 1. We assessed disk space requirements for storing genotype data using  
121 compressed threading instructions from Threads, PLINK2<sup>27</sup> *pgen* format, and compressed  
122 ARGs from tsinfer+tsdate in *tszip* format. We found Threads to provide efficient storage  
123 compared to these formats (Supplementary Figure 11), with variation across simulation settings  
124 and dataset size. Compressed threading instructions could be used to efficiently recover the  
125 input genotype data directly, without the need to convert into other ARG formats  
126 (Supplementary Figures 1, 12, 13), although working with an assembled ARG often led to  
127 faster data retrieval, depending on factors such as the genotyping error rate.

128 Taken together, these results demonstrate that Threads can infer large-scale genealogies using  
129 fewer computational resources than other ARG inference methods, while remaining highly  
130 accurate and producing a compact encoding of both genealogical and genotype data.

### 131 **Genome-wide genealogies for the UK Biobank and 1000 Genomes Project data sets**

132 We applied Threads to infer genome-wide ARGs using genotyped, imputed, and sequenced  
133 genomic variants for up to 487,409 individuals from the UK Biobank (UKB) and for 2,261  
134 individuals from the 1000 Genomes Project (1KGP). Consistent with simulations, we observed  
135 that the compressed threading instructions output by Threads efficiently stored both  
136 genealogical and genotypic information for both data sets (Table 1). Threads' output was  
137 particularly space-efficient for an ARG inferred from 9,992,478 imputed variants for the UKB  
138 data set, requiring only 7.1% of the space used to store the input *pgen* file. For an ARG inferred  
139 in a subset of 711,755 SNP array variants and 337,464 unrelated white British samples,  
140 threading instructions required 38.9% of the space required by the *pgen* input file. When  
141 applied to the 1KGP data set using 1,227,802 sequenced variants from chromosome 20,  
142 Threads produced threading instructions that required 30% of the disk space needed for the  
143 input *pgen* file. This also corresponded to 30.4% of the space to used store an ARG inferred

144 using *tsinfer*+*tsdate* and compressed using *tszip*, which took 10.8× more time to compute. By  
145 converting these threading instructions into ARGs and analyzing average coalescence times  
146 and genealogical relationships between different geographic groups, we found that these  
147 inferred genealogies effectively retained fine-grained ancestry information, recovering patterns  
148 of population structure across the UK (Supplementary Figure 14, Methods).

#### 149 **ARG-based ultra-rare variant imputation**

150 Genealogical relationships between samples can be utilized for genotype imputation<sup>28-31</sup>, where  
151 genomic variants not directly observed in a target individual are inferred using close genetic  
152 relatives found in a sequenced reference panel. Current genotype imputation strategies build  
153 on the Li-and-Stephens<sup>23</sup> approach, which is highly effective but relies on approximate  
154 genealogical modeling. In particular, this approach does not allow modeling scenarios where  
155 the age of variants being imputed is less than the time to the most recent common ancestor  
156 between reference and target individuals (Figure 3a), which may lead to a systematic  
157 overestimation of rare-variant genotype dosages<sup>13,21</sup>. We verified this phenomenon using  
158 simulations, where we observed inflated dosages for variants with a minor allele count (MAC)  
159 of up to about 10, approximately independent of panel size (Supplementary Figure 15).

160 We developed an ARG-based imputation strategy that improves accuracy for ultra-rare variants  
161 when used with ARGs inferred using Threads. Briefly, we first inferred the ARG of a sequenced  
162 reference panel and assigned one or more associated edges in the ARG to variants in the  
163 reference panel. We then inferred probabilistic threading instructions for the target genome  
164 relative to the reference panel, estimating the probability that the target genome inherited each  
165 variant in the panel (Methods). This approach allows for modelling the possibility that the  
166 ancestor shared by the target and reference genomes is older than the age of the mutation being  
167 imputed. Since we expected imputation approaches based on the Li-Stephens algorithm to only  
168 result in inflated dosages for the rarest variants (Figure 3, Supplementary Figure 15), we  
169 adopted the Li-Stephens forward-backward algorithm<sup>28,31</sup> for variants with a minor allele  
170 frequency exceeding 0.1% or minor allele count exceeding 20.

171 We assessed the accuracy of this threading-based genotype imputation approach through  
172 simulated and real datasets. We first simulated reference panels of up to 30,000 diploids and  
173 measured aggregate  $r^2$  for variants categorized by frequency in the panel. Threads showed a  
174 10-15% improvement in imputation  $r^2$  for singleton variants compared to IMPUTE5 and

175 Beagle 5.4. As allele counts increased, differences in accuracy decreased, with all methods  
176 reaching similar levels for variants of MAC=10 in the reference (Figure 3, Supplementary  
177 Figure 16). We next applied this method to an inferred ARG for up to 199,000 UKB exome-  
178 sequenced samples, obtaining accuracy improvements of ~5-10% for variants with MAC  
179 between 1 and 10 (Figure 3). Finally, we applied this approach to African and European  
180 ancestry samples from the 1KGP data set, observing gains starting at 5-6% for the rarest  
181 variants, with the accuracy of all methods evening out at around MAC=10 (Supplementary  
182 Figure 17).

### 183 **ARG-based association testing**

184 Recent work has shown that genealogy-wide association analyses, where genomic variants are  
185 inferred from a reconstructed ARG and tested for association against a heritable trait, can  
186 complement genotype imputation in the study of genomic variation that is not well represented  
187 in sequenced reference panels<sup>13</sup>. Inferred ARGs can also be used within a linear mixed model  
188 framework to estimate heritability<sup>13</sup> by constructing an ARG-based genetic relatedness matrix  
189 (ARG-GRM) that may better capture unobserved genomic variation. When these ARG-GRMs  
190 are built for a specific genomic region, such as a gene, this approach can be used to perform  
191 ARG-based variance component association testing<sup>14,32</sup>.

192 We performed ARG-based association testing within the UK Biobank dataset, comparing the  
193 use of ARGs inferred using either Threads or ARG-Needle (see Methods). We focused on  
194 individuals of non-European genetic ancestry, who are underrepresented in imputation  
195 reference panels<sup>33</sup>, analyzing 8,235 unrelated individuals of Central/South Asian ancestry  
196 (CSA) and 6,253 unrelated individuals of African ancestry (AFR), as defined in the Pan-UKB  
197 project<sup>34</sup>. We tested for association between 21,378 protein-coding genes and non-coding RNA  
198 regions (Supplementary Table 1) with 52 blood cell indices and blood biochemistry marker  
199 levels (Supplementary Table 2), comparing several association strategies (see Methods). We  
200 performed variance component testing using an ARG inferred from ~10 million HRC-imputed  
201 variants (Threads-HRC). Since ARG-Needle cannot be easily applied to data sets of this scale,  
202 we tested an ARG previously inferred using 711,754 SNP array markers<sup>13</sup> (ARG-Needle-SNP).  
203 We also included tests performed using an ARG inferred using Threads for the same subset of  
204 markers (Threads-SNP). Finally, we compared these ARG-based approaches to standard  
205 association testing based on genotype imputation. To this end, we applied the same variance  
206 component association test that we used in ARG-RHE to test HRC-imputed variants (HRC-

207 RHE), and ran standard mixed-model association testing of individual imputed variants using  
208 Regenie<sup>35</sup>.

209 In these analyses, ARG-based association effectively complemented imputation-based  
210 approaches (Figure 4; Supplementary Tables 3, 4; Supplementary Figure 18). Variance  
211 component association testing performed using ARGs inferred with Threads detected more  
212 gene-trait associations than when using ARGs inferred with ARG-Needle ( $N_{\text{Threads-HRC}} = 212$ ,  
213  $N_{\text{Threads-SNP}} = 182$ ,  $N_{\text{ARG-Needle-SNP}} = 173$ , Figure 4), and more than when we applied the same  
214 variance component association test to imputed genotype data alone ( $N_{\text{HRC-RHE}} = 155$ ). In  
215 addition, the signals detected using the ARG were complementary to those detected using  
216 imputation (Figure 4b,c, combined with HRC-RHE:  $N_{\text{Threads-HRC}} = 211$ ,  $N_{\text{Threads-SNP}} = 208$ ,  $N_{\text{ARG-Needle-SNP}} = 198$ , see Methods). This complementarity was also observed when comparing these  
217 associations to those detected using single-variant testing from a mixed-model analysis, with  
218 Threads again detecting a larger fraction of associations shared with imputation-based testing  
219 compared to ARG-Needle.  
220

221 We verified these associations by checking for overlap with larger, more powered association  
222 studies, in which we found most signals detected in the CSA and AFR subgroups to be reported.  
223 Previously established signals which we detected using the ARG but not using HRC-imputed  
224 data alone in our analyses included, in CSA, associations between the TFR2, ACTL6B, and  
225 HFE genes and Mean Corpuscular Hemoglobin<sup>36-41</sup>, as well as between the SLC12A3 gene and  
226 HDL Cholesterol<sup>39,42-44</sup>. We also detected associations that did not overlap with those reported  
227 in the Open Targets database (see URLs), particularly for the AFR subgroup, for which fewer  
228 ancestry-matched large-scale studies exist. Among these, we detected associations between  
229 reticulocyte measures and EMP2, which has previously been associated to red blood cell counts  
230 and other blood traits<sup>36,39,40,45</sup>, and between MMP26 and neutrophil percentage, previously  
231 associated to lymphocyte counts, monocyte counts, and other blood traits<sup>36,37,45-47</sup>.

232 Overall, ARG-based association testing effectively complemented imputation-based  
233 approaches in these analyses, with Threads outperforming ARG-Needle in both the number of  
234 detected associations and their overlap with imputation-based signals. Although Threads is  
235 more computationally efficient, its reliance on model simplifications may reduce the accuracy  
236 of ARGs inferred from SNP data alone at deeper time scales. Consistent with this, when we  
237 used the ARG-MLMA approach to perform genealogy-wide mixed-model association testing  
238 of individual ARG-derived variants<sup>13</sup>, ARGs inferred by Threads using SNP data yielded fewer



239 approximately independent associations than those inferred by ARG-Needle (see Methods,  
240 Supplementary Figure 18). Given that single-variant testing is less powered to detect  
241 associations with rare variants compared to variance-components-based testing<sup>48</sup>, this  
242 discrepancy likely reflects ARG-Needle's higher accuracy in inferring common variants, which  
243 tend to originate from deeper time scales. However, Threads' scalability allows it to be applied  
244 to denser collections of imputed variants, potentially improving the inferred ARG, particularly  
245 at deeper time scales. We note that genotype imputation may also introduce errors, which may  
246 affect the detection of long shared haplotypes. The effects of imputation on the accuracy of  
247 inferred ARGs may therefore vary across analyzed groups, reflecting variation in imputation  
248 accuracy for variants of different frequencies and ages.

## 249 **Discussion**

250 We developed Threads, a scalable algorithm for the inference of ancestral recombination  
251 graphs. Through extensive simulation and benchmarking, we verified that Threads performs  
252 ARG inference using fewer computational resources than other available methods, while  
253 retaining high accuracy. We applied Threads to infer genome-wide ARGs for 2,261 samples  
254 from the 1,000 Genomes Project, using ~57 million variants, and for 487,409 samples from the  
255 UK Biobank, using ~10 million high-quality imputed variants. The inferred ARGs, encoded  
256 using compressed threading instructions, compactly store both genotype and genealogical data  
257 while requiring significantly less disk space compared to standard genotype formats. We  
258 developed strategies to use these inferred ARGs to perform genotype imputation, observing  
259 accuracy improvements for imputed ultra-rare variants ( $MAC \leq 10$ ) in both simulated and UK  
260 Biobank exome sequencing data. Finally, we used the ARGs inferred in non-European samples  
261 to detect associations in 52 complex traits, where Threads increased the number of associated  
262 loci compared to ARG-Needle and detected signals that complement those obtained using  
263 standard imputation-based strategies.

264 Our work highlights connections and potential synergies between genealogical inference and  
265 genotype imputation. The initial step of the Threads algorithm uses the PBWT data structure  
266 to rapidly identify genetic relatives who share long-range haplotypes, a strategy also used in  
267 genotype imputation algorithms<sup>28</sup>. Additionally, to compute the threading instructions used to  
268 reconstruct the ARG, Threads estimates coalescence times between samples. Our analyses  
269 demonstrate that this step can also be leveraged in genotype imputation to improve the accuracy  
270 of the rarest imputed variants. Similar benefits are likely to be observed in the closely related

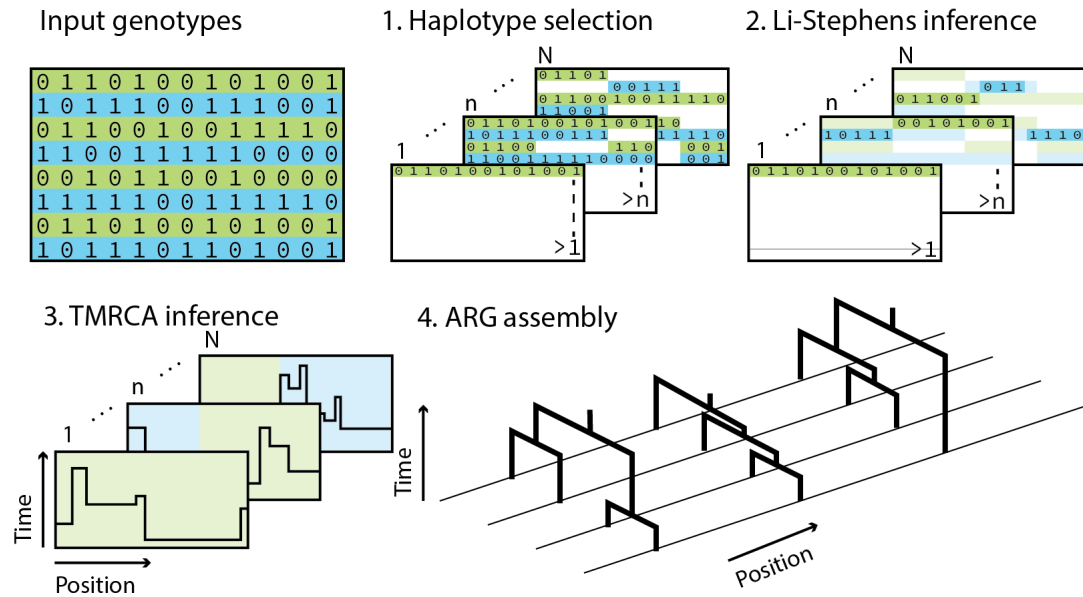
271 problem of haplotype phasing, where haplotype-based modeling of coalescence times has been  
272 shown to improve performance<sup>49</sup>. In a preliminary analysis, we used a simple algorithm that  
273 leverages inferred ARGs to phase ultra-rare variants. This approach was as accurate as  
274 SHAPEIT5 when applied to 1KGP data (Supplementary Figure 19b) but less accurate in other  
275 simulated scenarios (Supplementary Figure 19c,d), suggesting the need for further  
276 methodological development.

277 In addition to improving the imputation accuracy of rare variants, using an ARG to store the  
278 reference panel could facilitate the sharing of these panels. As recently suggested<sup>50</sup>, truncating  
279 an inferred ARG to remove genealogical connections within the most recent generations may  
280 preserve sufficient information to effectively perform analyses such as phasing and imputation,  
281 while offering some protection for the privacy of the individuals in the panel. This may  
282 facilitate data sharing, and the threading of new individuals onto these references could provide  
283 an effective algorithmic strategy towards building a more comprehensive global genealogical  
284 resource.

285 We identify several limitations and potential areas for future improvements. First, Threads  
286 estimates the age of genomic regions using a Li-Stephens algorithm, assuming these regions  
287 correspond to pairwise identical-by-descent (IBD) segments. While this approximation has  
288 proven reasonably accurate, particularly in large samples (see Supplementary Note), additional  
289 modeling may lead to improved accuracy for the ages inferred during the initial iterations of  
290 the algorithm, which typically influence the length of ARG edges in the deeper past.  
291 Furthermore, threading instructions inferred for the subset of samples that are initially threaded  
292 into an ARG can be easily replaced with another set of instructions for the same individuals.  
293 An inferred ARG may therefore be improved at a later stage by substituting these initial  
294 instructions with those derived from an ARG constructed using algorithms that are slower but  
295 more accurate in small data sets. Second, like other threading approaches, Threads  
296 incrementally adds new individuals to an existing ARG, resulting in slight variations in the  
297 ARG depending on the order of the individuals, which may be leveraged to obtain simple  
298 estimates of uncertainty<sup>13</sup>. The use of subtree pruning and regrafting operations<sup>11</sup> may lead to  
299 improved uncertainty quantification. Third, although the threading instructions inferred by  
300 Threads allow for efficient compression of genotypes, many ARG-based operations can be  
301 inefficient when performed directly on these instructions. In this case, it may be beneficial to  
302 convert the threading instructions into other formats or assemble them into ARG data

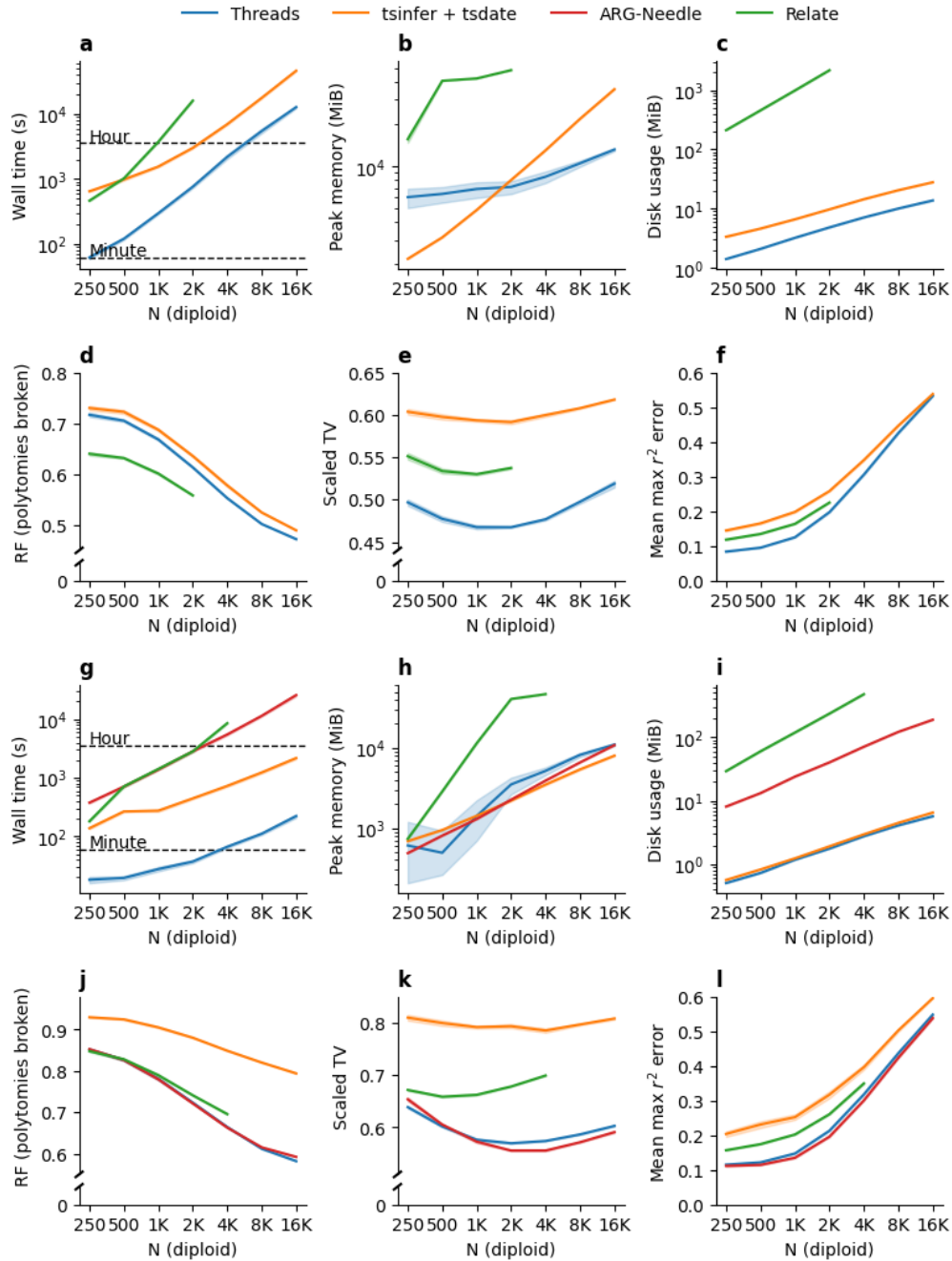
303 structures, which requires additional computational resources. Fourth, although we have shown  
304 that modeling of coalescence times can improve the imputation of ultra-rare variants, the  
305 current ARG-based imputation algorithm within the Threads package is at least an order of  
306 magnitude slower than standard, optimized algorithms. This underscores the need for further  
307 computational optimization of this approach. Fifth, Threads can more easily scale to analyses  
308 comprising large collections of imputed variants, but the quality of the inferred ARGs will  
309 depend on imputation accuracy, which in turn depends on several population- and sample-  
310 specific features. Despite these limitations and areas of future improvement, Threads provides  
311 a valuable new tool for the inference and analysis of genome-wide genealogies at biobank  
312 scales.

313



314

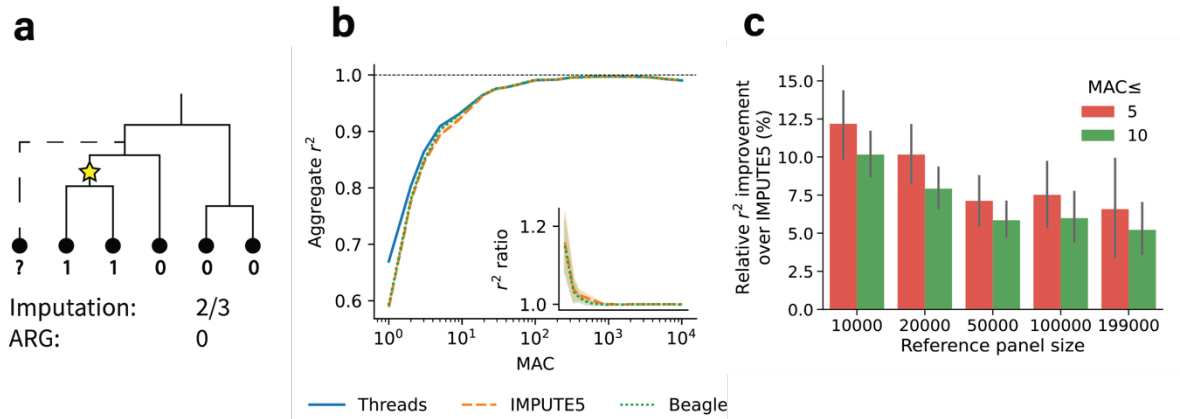
315 **Figure 1. Overview of the Threads algorithm.** Inference of threading instructions is  
316 performed in three steps. First, assuming a fixed order of input haplotypes, we select a subset  
317 of candidate genealogical closest cousins for each sample in each 0.5 centimorgan-sized  
318 window from among the samples of lower index. Second, we use the Li-Stephens algorithm to  
319 select a threading target (or genealogical closest cousin) from among the localized subsets.  
320 Third, we estimate the age of each segment inferred by the Li-Stephens algorithm. Finally, the  
321 inferred threading instructions may be used to assemble an ARG.



322

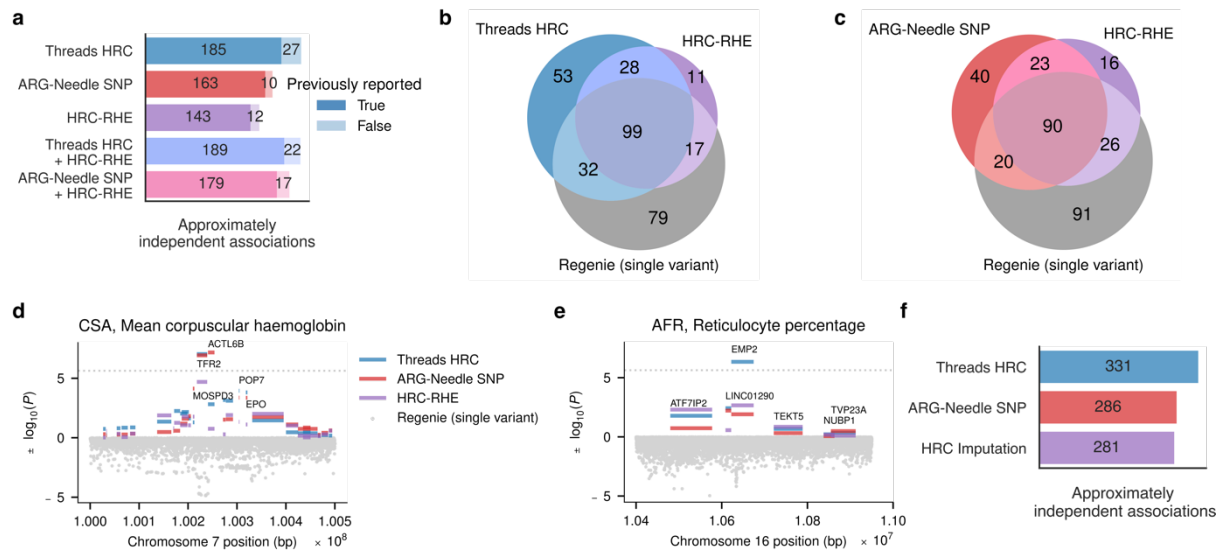
323 **Figure 2. Computational cost and inference accuracy of Threads, tsinfer+tsdate, ARG-**  
 324 **Needle, and Relate in simulated data.** We simulated 15 Mb regions under a European  
 325 demographic model with a constant recombination rate for sample sizes up to 16,000 diploid  
 326 individuals. Evaluation of Relate was truncated at 2,000 samples for sequencing simulations  
 327 and 4,000 samples for genotyping array simulations due to computational constraints. Shaded  
 328 regions show bootstrap 95% confidence intervals based on 15 random seeds. **a-c.**  
 329 Computational efficiency measured using runtime, memory consumption, and disk usage to  
 330 store results, excluding temporary files written by Relate. **d-f.** The Robinson-Foulds (RF)

331 metric, scaled total variation (TV), and mean  $\max\text{-}r^2$  error compared with the ground truth  
332 ARG. Polytomies were randomly broken when evaluating the RF metric (see Methods). **g-i.**  
333 Runtime, memory consumption, and disk usage, excluding temporary files written by Relate,  
334 for the same simulated regions with markers subsampled to genotyping array density. **j-l.** RF,  
335 TV, and mean  $\max\text{-}r^2$  error of ARGs inferred from simulated genotyping array data.



336

337 **Figure 3. ARG-based genotype imputation.** **a.** Toy example of how genotype dosages may  
338 be inflated even when the set of genealogical closest cousins is correctly inferred. Here, the  
339 target sample coalesces above the mutation (star symbol) carried by two of its closest cousins.  
340 **b.** Imputation quality for three methods using a simulated reference panel of 10,000  
341 individuals. The inset shows the ratio between Threads and other methods across the same  
342 allele count spectrum. Shaded regions represent 95% confidence intervals across 10 random  
343 seeds. **c.** Relative improvement (%) in imputation accuracy of variants with MAC up to 5 or  
344 10. We compare the results of applying Threads and IMPUTE5 to UK Biobank exome  
345 sequencing data. Error bars represent 95% confidence intervals computed using 260 regions of  
346 10 cM each.



347

348 **Figure 4. ARG-based association testing. a.** Total numbers of LD blocks containing genome-  
 349 wide significant gene-trait associations detected in the CSA and AFR subgroups using different  
 350 methods. Each bar is annotated with the subsets of LD blocks containing gene-trait associations  
 351 that were or were not previously reported on the Open Targets platform (see URLs). **b, c.**  
 352 Overlap in LD blocks containing gene-trait associations detected using HRC-RHE, Regenie,  
 353 and ARG-RHE applied to ARGs inferred using either Threads (**b**) or ARG-Needle (**c**).  
 354 Combined counts from the CSA and AFR subgroups; group-specific counts are reported in  
 355 Supplementary Figure 18. **d, e.** Manhattan plots for genomic regions containing the TFR2 and  
 356 ACTL6B genes (**d**) in the CSA subpopulation for mean corpuscular hemoglobin level and the  
 357 EMP2 gene (**e**) in the AFR subgroup for reticulocyte percentage. **f.** Approximately independent  
 358 single-variant associations detected through genealogy-wide testing of ARGs inferred using  
 359 Threads or ARG-Needle, as well as through genome-wide testing of HRC-imputed variants.



360  
361

(a)

| Algorithm/encoding | 1,000 Genomes Project<br>N=2,261; M=1,227,802 (Chr 20) |                     |
|--------------------|--------------------------------------------------------|---------------------|
|                    | Size (Mb)                                              | Inference time (h)* |
| Threads (.threads) | 46.7                                                   | 1.24                |
| tsinfer (.tsz)     | 153.7                                                  | 13.4                |
| vcf.gz             | 155.5                                                  | -                   |
| pgen               | 78.6                                                   | -                   |

362  
363

(b)

| Algorithm/encoding | HRC-imputed<br>N=487,409; M=9,992,478 |                      | SNP array<br>N=337,464; M=711,754 |                     |
|--------------------|---------------------------------------|----------------------|-----------------------------------|---------------------|
|                    | Size (Gb)                             | Inference time (h)** | Size (Gb)                         | Inference time (h)* |
| Threads (.threads) | 40.3                                  | 4,449.1              | 12.1                              | 192.2               |
| pgen               | 563.5                                 | -                    | 31.1                              | -                   |

364

365 \* 8 cpu threads

366 \*\* 16 cpu threads

367

368 **Table 1. Computational resources for the 1KG and UKB data sets.** We report the  
369 computing time and disk usage required to infer and store genealogical encodings, which  
370 include genealogical and genomic data, as well as the disk usage for storing the input  
371 genomic data. We present results for (a) Threads (.threads compressed threading instructions),  
372 tsinfer (.tsz ARG format), vcf.gz, and pgen formats for 1,227,802 variants from Chromosome  
373 20 for 2,261 individuals from the 1,000 Genomes Project; (b) Threads (.threads compressed  
374 threading instructions) and pgen formats for the UK Biobank data set, for 487,409 samples  
375 and 9,992,478 HRC-imputed genome-wide variants, as well as for a subset of 337,464  
376 unrelated white British samples and 711,754 SNP genome-wide array variants.

## 377 **Methods**

378 For a detailed description of the Threads algorithm, please refer to the Supplementary Note.

## 379 **Simulations**

380 To evaluate the accuracy of ARG inference methods, we simulated ARGs using msprime<sup>51</sup>  
381 under two demographic models, a model inferred using SMC++<sup>52,53</sup> for Northern Europeans  
382 from Utah (CEU) from the 1000 Genomes Project<sup>54</sup> and a constant demographic model of  
383 effective population size  $N_e = 10,000$  diploid individuals. All simulations used a fixed  
384 recombination rate of  $1.3 \times 10^{-8}$  per site per generation, approximately matching the average  
385 genome-wide recombination rate. We simulated mutations both using a realistic rate of  
386  $1.4 \times 10^{-8}$ , compatible with recent estimates<sup>55</sup>, and an artificially low mutation rate of  
387  $1.4 \times 10^{-9}$ . In each case we simulated sequences of 15 Mb in length with sample sizes ranging  
388 from 250 to 16,000 diploid individuals. To avoid boundary effects, all accuracy metrics were  
389 evaluated only on the central 5 Mb region. Each experiment was repeated for 15 different  
390 random seeds. Genotyping array data sets were simulated by subsampling the observed  
391 genotypes to match in frequency with genotyping array data from the UK Biobank. In  
392 experiments involving genotype imputation, we simulated regions of 20 Mb in length and  
393 evaluated accuracy on the central 10 Mb for 10 random seeds, but otherwise kept simulation  
394 parameters unchanged<sup>54</sup>.

## 395 **ARG metrics and ARG benchmarks**

396 We evaluated seven tree- or ARG-based accuracy metrics for Threads, ARG-Needle<sup>13</sup> (v1.0.1),  
397 Relate<sup>8</sup> (v1.2.1), tsinfer<sup>19</sup> (v0.3.1) and tsdate<sup>10</sup> (v0.1.5). Evaluation of Relate was truncated at  
398 2,000 samples for sequencing experiments and at 4,000 samples for genotyping array  
399 experiments due to computational constraints. We evaluated the Robinson-Foulds metric<sup>24</sup>  
400 (RF), the ARG total variation<sup>13</sup> (TV) and the max- $r^2$  score for all sample sizes. In addition, we  
401 computed the topological Kendall-Colijn metric<sup>25</sup> (KC), the root-mean-square error metric<sup>13</sup>  
402 (RMSE), and the split-size vector metric<sup>26</sup> (SV) for sample sizes up to 4,000. The mapping  
403 score was evaluated only at sample size 2,000, the largest sequencing simulations where all  
404 methods were included.

405 The max- $r^2$  score is an ARG-based accuracy metric measuring the correlation between variants  
406 in an inferred ARG and underlying true variants, which reflects the extent to which an inferred

407 ARG is expected to tag underlying variation in a genealogy-wide association study. Given a  
408 ground-truth ARG A and an inferred ARG B, we simulate  $M_A$  and  $M_B$  mutations on A and B  
409 respectively, by uniformly distributing them along the ARG volume. Then, for each mutation  
410  $m_A$  on A, represented as a bit-set, we compute  $\max_{m_B \in B} r^2(m_A, m_B)$ , thus finding the mutation  
411 on B that is most closely correlated with  $m_A$  from among the  $M_B$  mutations on B. The max- $r^2$   
412 score is defined as the mean over all such correlations,

$$413 \quad \max r^2(A, B) = \frac{1}{M_A} \sum_{m_A \in A} \max_{m_B \in B} r^2(m_A, m_B).$$

414 For the experiments of Figure 2 and Supplementary Figures 6-9, we set  $M_A = M_B = 10,000$ ,  
415 so that the computational time required to evaluate this metric remained independent of sample  
416 size. Note, however, that this may cause the max- $r^2$  score to increase with sample size, due to  
417 a corresponding increase in ARG volume. An alternative approach, which would be less  
418 computationally efficient but would allow keeping the sampling rate constant over ARG edges,  
419 is to set the number of resampled mutations to be proportional to ARG volume.

420 We define the mapping score (MS) of an observed variant, given an inferred marginal tree, as  
421 the minimum number of edges in the tree where mutations need to be added to cover all and  
422 only the carriers of the variant, divided by the total number of carriers. We compute the  
423 minimum number of edges using the following iterative approach. First, we select any variant  
424 carrier for which an ancestral mutation has not yet been found. We then consider the sequence  
425 of edges connecting this node to the root and add a mutation to the last edge in this sequence  
426 that only subtends individuals who carry the variant. We consider these subtended carriers to  
427 be covered by the mutation and iterate until all remaining carriers have been covered. We then  
428 compute the MS as the ratio between the number of added mutations and the number of carriers.  
429 In simulations, when a ground-truth ARG is known, we compute the mapping score by  
430 randomly generating new mutations under the ground-truth ARG and aggregating their  
431 mapping score by ancestral allele frequency.

### 432 **ARG inference in the UK Biobank and 1000 Genomes data sets**

433 To infer genome-wide genealogies for the UK Biobank from markers imputed using the  
434 Haplotype Reference Consortium (HRC) reference panel<sup>30,56</sup>, we first extracted bi-allelic SNPs  
435 with INFO score  $\geq 0.95$  and minor allele frequency (MAF)  $\geq 0.0001$ . We rounded all dosages

436 within 0.1 of the nearest integer, otherwise setting genotypes to missing, and discarded all  
437 variants with missingness greater than 10%. All filtering was performed using PLINK2<sup>27</sup>. We  
438 divided the genome into chunks of 15 cM along with 1 cM of padding on each end and phased  
439 using SHAPEIT5-common<sup>49</sup>. We then joined pairs of chunks using SHAPEIT5-ligate, to create  
440 chunks of 30 cM, with 1 cM of padding. This procedure gave a total of 136 chunks, which were  
441 used to infer ARGs using Threads using parameters `--query_interval 0.02` and `--`  
442 `match_group_interval 1.0` instead of the default `0.01` and `0.5`, respectively. These  
443 parameters determine the density of queries for candidate matches in the haplotype matching  
444 step of the algorithm and may be raised to reduce memory usage; in our experiments Threads  
445 proved robust to the choice of matching algorithm parameters. To infer ARGs for samples of  
446 African and Central/South-Asian ancestry, we applied the same procedure, with a minor allele  
447 count (MAC) threshold of 1 and without any chunking, inferring ARGs for whole chromosome  
448 arms at a time using Threads with default parameters. For ARGs based on genotyping arrays,  
449 we followed the procedure described by Zhang et al.<sup>13</sup> and phased the arrays using Beagle 5.1<sup>57</sup>  
450 and divided the genome up into 166 chunks of equal size, inferring ARGs using Threads in  
451 genotyping array mode. To infer ARGs for 2,261 individuals from the 1,000 Genomes Project,  
452 we downloaded phased, curated data comprising 56,935,222 variants (see URLs) and applied  
453 Threads directly to each chromosome arm without any further chunking.

#### 454 **ARG-based analysis of population structure**

455 To extract regional structure (Supplementary Figure 14), we subsampled the ARG to 100  
456 samples from each of 122 UK postcodes, keeping only self-identified white British  
457 individuals<sup>56</sup>. For each sample, we then counted the occurrences of each postcode within the  
458 set of genealogical closest cousins genome-wide, querying marginal trees of the ARG at 10-  
459 kilobase intervals. We averaged these observations across individuals within each postcode.  
460 Regional heatmaps show the proportion of closest genetic relatives from each postcode for a  
461 single, focal postcode. To visualize regional structure<sup>19,58</sup>, we focused on the ARG of 49,354  
462 self-reported white British individuals<sup>56</sup> from five postcodes in north-east England (DH,  
463 Durham; DL, Darlington; NE, Newcastle; SR, Sunderland; TS, Middlesbrough) and evaluated  
464 the proportion of genome shared between individuals within 10 generations. We averaged these  
465 observations across individuals and divided by the number of observations to obtain an affinity  
466 matrix that was used for average-linkage hierarchical agglomerative clustering algorithm as  
467 implemented in scikit-learn<sup>59</sup> (v.1.3.0). We truncated the clustering operation at different

468 relatedness thresholds, ranging from  $\alpha = 0.1$  down to  $\alpha = 1 \times 10^{-5}$ , extracting the largest  
469 clusters for each  $\alpha$ .

#### 470 **Threading-based compression**

471 For experiments quantifying the disk space required by Threads' output, we measured the disk  
472 space required to store compressed threading instructions. These threading instructions consist  
473 of a partition of a genomic region into discrete segments, and attached to each such segment, a  
474 threading target, a coalescence time, and a list of sites on the segment where the sequence and  
475 the threading target are heterozygous. These quantities are serialized and then compressed  
476 using the HDF5 format<sup>60</sup>. In addition to uniquely determining the ARG<sup>13</sup>, these threading  
477 instructions can be used to recover the input genotype variants, as they define, at each site, a  
478 directed acyclic graph that can be used to extract genotypes, as illustrated in Supplementary  
479 Figure 1.

#### 480 **Threading-based imputation**

481 We evaluated the accuracy of threading-based genotype imputation in both simulations and  
482 real data. We simulated reference panels and genotyping array data sets of size ranging from  
483 100 to 30,000 samples for 10 random seeds using the parameters described above and used  
484 these to impute 10 held-out target samples. We estimated accuracy using the aggregate  $r^2$   
485 binned by minor allele count in the panel. In experiments involving the 1,000 Genomes Project  
486 data, we used 10 held-out samples of European ancestry (2 from each of 5 sub-populations:  
487 CEU, Northern Europeans from Utah; FIN, Finnish; GBR, British; IBS, Iberian; TSI Tuscans  
488 from Italy) and 7 held-out samples of African ancestry (1 from each of 7 sub-populations: ACG,  
489 African Caribbean in Barbados; ASW, African-American in Soth West USA; ESN, Esan in  
490 Nigeria; YRI, Yoruba, Nigeria; LWK, Luhya, Kenya; GWD, Gambian; MSL, Mende, Sierra  
491 Leone) from a set of 2,261 unrelated samples from the 1,000 Genomes Project<sup>54,61</sup>. In  
492 experiments using the UK Biobank exome sequencing data, we evaluated imputation accuracy  
493 on 100 held-out samples using a reference panel of up to 199,000 exome sequenced samples.

494 To perform threading-based imputation, we first inferred an ARG for the reference panel using  
495 Threads. We then mapped variants to edges in the ARG using a method previously described  
496 by Speidel et al.<sup>8</sup>. This approach can be applied in cases where the ancestral allele status is  
497 unknown and accounts for errors in the input genotypes and the inferred ARG by allowing for  
498 imperfect overlap between the leaves subtended by lineages assigned to mutations and their

499 observed carriers. This enabled us to estimate the age  $s$  of a mutation as the midpoint of its  
500 mapped lineage. As we only applied ARG-based imputation to rare variants, we heuristically  
501 assumed the major allele to be ancestral. For variants that did not map to the ARG, we  
502 performed standard imputation using the Li-Stephens model. Next, we inferred threading  
503 instructions with respect to the reference ARG for each sample being imputed. We modeled the  
504 age of a shared haplotype,  $t$ , as an Erlang-2 distribution with parameter  $\lambda$ , a function of segment  
505 length and population size (see Supplementary Note). We used this to compute the threading  
506 dosage:

$$507 \quad p(t < s) = 1 - e^{-\lambda s}(1 + \lambda s),$$

508 which estimates the probability that the lineage being threaded inherited the mutation. To  
509 account for uncertainty in the threading target, we used the PBWT to recover all sequences that  
510 are identical to the threading target along the segment being imputed. Finally, to account for  
511 uncertainty on the boundaries of the segment, and to further account for uncertainty in the  
512 threading target, we also ran a forward-backward algorithm to get a full posterior matrix  $P$ ,  
513 with  $P_{ij}$  denoting the probability that sample  $i$  is the genealogical closest cousin for the target  
514 at site  $j$ . To obtain the dosage for each site, we computed the expectation

$$515 \quad E[g_j] = P(g_j = 1) = \sum_{i=1}^n P(c_{ij} < s_{ij}) \cdot P_{ij} \cdot g_{ij},$$

516 where  $n$  is the number of sequences in the reference panel. Here,  $P(c_{ij} < s_{ij})$  denotes the  
517 probability that at site  $j$  the target sequence coalesces to reference sample  $i$  at time  $c_{ij}$  lower  
518 than the age  $s_{ij}$  of the mutation carried by sample  $i$  at site  $j$ , with  $s_{ij} = 0$  if it carries no  
519 mutation. The term  $g_{ij}$  denotes the observed genotype of sample  $i$  at site  $j$  in the panel.

## 520 **ARG-based phasing**

521 To phase rare variants using the ARG, we traversed a marginal coalescence tree upwards,  
522 starting from each of the two possible haplotypes carrying an unphased variant, until at least  
523 one haplotype of all heterozygous carriers and both haplotypes of all homozygotes were found.  
524 We defined the average tree-based distance to all such carriers as the *phasing distance* for the  
525 haplotype and placed the variant on the haplotype of lower phasing distance. More formally, if

526 for each unphased carrier  $c$ ,  $c_1$  and  $c_2$  are its two haplotypes, then the phasing distance can be  
527 written as

$$528 \quad d_{\text{phase}}(c'_a) = \sum_{c \text{ heterozygous}} \min(d(c'_a, c_1), d(c'_a, c_2)) + \sum_{c \text{ homozygous}} d(c'_a, c_1) + d(c'_a, c_2)$$

529 for  $a = 1, 2$ . If  $c'$  is the only carrier of the mutation, we set  $d_{\text{phase}}(c'_a)$  to be the negative of the  
530 age of its most recent common ancestor. The singleton case prioritizes the haplotype that is less  
531 related to the rest of the sample and more likely to carry a singleton mutation. A similar step  
532 based on identical-by-state tracts is performed by SHAPEIT5<sup>49</sup>. For all ARG-based analyses,  
533 we rely on common-variant scaffolds inferred using SHAPEIT5 for variants with MAF above  
534 0.1%.

### 535 **Association analyses**

536 We tested 52 quantitative blood traits including 27 blood cell indices and 25 blood biochemistry  
537 marker levels (see Supplementary Table 1). We applied standard preprocessing steps<sup>45</sup> to these  
538 phenotypes separately for the AFR and CSA subgroups. We first stratified samples based on  
539 sex and menopause status and applied a rank-inverse-normal transformation (RINT). We then  
540 regressed out covariates, which included alcohol use, smoking status, age, age squared, height,  
541 body mass index (BMI), assessment center, genotyping array, and 20 genetic principal  
542 components, and applied RINT a second time. We considered 21,378 protein-coding and non-  
543 coding RNA regions spanning more than 1,000 bp from the University of California, Santa  
544 Cruz (UCSC) genome browser Gene Nomenclature Committee (HGNC) table (see  
545 Supplementary Table 2 and URLs).

546 We performed gene-based, variance component association testing using the ARG-RHE  
547 algorithm<sup>32</sup> implemented in the arg-needle-lib library (v1.0.3). We modeled the vector of  
548 phenotypes  $y$  for  $n$  individuals as  $y \sim N(0, \sigma_g^2 K + \sigma_e^2 I)$ , where  $K$  is a GRM formed by variants  
549 within a region,  $I$  is an  $n \times n$  identity matrix, and  $\sigma_g^2 + \sigma_e^2 = 1$  are scalars reflecting genetic  
550 and environmental variance components. In the case of ARG-based association testing, the  
551 GRM  $K$  is a local Monte Carlo ARG-GRM from normalized resampled mutations within a  
552 region, built as in as in Zhang et al.<sup>13</sup> by sampling variants with  $\text{MAC} \geq 5$  from the ARG at a  
553 rate  $\mu = 10^{-6}$ . For HRC-RHE,  $K$  is the local GRM using normalized imputed genotypes in a  
554 gene region having  $\text{MAC} \geq 5$ , INFO score  $> 0.3$ , missingness  $< 0.1$ , and Hardy-Weinberg

555 equilibrium  $p > 10^{-15}$ . We used the ARG-RHE algorithm to perform variance component  
556 association testing, using GRMs built from ARG-derived (referred to as ARG-RHE) or imputed  
557 (referred to as HRC-RHE) variants within each of the considered gene regions (Supplementary  
558 Table 1). We applied a Bonferroni correction at a 5% family-wise error rate, obtaining a  
559 genome-wide significance threshold of  $0.05/21,378 \approx 2.3 \times 10^{-6}$  for each trait. We also report  
560 the results of combining the ARG-RHE and HRC-RHE tests on both ARG-derived and imputed  
561 variants, by comparing the smaller p-value for the gene given by the two tests with a genome-  
562 wide significance threshold of  $0.05/21,378/2 \approx 1.2 \times 10^{-6}$ .

563 As a baseline, we also performed single-variant association testing using Regenie<sup>35</sup>, testing  
564 HRC-imputed dosages on the same 52 preprocessed phenotypes with covariates regressed out.  
565 We applied the same variant filtering criteria as in the HRC-RHE analysis, which led to ~30  
566 million imputed variants, and adopted a genome-wide significance threshold obtained through  
567 resampling-based testing<sup>13,62</sup> (see Supplementary Table 5). To count associated regions for each  
568 approach, we assigned genome-wide significant signals into approximately independent LD  
569 blocks provided by LDetect<sup>63</sup>, which include 2,583 and 1,443 LD blocks for the AFR and CSA  
570 subgroups respectively.

571 We also performed genome-wide and genealogy-wide association analyses for the same  
572 preprocessed traits, following the methodology of Zhang et al.<sup>13</sup>. For GWAS, we kept variants  
573 with  $MAC \geq 5$ , missingness  $< 10\%$ , and INFO score  $> 0.3$  and then performed association  
574 using BOLT-LMM (v.2.3.4), without excluding related individuals. For genealogy-wide  
575 association, we inferred ARGs from SNP array and imputed data as described above. We then  
576 tested for association using the same parameters as Zhang et al.<sup>13</sup>, including a sampling rate of  
577  $\mu = 10^{-5}$  and calibration factors estimated by BOLT-LMM, while filtering out clades with  
578 derived allele count  $< 5$ . Genome-wide significance thresholds for this analysis, reported in  
579 Supplementary Table 5, were established using resampling-based testing<sup>13,62</sup>. To estimate the  
580 number of approximately independent associations, we first performed stringent clumping  
581 using PLINK<sup>27</sup> (v.1.90) with  $r^2=0.01$ . We then ran a conditional-joint analysis on the clumped  
582 SNPs using GCTA-COJO<sup>64,65</sup> with default parameters, setting the `--cojo-slct` flag with `-`  
583 `-cojo-p 1e-7`.



584 **Data availability statement:**

585 UK Biobank data can be accessed by approved researchers through  
586 <https://www.ukbiobank.ac.uk/>. The 1000 Genomes Project data set can be accessed through  
587 [http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data\\_collections/1000G\\_2504\\_high\\_coverage/work](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/1000G_2504_high_coverage/working/20220422_3202_phased_SNV_INDEL_SV/)  
588 [ing/20220422\\_3202\\_phased\\_SNV\\_INDEL\\_SV/](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/1000G_2504_high_coverage/working/20220422_3202_phased_SNV_INDEL_SV/). Variant annotations were obtained through  
589 the genome aggregation data base (<https://gnomad.broadinstitute.org/>), coding and non-coding  
590 RNA regions were obtained through the UCSC genome browser ([https://genome.ucsc.edu/cgi-](https://genome.ucsc.edu/cgi-bin/hgTables)  
591 [bin/hgTables](https://genome.ucsc.edu/cgi-bin/hgTables)), known variant associations were obtained through the Open Targets GWAS  
592 association database (<https://genetics.opentargets.org>), ancestry allocations for the UK  
593 Biobank were downloaded from the Pan-UKB project (<https://pan.ukbb.broadinstitute.org/>).

594

595 **Code availability statement:**

596 The Threads software is available at <https://palamaralab.github.io/software/threads>. External  
597 software can be downloaded from the following URLs: msprime (v.1.2.0  
598 <https://pypi.org/project/msprime/>), tsinfer (v.0.3.1, <https://pypi.org/project/tsinfer/>), tsdate  
599 (v.0.1.5, <https://pypi.org/project/tsdate/>), Relate (v.1.2.1, <https://myersgroup.github.io/relate/>),  
600 ARG-Needle (v1.0.2, <https://pypi.org/project/arg-needle/>), IMPUTE5 (v.1.1.5,  
601 <https://jmarchini.org/software/#impute-5>), Beagle 5.4 (v.22Jul22.46e,  
602 <https://faculty.washington.edu/browning/beagle/beagle.html>), SHAPEIT5 (v.5.1.1,  
603 <https://odelaneau.github.io/shapeit5/>), GCTA (v.1.94.1,  
604 <https://yanglab.westlake.edu.cn/software/gcta/>), PLINK2 (v1.90b6.26 and v2.00a3.7LM,  
605 <https://www.cog-genomics.org/plink/2.0/>).

606

607 **Competing interests:**

608 Á.F.G. is an employee of deCODE genetics/Amgen; B.C.Z. is an employee of Adaptive  
609 Biotechnologies.

610

611 **Acknowledgements:**

612 We thank R. Fournier for discussions and suggestions. This work was conducted using the UK  
613 Biobank resources (application no. 43206) and supported by the Clarendon Scholarship  
614 (Á.F.G., B.C.Z.); the Keble College de Breyne Clarendon Scholarship (Á.F.G.); Wellcome  
615 Trust ISSF grant no. 204826/Z/16/Z (P.F.P); Wellcome Trust grant no. 222336/Z/21/Z (Á.F.G.);  
616 ERC Starting Grant ARGPHENO no. 850869 (P.F.P, B.C.Z.); EPSRC grant EP/S023151/1

617 (Z.T.); EPSRC Centre for Doctoral Training in Health Data Science (EP/S02428X/1) (J.Z.).  
618 Computation used the Oxford Biomedical Research Computing (BMRC) facility, a joint  
619 development between the Centre for Human Genetics and the Big Data Institute supported by  
620 Health Data Research UK and the NIHR Oxford Biomedical Research Centre. Financial  
621 support was provided by the Wellcome Trust Core Award Grant Number 203141/Z/16/Z. The  
622 views expressed are those of the author(s) and not necessarily those of the NHS, the NIHR or  
623 the Department of Health.

624

625 **Contributions:**

626 Á.F.G. and P.F.P. designed the Threads algorithm. Á.F.G implemented algorithms and  
627 performed simulations and analyses of 1000 Genomes Project data. Á.F.G, J.Z., and Z.T.  
628 performed analyses of UK Biobank data. B.C.Z and A. A. provided software tools. Á.F.G. and  
629 P.F.P. wrote the manuscript.

630

631 **Correspondence:**

632 palamara@stats.ox.ac.uk

## 633 References

634

- 635 1. Hudson, R.R. Properties of a neutral allele model with intragenic recombination.  
636 *Theoretical Population Biology* **23**, 183-201 (1983).
- 637 2. Griffiths, R.C. & Marjoram, P. Ancestral Inference from Samples of DNA Sequences with  
638 Recombination. *Journal of Computational Biology* **3**, 479-502 (1996).
- 639 3. Griffiths, R.C. & Marjoram, P. An ancestral recombination graph. *Progress in population*  
640 *genetics and human evolution*, 257-270 (1997).
- 641 4. Hudson, R.R. Generating samples under a Wright-Fisher neutral model of genetic  
642 variation. *Bioinformatics* **18**, 337-8 (2002).
- 643 5. Shlyakhter, I., Sabeti, P.C. & Schaffner, S.F. Cosi2: an efficient simulator of exact and  
644 approximate coalescent with selection. *Bioinformatics* **30**, 3427-9 (2014).
- 645 6. Palamara, P.F. ARGON: fast, whole-genome simulation of the discrete time Wright-  
646 fisher process. *Bioinformatics* **32**, 3032-3034 (2016).
- 647 7. Kelleher, J., Etheridge, A.M. & McVean, G. Efficient Coalescent Simulation and  
648 Genealogical Analysis for Large Sample Sizes. *PLoS Comput Biol* **12**, e1004842 (2016).
- 649 8. Speidel, L., Forest, M., Shi, S. & Myers, S.R. A method for genome-wide genealogy  
650 estimation for thousands of samples. *Nature Genetics* **51**, 1321-1329 (2019).
- 651 9. Fan, C., Mancuso, N. & Chiang, C.W.K. A genealogical estimate of genetic relationships.  
652 *The American Journal of Human Genetics* **109**, 812-824 (2022).
- 653 10. Wohns, A.W. *et al.* A unified genealogy of modern and ancient genomes. *Science* **375**,  
654 eabi8264-eabi8264 (2023).
- 655 11. Rasmussen, M.D., Hubisz, M.J., Gronau, I. & Siepel, A. Genome-Wide Inference of  
656 Ancestral Recombination Graphs. *PLOS Genetics* **10**(2014).
- 657 12. Minichiello, M.J. & Durbin, R. Mapping Trait Loci by Use of Inferred Ancestral  
658 Recombination Graphs. *The American Journal of Human Genetics* **79**, 910-922 (2006).
- 659 13. Zhang, B.C., Biddanda, A., Gunnarsson, Á.F., Cooper, F. & Palamara, P.F. Biobank-scale  
660 inference of ancestral recombination graphs enables genealogical analysis of complex  
661 traits. *Nature Genetics* **55**, 768-776 (2023).
- 662 14. Link, V. *et al.* Tree-based QTL mapping with expected local genetic relatedness  
663 matrices. *Am J Hum Genet* **110**, 2077-2091 (2023).
- 664 15. Salehi Nowbandegani, P. *et al.* Extremely sparse models of linkage disequilibrium in  
665 ancestrally diverse association studies. *Nature Genetics* **55**, 1494-1502 (2023).
- 666 16. McVean, G.A.T. & Cardin, N.J. Approximating the coalescent with recombination.  
667 *Philosophical Transactions of the Royal Society B: Biological Sciences* **360**, 1387-1393  
668 (2005).
- 669 17. Song, Y.S. & Hein, J. Constructing minimal ancestral recombination graphs. *J Comput*  
670 *Biol* **12**, 147-69 (2005).
- 671 18. Mirzaei, S. & Wu, Y. RENT+: an improved method for inferring local genealogical trees  
672 from haplotypes with recombination. *Bioinformatics* **33**, 1021-1030 (2017).
- 673 19. Kelleher, J. *et al.* Inferring whole-genome histories in large population datasets. *Nature*  
674 *Genetics* **51**, 1330-1338 (2019).
- 675 20. Schaefer, N.K., Shapiro, B. & Green, R.E. An ancestral recombination graph of human,  
676 Neanderthal, and Denisovan genomes. *Sci Adv* **7**(2021).
- 677 21. Si, Y., Vanderwerff, B. & Zöllner, S. Why are rare variants hard to impute? Coalescent  
678 models reveal theoretical limits in existing algorithms. *Genetics* **217**(2021).

- 679 22. Durbin, R. Efficient haplotype matching and storage using the positional Burrows–  
680 Wheeler transform (PBWT). *Bioinformatics* **30**, 1266-1272 (2014).
- 681 23. Li, N. & Stephens, M. Modeling Linkage Disequilibrium and Identifying Recombination  
682 Hotspots Using Single-Nucleotide Polymorphism Data. *Genetics* **165**, 2213-2233  
683 (2003).
- 684 24. Robinson, D.F. & Foulds, L.R. Comparison of phylogenetic trees. *Mathematical*  
685 *Biosciences* **53**, 131-147 (1981).
- 686 25. Kendall, M. & Colijn, C. Mapping Phylogenetic Trees to Reveal Distinct Patterns of  
687 Evolution. *Molecular Biology and Evolution* **33**, 2735-2743 (2016).
- 688 26. Smith, M.R. Robust Analysis of Phylogenetic Tree Space. *Systematic Biology* **71**, 1255-  
689 1270 (2022).
- 690 27. Chang, C.C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer  
691 datasets. *GigaScience* **4**, s13742-8 (2015).
- 692 28. Rubinacci, S., Delaneau, O. & Marchini, J. Genotype imputation using the Positional  
693 Burrows Wheeler Transform. *PLoS Genetics* **16**(2020).
- 694 29. Marchini, J., Howie, B., Myers, S., McVean, G. & Donnelly, P. A new multipoint method  
695 for genome-wide association studies by imputation of genotypes. *Nat Genet* **39**, 906-  
696 913 (2007).
- 697 30. McCarthy, S. *et al.* A reference panel of 64,976 haplotypes for genotype imputation.  
698 *Nature Genetics* **48**, 1279-1283 (2016).
- 699 31. Browning, B.L., Zhou, Y. & Browning, S.R. A One-Penny Imputed Genome from Next-  
700 Generation Reference Panels. *American Journal of Human Genetics* **103**, 338-348  
701 (2018).
- 702 32. Zhu, J. *et al.* Fast variance component analysis using large-scale ancestral  
703 recombination graphs. *bioRxiv*.
- 704 33. Zhou, W. *et al.* Global Biobank Meta-analysis Initiative: Powering genetic discovery  
705 across human disease. *Cell Genom* **2**, 100192 (2022).
- 706 34. Karczewski, K.J. *et al.* Pan-UK Biobank GWAS improves discovery, analysis of genetic  
707 architecture, and resolution into ancestry-enriched effects. *medRxiv*, 2024.03.  
708 13.24303864 (2024).
- 709 35. Mbatchou, J. *et al.* Computationally efficient whole-genome regression for  
710 quantitative and binary traits. *Nature Genetics* **53**, 1097-1103 (2021).
- 711 36. Chen, M.H. *et al.* Trans-ethnic and Ancestry-Specific Blood-Cell Genetics in 746,667  
712 Individuals from 5 Global Populations. *Cell* **182**, 1198-1213 e14 (2020).
- 713 37. Vuckovic, D. *et al.* The Polygenic and Monogenic Basis of Blood Traits and Diseases. *Cell*  
714 **182**, 1214-1231 e11 (2020).
- 715 38. Astle, W.J. *et al.* The Allelic Landscape of Human Blood Cell Trait Variation and Links to  
716 Common Complex Disease. *Cell* **167**, 1415-1429 e19 (2016).
- 717 39. Sakaue, S. *et al.* A cross-population atlas of genetic associations for 220 human  
718 phenotypes. *Nature Genetics* **53**, 1415-1424 (2021).
- 719 40. Kichaev, G. *et al.* Leveraging Polygenic Functional Enrichment to Improve GWAS Power.  
720 *Am J Hum Genet* **104**, 65-75 (2019).
- 721 41. Hodonsky, C.J. *et al.* Ancestry-specific associations identified in genome-wide  
722 combined-phenotype study of red blood cell traits emphasize benefits of diversity in  
723 genomics. *BMC Genomics* **21**, 228 (2020).
- 724 42. Surakka, I. *et al.* The impact of low-frequency and rare variants on lipid levels. *Nat*  
725 *Genet* **47**, 589-97 (2015).

- 726 43. Willer, C.J. *et al.* Discovery and refinement of loci associated with lipid levels. *Nat Genet*  
727 **45**, 1274-1283 (2013).
- 728 44. Teslovich, T.M. *et al.* Biological, clinical and population relevance of 95 loci for blood  
729 lipids. *Nature* **466**, 707-13 (2010).
- 730 45. Barton, A.R., Sherman, M.A., Mukamel, R.E. & Loh, P.-R. Whole-exome imputation  
731 within UK Biobank powers rare coding variant association and fine-mapping analyses.  
732 *Nature Genetics* **53**, 1260-1269 (2021).
- 733 46. Southam, L. *et al.* Whole genome sequencing and imputation in isolated populations  
734 identify genetic associations with medically-relevant complex traits. *Nat Commun* **8**,  
735 15606 (2017).
- 736 47. Mikhaylova, A.V. *et al.* Whole-genome sequencing in diverse subjects identifies genetic  
737 correlates of leukocyte traits: The NHLBI TOPMed program. *Am J Hum Genet* **108**,  
738 1836-1851 (2021).
- 739 48. Wu, Michael C. *et al.* Rare-Variant Association Testing for Sequencing Data with the  
740 Sequence Kernel Association Test. *The American Journal of Human Genetics* **89**, 82-93  
741 (2011).
- 742 49. Hofmeister, R.J., Ribeiro, D.M., Rubinacci, S. & Delaneau, O. Accurate rare variant  
743 phasing of whole-genome and whole-exome sequencing data in the UK Biobank.  
744 *Nature Genetics* **55**, 1243-1249 (2023).
- 745 50. Harris, K. Using enormous genealogies to map causal variants in space and time. *Nat*  
746 *Genet* **55**, 730-731 (2023).
- 747 51. Baumdicker, F. *et al.* Efficient ancestry and mutation simulation with msprime 1.0.  
748 *Genetics* **220**(2022).
- 749 52. Terhorst, J., Kamm, J.A. & Song, Y.S. Robust and scalable inference of population history  
750 from hundreds of unphased whole genomes. *Nature Genetics* **49**, 303-309 (2017).
- 751 53. Spence, J.P. & Song, Y.S. Inference and analysis of population-specific fine-scale  
752 recombination maps across 26 diverse human populations. *Science Advances* **5**(2023).
- 753 54. Auton, A. *et al.* A global reference for human genetic variation. *Nature* **526**, 68-74  
754 (2015).
- 755 55. Tian, X., Browning, B.L. & Browning, S.R. Estimating the Genome-wide Mutation Rate  
756 with Three-Way Identity by Descent. *The American Journal of Human Genetics* **105**,  
757 883-893 (2019).
- 758 56. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data.  
759 *Nature* **562**, 203-209 (2018).
- 760 57. Browning, B.L., Tian, X., Zhou, Y. & Browning, S.R. Fast two-stage phasing of large-scale  
761 sequence data. *The American Journal of Human Genetics* **108**, 1880-1890 (2021).
- 762 58. Nait Saada, J. *et al.* Identity-by-descent detection across 487,409 British samples  
763 reveals fine scale population structure and ultra-rare variant associations. *Nature*  
764 *Communications* **11**(2020).
- 765 59. Pedregosa, F. *et al.* Scikit-learn: Machine Learning in Python. *Journal of Machine*  
766 *Learning Research* **12**, 2825-2830 (2011).
- 767 60. Koranne, S. *Handbook of Open Source Tools*, 191-200 (Springer US, Boston, MA, 2011).
- 768 61. Gazal, S., Sahbatou, M., Babron, M.-C., Génin, E. & Leutenegger, A.-L. High level of  
769 inbreeding in final phase of 1000 Genomes Project. *Scientific Reports* **5**, 17453-17453  
770 (2015).

- 771 62. Kanai, M., Tanaka, T. & Okada, Y. Empirical estimation of genome-wide significance  
772 thresholds based on the 1000 Genomes Project data set. *Journal of Human Genetics*  
773 **61**, 861-866 (2016).
- 774 63. Berisa, T. & Pickrell, J.K. Approximately independent linkage disequilibrium blocks in  
775 human populations. *Bioinformatics* **32**, 283-5 (2016).
- 776 64. Yang, J. *et al.* Conditional and joint multiple-SNP analysis of GWAS summary statistics  
777 identifies additional variants influencing complex traits. *Nature Genetics* **44**, 369-375  
778 (2012).
- 779 65. Yang, J., Lee, S.H., Goddard, M.E. & Visscher, P.M. GCTA: A Tool for Genome-wide  
780 Complex Trait Analysis. *The American Journal of Human Genetics* **88**, 76-82 (2011).  
781