

## 1 A Genomics England haplotype reference panel and the imputation of the UK Biobank

2  
3 Sinan Shi<sup>1†</sup>, Simone Rubinacci<sup>2</sup>, Sile Hu<sup>3</sup>, Loukas Moutsianas<sup>4,5</sup>, Alex Stuckey<sup>4</sup>, Anna C Need<sup>4</sup>,  
4 The Genomics England Research Consortium, Mark Caulfield<sup>4,5</sup>, Jonathan Marchini<sup>6</sup>, Simon  
5 Myers<sup>1†</sup>

- 6 1. Department of Statistics, University of Oxford, Oxford, United Kingdom
- 7 2. Harvard medical school, Harvard University, Boston, United States
- 8 3. Novo Nordisk Research Centre, Oxford, United Kingdom
- 9 4. Genomics England, London, United Kingdom
- 10 5. Queen Mary University of London, London, United Kingdom
- 11 6. Regeneron Genetic Center, Tarrytown, New York, United States

12 † Corresponding authors.

13

### 14 Abstract

15 The choice of reference panels significantly impacts phasing, imputation and GWAS results. In  
16 this study, we built a haplotype reference panel using the Genomics England (GEL) high-  
17 coverage sequencing dataset, one of the largest genetic variation resources ever collected in the  
18 UK. The resulting reference panel consists of 156,390 haplotypes and 342 million autosomal  
19 variants. The GEL reference panel demonstrates reliable imputation of variants as rare as 1 in  
20 10,000 within the White British population, with an imputation  $r^2$  value of 0.75. The resulting  
21 imputed UKB data (GEL-UKB) contains three times more variants, predominantly rare variants,  
22 compared to the UKB data previously imputed using the HRC and UK10K reference panel. The  
23 GEL-UKB presents a unique opportunity for the reliable discovery of rare associations across the  
24 whole genome, especially within the regions not covered by the exome sequencing data. Rare  
25 variant signals with high confidence are predominantly from rare coding variants, implying  
26 firstly, a probable tendency for existing rare non-coding mutations to not reach a disruptive level  
27 comparable to that of coding variants. Secondly, it raises the possibility that the current sample  
28 size of UK Biobank may be insufficient for detecting rare variants with a moderate effect size,  
29 even with the whole genome sequencing. The resulting GEL phased haplotype reference panel  
30 has been made available on the GEL platform and widely used by GEL users. Our GEL imputed  
31 UKB data has been adopted as one of the UKB official imputed data resources (Data Field  
32 21008).

33

34

35

36

37

## 38 Main

39 Genomics England (GEL) has carried out whole genome sequencing (WGS) of over 120,000  
40 genomes from over 80,000 individuals taking part in the 100,000 Genomes Project, using an  
41 average sequencing coverage depth of  $\sim 30\times^1$ . The recruitment strategy focussed on patients with  
42 rare disease (disorders affecting  $< 1$  in 2000 people) and cancer, and their close relatives, across  
43 hospitals in England. We constructed a GEL phased reference panel based on 78,195 high-  
44 coverage sequencing germline genomes, with a diverse ethnic representation. The high degree of  
45 relatedness among the samples enhances the power of filters, such as the Mendel error filter, for  
46 eliminating false positive variant sites identified in the sequencing data, and also leads to more  
47 accurate phasing and imputation of rare variants. In particular, it enables even variants found in  
48 only one or two individuals to be phased through transmission, a task which is more difficult in  
49 the absence of related samples or phase information in sequencing reads<sup>2</sup>.

50  
51 The resulting GEL reference panel consists of 341,922,205 autosomal variants, with 31,502,703  
52 (9.26%) being INDELs with an average length of 5bp and a maximum length of 50bp. The  
53 majority of the variants in the GEL reference panel are rare. 287.2 million (84.1%) of identified  
54 variants possess an allele frequency lower than 0.0001, including 66.7 million (19.5%) singletons  
55 and 91.1 million (26.7%) doubletons. We compared the variants in GEL reference panel to the  
56 widely used TOPMed r2<sup>3</sup> and HRC<sup>4</sup> panels and found GEL has 8 times and 1.1 times more  
57 variants than the HRC and TOPMed panels respectively (**Figure 1b** and **Supplementary Figure**  
58 **1**). Due to the use of mostly low coverage sequencing technology, the HRC dataset has limited  
59 numbers of rare variants, especially those with  $AF \leq 10^{-4}$ . While the numbers of rare variants  
60 captured in TOPMed and GEL are similar, around half of the ultra-rare variants ( $AF \leq 10^{-4}$ )  
61 from GEL and TOPMed are non-shared across the panels (**Supplementary Figure 1**). As  
62 expected, all three panels capture a similar set of more common ( $AF > 10^{-2}$ ) variants, with less  
63 than 4% unique to each panel (**Supplementary Figure 1**), indicating common variants are  
64 largely saturated.

65  
66 The GEL reference panel can be used as a powerful resource for phasing European and South  
67 Asian samples, due to their strong representation in the dataset. We compared the phasing  
68 accuracy achievable using the GEL and HRC reference panel across 26 diverse populations from  
69 the 1000 Genomes project (**Methods**). GEL phasing of these samples achieved lower switch  
70 error rates than HRC phasing, across the CEU (Northern European from Utah), African, South  
71 Asian and East Asian ancestry populations (**Figure 1a**), with HRC only showing improved  
72 performance for South American samples, which are not significantly represented in GEL. GEL  
73 phasing switch error rates are 0.18%, 0.33%, 0.31% and 0.73% for European, African, South  
74 Asian and East Asian samples respectively.

75  
76 A primary use of the GEL will be as a reference panel for genotype imputation of other datasets.  
77 We assessed the imputation accuracy among 2,405 1,000 Genomes samples, using the GEL,  
78 TOPMed and HRC reference panels. We used genotypes at the 716,473 autosomal bi-allelic SNP  
79 positions on the UK Biobank Axiom array<sup>5</sup> to impute all non-array sites using each reference  
80 panel (**Methods**). Squared correlation  $r^2$  between the imputed allele dosages and true genotypes  
81 were calculated, stratified by the independently estimated gnomAD (v3.3.1) minor allele  
82 frequency<sup>6</sup>. As we focus on showing the overall performance of the reference panel across  
83 different allele frequencies, only variants present within gnomAD are shown. As a result, the

84 number of tested variants differs across reference panels. GEL achieved higher  $r^2$  than HRC in  
85 all allele frequency bins for all ethnicities (**Supplementary Figure 4**) and outperforms the  
86 TOPMed panel in White British (GBR) and South Asian (SAS) samples, especially for rarer  
87 variants: at  $MAF < 10^{-5}$ , the GEL imputation  $r^2$  for GBR samples is 0.6, compared to 0.3 and  
88 0.29 using TOPMed and HRC, respectively (**Figure 1c**). The TOPMed panel outperforms GEL  
89 in African, American and East Asian samples due to its better representation from these groups  
90 (**Supplementary Figure 4**).

91  
92 We used the GEL panel to impute 488,315 UK Biobank samples at 342,573,817 variants,  
93 producing a “GEL-UKB” dataset; we compared to the corresponding HRC and UK10K-imputed  
94 “HRC-UKB”<sup>5</sup>. GEL-UKB has around 3 times more variants than HRC-UKB, 3.5 times more  
95 missense variants, and 6.6 times more “high impact consequence” variants (**Supplementary**  
96 **Table 5**). The imputed information scores (**Method**) were higher for GEL-UKB than HRC-UKB  
97 for 87% of the variants that are in common, while 98% (78%) of GEL-imputed variants at  
98 frequency below  $10^{-4}$  ( $10^{-5}$ ) exceeded a threshold of 0.3, vs 78% (54%) for HRC  
99 (**Supplementary Figure 2-3**).

100  
101 To demonstrate the use of GEL-UKB, exemplar GWAS were carried out on four quantitative  
102 traits, including standing height (HEIGHT), body mass index (BMI), systolic (SBP) and diastolic  
103 (DBP) blood pressure, with variant testing using REGENIE<sup>7</sup>. Across all four traits, we found  
104 31,699 and 30,711 significant ( $P\text{-value} < 5 \times 10^{-8}$ ) rarer variant associations ( $MAF < 0.05$ )  
105 from GEL-UKB and HRC-UKB, respectively. The GEL-UKB common variant associations  
106 were also less likely to be subjected to false associations than HRC-UKB (**Supplementary**  
107 **Notes; Supplementary Table 2; Supplementary Figure 6-8**). A recent exome-sequencing  
108 based association study reported 31, 0, 1, and 2 rarer ( $MAF < 0.05$ ) genome-wide significant ( $P\text{-}$   
109  $value < 2.18 \times 10^{-11}$ ) variant-trait associations across HEIGHT, BMI, SBP and DBP,  
110 respectively<sup>8</sup>. We discovered 70% of these associations using GEL-UKB, compared to 56%  
111 using HRC-UKB at the same p-value threshold. Relaxing the GEL p-value threshold to  
112  $5 \times 10^{-8}$ , GEL-UKB identifies 76% of these associations (**Supplementary Table 3**). When we  
113 compare to the UKB whole exome imputation results<sup>9</sup>, all but 4 out of the 28 exome imputation  
114 likely-causal rare coding variants associated with standing height ( $p\text{-value} < 5 \times 10^{-8}$ ) are found  
115 to be significant using GEL-UKB, whereas all but 9 of such variants are found to be significant  
116 using HRC-UKB (**Supplementary Figure 9**).

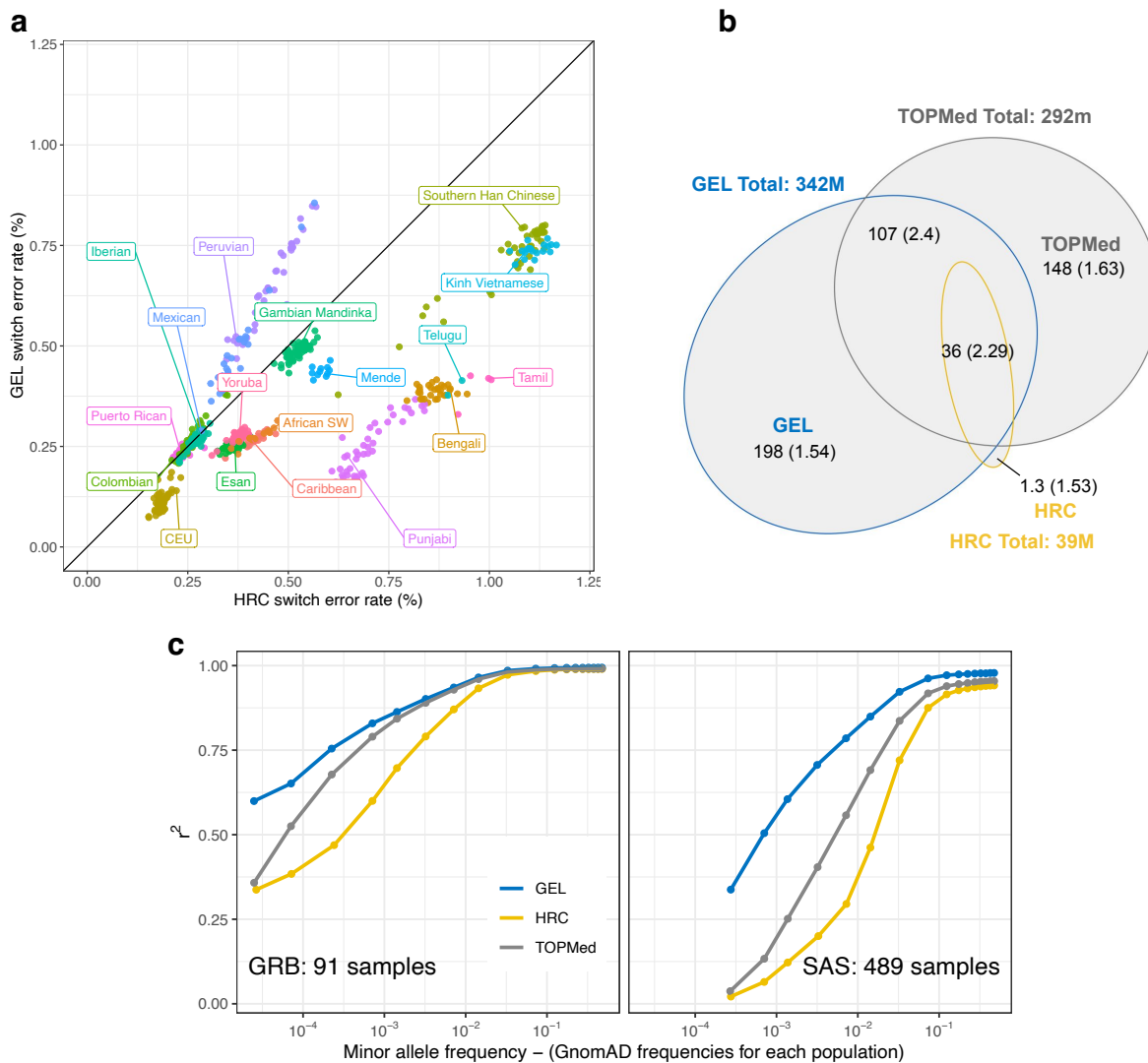
117  
118 This comparison in the exonic portion of the genome provides confidence that whole-genome  
119 imputation using GEL can identify most associations directly observed using sequencing. We  
120 next compared the performance of GEL-UKB to the widely used imputed genotypes available  
121 for the full set of UKB samples, HRC-UKB, and we examined those novel associations  
122 identified using GEL-UKB. First examining shared signals (mainly at common sites), we saw a  
123 useful improvement in fine-mapping (**method**) using GEL-UKB vs. HRC-UKB. 44% of the  
124 GEL-UKB based 95% credible sets contain fewer SNPs, while 25% contain more SNPs (**Figure**  
125 **2b; Supplementary Table 4**), with the remainder identical in size.

126  
127 A more dramatic difference is observed for rare variants: independent rare variant associations  
128 ( $MAF < 5 \times 10^{-4}$ ), accompanied by high estimated effect sizes (**Figure 2a**) required to reach  
129 statistical significance at these frequencies, are almost exclusively discovered by GEL-UKB

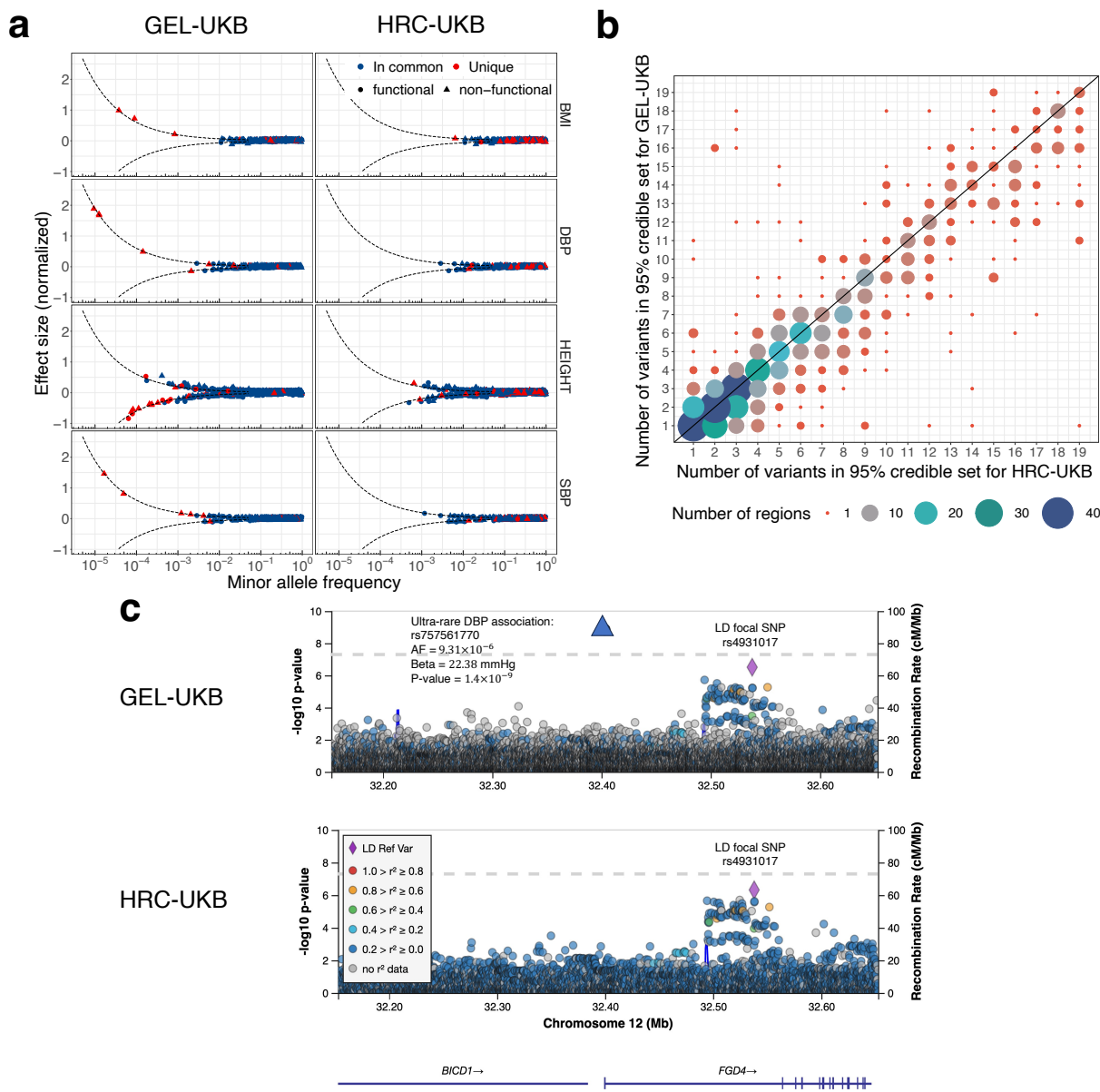
130 **(Figure 2a)**. For example, GEL-UKB detected a novel ultra-rare association signal for DBP at  
131 rs757561770 in FGD4, with allele frequency  $9.31 \times 10^{-6}$ . Common variants in FGD4 has  
132 previously been reported to be associated with hypertension<sup>10</sup> **(Figure 2c)**. Interestingly, this  
133 SNP is intronic and does not show strong linkage disequilibrium ( $r^2 > 0.7$ ) with any coding  
134 variant within the GEL panel **(Supplementary Table 6)**.  
135

136 Because we test the entire genome, our results allow us to investigate whether large-effect  
137 mutations (which in our example GWAS are only found at low frequency; **Figure 2b**) occur in  
138 coding or non-coding DNA. We identified 27 independent large-effect/rare-variant signals ( $AF <$   
139  $0.001$ ), across the four traits using step-wise regression **(Method)**. Of these, 17 were either  
140 coding ( $n=7$ ) or in strong LD ( $r^2 > 0.7$ ) with a coding variant. An additional 1 is associated with  
141 splice-site variants and 2 with variants in 5' UTRs or 3' UTRs of genes **(Supplementary Table**  
142 **6)**. In total, 62% of all the rare variant associations and 88% of the strongest associations ( $p$ -  
143 value  $< 2.18 \times 10^{-11}$ ) were associated with genic sequences **(Supplementary Table 6)**. If  
144 replicated for other phenotypes, this implies that it could be likely rare for variation in other non-  
145 coding regions such as enhancers to achieve dramatic trait effects – despite such regions  
146 dominating GWAS signals overall<sup>11</sup>. Because it seems likely that non-coding mutations *are* able  
147 to strongly disrupt the binding of individual transcription factors, this might imply that (except in  
148 5' UTR and 3' UTR regions) no one transcription factor plays an essential role in the  
149 overwhelming majority of cases. Nonetheless, we still observed several cases implicating only  
150 non-genic sites, for example an intronic signal for decreasing height (rs570873498;  $AF=0.0002$ )  
151 at SLC12A1, a gene known to be associated with height and Bartter syndrome, whose symptoms  
152 include growth retardation<sup>12</sup>. We anticipate that despite their modest effect sizes and limiting  
153 power at present (likely, even if genomes are fully sequenced), the number of non-coding  
154 associations will likely increase rapidly in future, once sample sizes become larger. Moreover,  
155 our results imply imputation will be highly effective in identifying such associations, even for  
156 rare variants.  
157

158 One unexpected finding for height from our analysis was a cluster of five independent low-  
159 frequency associations with height on chromosome 6 **(Supplementary Table 6; Extended data**  
160 **table)**, including the rare missense variant rs957675208, in a region not reported by the previous  
161 exome sequencing<sup>8</sup> and exome imputation<sup>9</sup> analyses, or by HRC-UKB (low imputation INFO).  
162 Strikingly, rs957675208 in HMGA1 shows the strongest height-increasing impact of any SNP in  
163 the whole dataset, equivalent to gaining 3.5 cm of height. On further examination, three of these  
164 four variants are missense mutations and the remaining two 5' UTR variants are in a gene not  
165 annotated in the exome studies. This gives one example of how the complete genome-wide  
166 coverage of the GEL-UKB data allows for additional findings compared to previous approaches.  
167



168  
 169 **Figure 1:** a) Phasing quality for 589 high coverage 1,000 Genome children from mother-father-child trio families,  
 170 using HRC and GEL reference panels. b) Venn diagram comparing numbers of variants from the GEL, HRC and  
 171 TOPMed reference panels. The numbers show the variant count (in millions of variants), followed by the Ts/Tv ratio  
 172 of these variants in brackets. c) Imputation performance, measured by  $r^2$  (**Methods**), for imputation of 1000  
 173 genomes samples from the *White British* (left) and South Asian (right) groups, using three different reference panels  
 174 (labels). The variants are stratified by GnomAD allele frequency ( $v3.3.1$ )<sup>6</sup> of their corresponding population.  
 175  
 176



177  
 178 Figure 2: a) A set of independent genome-wide significant ( $p < 5 \times 10^{-8}$ ) associations identified by step-wise  
 179 regressions (conditioned joint analysis), and with INFO > 0.8, are plotted versus their imputed allele frequency (x-  
 180 axis). The blue colour represent variants that were flagged by step-wise regressions in one dataset and also showed a  
 181 significant GWAS association in the other dataset; The red colour indicates that the variant is unique to each dataset.  
 182 The shape of the data points reflects the predicted consequences of the variants as determined by VEP. Dots  
 183 represent functional variants, including stop gained, stop lost, splice donor/acceptor, frameshift, in-frame  
 184 insertion/deletion, and missense and the triangles indicate non-functional variants. The dotted lines indicates the  
 185 smallest effect sizes that can be captured by the p-value threshold ( $p < 5 \times 10^{-8}$ ). b) Comparison of the number of  
 186 variants in the 95% credible sets for GEL-UKB and HRC-UKB fine-mapping results for standing height (capped at  
 187 20 variants; **Methods**). The circle sizes represent the number of fine-mapping regions showing each combination;  
 188 plots below the diagonal correspond to GEL-UKB having fewer variants in the credible set compared to HRC-UKB.  
 189 c) The LocusZoom plot of ultra-rare variant association (rs757561770) detected by GEL-UKB. The color indicates  
 190 the LD between SNPs and the focal SNP rs4931017, showing that rs757561770 is in low LD with the focal SNP  
 191 ( $r^2 = 6.57 \times 10^{-6}$ ). The blue lines show the recombination rate of the region.

## 192 **Methods**

### 193 **Genomics England high coverage sequencing data**

194 The Genomics England 100,000 Genomes Project was launched in 2013, focusing on rare  
195 diseases and cancer. Over 120,000 genomes have been sequenced. It comprises genomes from  
196 73,700 rare disease (disorders affecting  $\leq 1$  in 2000 persons) patients and their close relatives, and  
197 46,539 genomes from cancer patients<sup>1</sup>. The GEL reference panel described in this paper is built  
198 on the aggregated dataset (aggV2), comprising 78,195 samples from both rare disease and cancer  
199 germline genomes. Samples were sequenced with 150bp paired-end reads on the IlluminaHiSeq  
200 X platform and processed with the Illumina North Star Version 4 Whole Genome Sequenced  
201 Workflow (iSAAC Aligner v03.16.02.19 and Starling small variant caller v2.4.7), and aligned to  
202 the GRCh38 human reference genome. The individual gVCF files were aggregated into multi-  
203 sample VCF files using Illumina gVCF genotyper and normalised with vt v0.57721. The  
204 aggregated multi-sample VCF dataset (aggV2) comprises over 722 million initial called SNPs  
205 and short indels ( $\leq 50$ bp). Multi-allelic variants were decomposed into biallelic variants. The  
206 dataset includes 49,641 samples (63.48%) from individuals self-identifying as White British, 4,100  
207 (5.24%) as “Other White”, 2,885 (3.69%) as Pakistani, 1,860 (2.3%) as Black, 1,751 (2.24%) as  
208 Indian, and 12,277 samples (15.7%) as “Unknown”. The large White British and relatively large  
209 South Asian sample size made GEL an ideal reference panel for phasing and imputing UK  
210 Biobank, which has a similar ethnic composition<sup>5</sup>. According to the self-reported data, only  
211 27,346 samples (34.97%) have no relatives in the reference panel. 11,584 (14.81%), 32,679  
212 (41.79%), and 6,586 (8.43%) samples are one of 2, 3 and  $>3$  family members in the dataset  
213 respectively. We identified 12,816 (16.39%) samples as members of duo families and 35,106  
214 (44.9%) as members of trio families, while 30,273 (38.71%) samples are treated unrelated for  
215 phasing (**Supplementary Notes**).

### 218 **Quality Control**

219 Prior to the quality control (QC) described here, sample level QC was carried out by Genomics  
220 England informatics team on variants called one sample at a time. We conducted additional  
221 quality control by pooling information across samples, to remove false positive sites. Specifically  
222 we utilised aggregated VCFs, considering genotype quality, depth, missingness, allele balance,  
223 Mendel errors, Hardy-Weinberg equilibrium, and gnomAD<sup>6</sup> allele frequency concordance.  
224 Because singletons observed in unrelated samples are very hard to phase accurately these sites  
225 were removed. We applied two sets of QC rules. First, we applied a stringent rule set applied to  
226 all sites, including those *de novo* in Genomics England and very rare sites. Second, we applied a  
227 more lenient group of filters for relatively common sites ( $AF > 0.001$ ) that additionally showed  
228 support from independent external datasets (TOPMed, HRC, 1000 Genomes, GnomAD), to  
229 avoid removing a proportion of genuine sites (e.g. for a modest number of Mendel errors). For  
230 these sites, if they failed our stringent filters but passed with somewhat less stringent  
231 missingness, Mendel error and gnomAD frequency concordance thresholds, we included them,  
232 after separate phasing conditional on the phase of sites passing the more stringent thresholds, i.e.  
233 in a manner which did not impact the stringent sites. These sites were incorporated in the final  
234 dataset, but with a QC flag indicating their slightly lower reliability. Overall, our filters reduced  
235 the initial number of sites from 722 million to 342 million. (**Supplementary Notes and**  
236 **Supplementary Table 1**)

237

## 238 **Phasing the GEL reference panel**

239 We used a multi-stage phasing strategy leveraging the relatedness within GEL, in particular  
240 allowing phasing of singletons were possible.

- 241 1. We used the makeScaffold software (<https://github.com/odelaneau/makeScaffold>) to  
242 determine the phase of duo and trio samples (**Supplementary Notes**) by direct  
243 transmission information (this phases most sites in these samples).
- 244 2. For remaining unphased genotypes in these related samples, with phases undetermined  
245 due to heterozygosity or missing data, phases were inferred using SHAPEIT4.2.2<sup>13</sup>, with  
246 the phased genotypes from step 1 as a scaffold.
- 247 3. To phase genotypes in the unrelated samples, we first phased the common variants ( $AF >$   
248  $0.01$ ) one chromosome at a time, using SHAPEIT4.2.2 and now using the genotypes (at  
249 these common sites) from step 1 and 2 in the related samples as a reference panel.
- 250 4. Finally, to phase the remaining sites: genotypes at rare variants in unrelated samples, we  
251 using SHAPEIT4.2.2 with the phased samples from steps 1-2 as a reference panel, and  
252 the phased common variants from step 3 as a scaffold for these samples.
- 253 5. For sites only passing our lenient filters (see **“Quality Control”** section above and  
254 **Supplementary Notes**) we used the results of step 4, for the sites on the UKB Axiom  
255 array sites passing the stringent filters, as a scaffold, and then used SHAPEIT4.2.2 on the  
256 remaining genotypes.

257 Phasing for steps 1 and 3 was done at the entire chromosome level; for steps 2 and 4 was carried  
258 out in regions of approximately 300,000 sites, with 30,000 sites on each side as buffer. The  
259 resulting phased regional segments were merged and concatenated using bcftools<sup>14</sup>. These  
260 phasing steps were computationally intensive, and took about 6,500 CPU days in total to  
261 accomplish. The phased reference panel is stored in VCF format and has been made available for  
262 all Genomics England registered users on the GEL trusted research environment.

263

## 264 **Estimation of 1000 Genome trio phasing switch error rate**

265 Phasing accuracy is important for direct biological interpretation of variants within GEL, as well  
266 as ensuring high-quality imputation in other samples and other downstream applications. We  
267 assessed the ability of the GEL panel to phase such external samples. Specifically, we phased the  
268 parents of mother-father-child trios included in the 1000 Genomes Project (but not HRC or GEL)  
269 using the reference panels from HRC and GEL. We then assessed the resulting phase accuracy,  
270 by comparing phased haplotypes to those directly inferred using inheritance patterns to the child  
271 in each trio. The HRC reference panel was lifted over from the GRCh37 to the GRCh38  
272 reference genome using GATK Picard LiftoverVCF<sup>15</sup>. The original GRCh37 HRC reference  
273 panel has 39,131,578 autosomal variants. 13,813 variants were removed either due to the  
274 incompatibility between reference genomes or mismatching chromosome between the two  
275 reference genomes. The resulting autosomal GRCh38 HRC reference panel contain 39,115,765  
276 variants and 27,165 samples. 1000 Genome samples within the HRC reference panel were  
277 removed.

278

279 We analysed only sites passing 1000 Genome data<sup>16</sup> filters. The phasing test was carried out on  
280 589 trio families from diverse ethnic backgrounds, using SHAPEIT 4.2.2<sup>13</sup>. We tested all the



281 heterozygous 1000G sites for each individual reference panel, yielding a total of  $1.04 \times 10^9$   
282 heterozygous sites (1.76 million per trio family) for the HRC panel and  $1.16 \times 10^9$  (1.9 million  
283 per trio family) for the GEL panel.

284

### 285 **Imputation testing of 1000 Genomes samples**

286 We used 2,405 1000 Genomes samples to test the relative performance of imputation based on  
287 the GEL, TOPMed and HRC imputation panels. We first performed quality control on the 1000  
288 Genomes data, by removing sites which either possess a missingness larger than 5% or failed a  
289 Hardy Weinberg equilibrium test, by having a p-value smaller than  $10^{-10}$  in any of the 26 1000  
290 Genome populations. We then masked genotypes in 1000 Genomes sequencing samples, except  
291 the sites existing in the UK Biobank Axiom array, to mimic imputation using this array. This  
292 gave 716,473 bi-allelic SNPs across all autosomes. The pseudo-SNP array dataset was then  
293 phased one chromosome at a time using SHAPEIT4.1.2<sup>13</sup>. TOPMed imputation was carried out  
294 using the TOPMed imputation server with the TOPMed r2 reference panel and the imputation  
295 software minimac4 1.5.7<sup>17</sup>. IMPUTE5<sup>18</sup> was used to impute from the GEL and HRC reference  
296 panels. We stratified imputation results into 6 groups : 661 African (AFR), 347 American  
297 (AMR), 504 Eastern Asian (EAS), 489 South Asian (SAS), 313 non-Finnish European (NFE)  
298 samples and 91 British (GBR) samples.

299

### 300 **UK Biobank imputation using the GEL reference panel**

301 The UK Biobank SNP array data consists of 784,256 autosomal variants. We removed the set of  
302 113,515 sites identified by the previous centralized UK Biobank analysis as failing quality  
303 control<sup>5</sup> and an additional set of 39,165 sites failing a test of Hardy-Weinberg equilibrium on  
304 409,703 White British samples, with the p-value threshold of  $10^{-10}$ . The resulting UK Biobank  
305 SNP array data was mapped from the GRCh37 to GRCh38 genome build, using the GATK  
306 Picard LiftOver tool. Alleles with mismatching strand but matching alleles were flipped. 495  
307 sites were removed due to incompatibility between the two reference genomes, resulting in a  
308 final SNP array incorporating 631,081 autosomal variants that we used for phasing and  
309 imputation. Haplotype estimation of the SNP array data is a prerequisite for imputation. Phasing  
310 was carried out one chromosome at a time using SHAPEIT4.2.2 without a reference panel, using  
311 the full set of UK Biobank samples. We ran SHAPEIT4 using its default 15 MCMC iterations  
312 and 30 threads. The runtime varied from 2 hours to 30 hours for each chromosome. Imputation  
313 of normal filter set and lenient filter set SNPs was carried out independently. Autosomal  
314 imputation using the GEL reference panel was performed using IMPUTE5 (v1.1.4). The SNP  
315 array data was divided into 408 consecutive and overlapping chunks with roughly 5mb for each  
316 chunk and 2.5mb buffer across the genome, using the Chunker program in IMPUTE5<sup>18</sup> and each  
317 chunk was further divided into 24 sample batches with each batch containing 20,349 samples.  
318 IMPUTE5 was run on each of the 9,792 subsets using a single thread and default settings, at a  
319 speed less than 4 minutes per genome, resulting in a total time of around 1,200 CPU days to  
320 impute all UK Biobank samples.

321

### 322 **Genome-wide association studies**

323 We selected four quantitative traits to demonstrate the GWAS performance of the GEL imputed  
324 UK Biobank data (GEL-UKB), compared to the HRCUK10K imputed UKB (HRC-UKB) data  
325 on 429,460 white British samples. These traits are standing height (HEIGHT), body mass index

326 (BMI), systolic (SBP) and diastolic (DBP) blood pressure. Variants with minor allele count  
327 lower than 5 are not included in testing. The trait measures are transformed using rank inverse  
328 normal transformation (RINT) within sexes to ensure normally distributed input phenotypes and  
329 reduce the likelihood of false positives due to outliers.

330  
331 Samples between 40 to 70 years old are included and for each data point, outliers that are above  
332  $\pm 4$  standard deviation from the mean value were removed<sup>5</sup>. SBP and DBP values are based on  
333 automated blood pressure readings, substituting in manual reading values when automated  
334 readings are not available. We calculated the mean SBP and DBP values from two automated (n  
335 = 418,755) or two manual (n = 25,888) blood pressure measurements. For individuals with one  
336 manual and one automated blood pressure measurement (n = 13,521), we used the mean of these  
337 two values. For individuals with only one available blood pressure measurement (n = 413), we  
338 used this single value. After calculating blood pressure values, we adjusted for blood pressure-  
339 lowering medication (n=94,289) use by adding 15 and 10 mmHg to SBP and DBP,  
340 respectively<sup>19</sup>, for individuals on such medication.

341  
342 GWAS effect size estimates and p-values were obtained using REGENIE<sup>7</sup>. We used the UKB  
343 SNP array data to estimate the LOCO predictors in REGENIE Step 1 and the imputed data for  
344 Step 2, accounting for sex, age, sex squared, sex  $\times$  age, and 20 principal components as  
345 covariates<sup>7</sup>. The association tests for GEL imputed UKB (GEL-UKB) and HRCUK10K imputed  
346 UKB (HRC-UKB) used the identical setup. The HRC-UKB summary statistics of the association  
347 tests were mapped using Picard LiftOver from GRCh37 to GRCh38 to compare the results with  
348 GEL-UKB. In all analysis, we used an INFO threshold of 0.3 for common imputed variants  
349 (MAF>0.05) and 0.8 for rare imputed variants (MAF $\leq$ 0.05). **Supplementary Figure 5** shows  
350 higher INFO threshold are effective for detecting false positive rare associations.

351

### 352 **Bayesian fine-mapping**

353 Bayesian fine-mapping credible set size comparison was carried out on 1,660, 711, 505 and 546  
354 non-overlapping regions for HEIGHT, BMI, SBP and DBP respectively based on HRC-UKB  
355 GWAS summary statistics. These regions were defined by the following procedure. First,  
356 candidate regions were identified with width 0.125 centiMorgans plus 25 kb on each side of a  
357 significant marker. Overlapping candidate regions were successively merged until there are no  
358 remaining regions overlapping. We removed 60, 30, 33, and 51 regions for above traits  
359 respectively, in which GEL-UKB showed no significant sites (p-value  $< 5 \times 10^{-8}$  in GWAS) for  
360 each trait. The recombination rate is based on the HapMap genetic map<sup>20</sup>. The detail description  
361 of this approach can be found in Maller et al., and Bycroft et al.<sup>5,21</sup>

362

363 For each region, we assume a single causal variant – call this model  $M$ . Given this, define model  
364  $M_i$  to be the model where SNP  $i$  is the causal variant. We seek the probability of  $M_i$  given the  
365 data and that model  $M$  is true. This posterior  $Pr(M_i|\mathbf{X}, M)$  can be written in terms of the Bayes  
366 factor relating the probability of the data given  $M_i$  versus the probability of the data under the  
367 null model with no associated SNP in the region,  $BF_i$ . Further,  $BF_i$  can be approximated by an  
368 asymptotic Bayesian factor ( $ABF_i$ ):

369

370 
$$Pr(M_i|X, M) = \frac{BF_i}{\sum_{i=1}^k BF_i} \approx \frac{ABF_i}{\sum_{i=1}^k ABF_i}$$

371  
372  $ABF_i$  can be calculated using the standard error ( $V_i$ ) and Z score ( $z$ ) estimated by REGENIE<sup>5</sup>. In  
373 each region, the smallest possible 95% credible set of potential causal markers can be obtained  
374 by successively including the sites with the highest probabilities, to accumulatively reach 0.95.  
375 Model  $M$  requires a prior for the (Gamma distribution) on effect sizes; we choose this prior  $W$  to  
376 have parameters  $0.2^2$  and  $0.02^2$ , but found the results are not particularly sensitive to the choice  
377 of the prior.

378

### 379 **Conditional joint analysis: step-wise regression**

380 A standard GWAS uses marginal model considering one variant at a time, while a joint model  
381 considers all the selected variants and estimates their joint effect simultaneously. In order to  
382 remove rare variant signals that are explained by stronger signals at more common nearby  
383 SNPs<sup>8</sup>. We performed a conditional joint analysis via a stepwise forward selection procedure,  
384 considering each chromosome separately. First we defined the set  $S$  of genome-wide significant  
385 variants in one chromosome ( $P$ -value  $< 5 \times 10^{-8}$ ) in the marginal regression using REGENIE.  
386 We initialized a set of variants  $R$  as the most significant variant in the marginal regression.  
387 Given the current value of  $R$ , we calculate the  $P$ -value of all the remaining variants in  $S$  one at a  
388 time, conditioned on  $R$  and the covariates used for the initial GWAS. We then move the variant  
389 with the smallest conditional  $P$ -value from  $S$  to  $R$ , until this smallest  $P$ -value is no longer  
390 genome-wide significant. This approach identifies a set of variants that are independently  
391 significant, and account for all the genome-wide association signals (note that this set is not  
392 unique), while also accounting for linkage disequilibrium between sites. To identify rare causal  
393 variants within UKBB found using GEL-UKB imputation, we considered only those variants  
394 found by this stepwise forward selection approach. The full conditional joint analysis results can  
395 be found in the **Extended data table**.

### 396 **Data availability**

397 The GEL haplotype reference panel is available within the GEL trusted research environment to  
398 approved researchers only. The imputed UK Biobank data imputed using the GEL haplotype  
399 reference panel is available to those with approved access to the UK Biobank resource and  
400 described on the UK Biobank showcase here

401 <https://biobank.ctsu.ox.ac.uk/crystal/field.cgi?id=21008>

### 402 **Acknowledgements**

403 We thank the Wellcome Trust for funding (200186/Z/15/Z to JM, SM) and (212284/Z/18/Z to  
404 SM). Work conducted under UKB applications (48031 and 27960). This research was made  
405 possible through access to data in the National Genomic Research Library, which is managed by  
406 Genomics England Limited (a wholly owned company of the Department of Health and Social  
407 Care). The National Genomic Research Library holds data provided by patients and collected by  
408 the NHS as part of their care and data collected as part of their participation in research. The  
409 National Genomic Research Library is funded by the National Institute for Health Research and  
410 NHS England. The Wellcome Trust, Cancer Research UK and the Medical Research Council  
411 have also funded research infrastructure. This work is part of the research portfolio of the

412 National Institute for Health and Social Care Research Barts Biomedical Research Centre. MC is  
413 funded by the Barts Charity and is an NIHR Senior Investigator alumnus.  
414

415 **References**

416

417 1. 100,000 Genomes Pilot on Rare-Disease Diagnosis in Health Care — Preliminary Report. *N.*

418 *Engl. J. Med.* **385**, 1868–1880 (2021).

419 2. Delaneau, O., Howie, B., Cox, A. J., Zagury, J.-F. & Marchini, J. Haplotype Estimation Using

420 Sequencing Reads. *Am. J. Hum. Genet.* **93**, 687–696 (2013).

421 3. Taliun, D. *et al.* Sequencing of 53,831 diverse genomes from the NHLBI TOPMed Program.

422 *bioRxiv* 563866 (2019) doi:10.1101/563866.

423 4. McCarthy, S. *et al.* A reference panel of 64,976 haplotypes for genotype imputation. *Nat.*

424 *Genet.* **48**, 1279–1283 (2016).

425 5. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data. *Nature*

426 **562**, 203–209 (2018).

427 6. Karczewski, K. J. *et al.* The mutational constraint spectrum quantified from variation in

428 141,456 humans. *Nature* **581**, 434–443 (2020).

429 7. Mbatchou, J. *et al.* Computationally efficient whole-genome regression for quantitative and

430 binary traits. *Nat. Genet.* **53**, 1097–1103 (2021).

431 8. Backman, J. D. *et al.* Exome sequencing and analysis of 454,787 UK Biobank participants.

432 *Nature* **599**, 628–634 (2021).

433 9. Barton, A. R., Sherman, M. A., Mukamel, R. E. & Loh, P.-R. Whole-exome imputation within

434 UK Biobank powers rare coding variant association and fine-mapping analyses. *Nat. Genet.*

435 **53**, 1260–1269 (2021).

- 436 10. Takeuchi, F. *et al.* Interethnic analyses of blood pressure loci in populations of East  
437 Asian and European descent. *Nat. Commun.* **9**, 5052 (2018).
- 438 11. Watanabe, K. *et al.* A global overview of pleiotropy and genetic architecture in complex  
439 traits. *Nat. Genet.* **51**, 1339–1348 (2019).
- 440 12. Yengo, L. *et al.* A saturated map of common genetic variants associated with human  
441 height. *Nature* **610**, 704–712 (2022).
- 442 13. Delaneau, O., Zagury, J.-F., Robinson, M. R., Marchini, J. L. & Dermitzakis, E. T.  
443 Accurate, scalable and integrative haplotype estimation. *Nat. Commun.* **10**, 5436 (2019).
- 444 14. Li, H. A statistical framework for SNP calling, mutation discovery, association mapping  
445 and population genetical parameter estimation from sequencing data. *Bioinformatics* **27**,  
446 2987–2993 (2011).
- 447 15. DePristo, M. A. *et al.* A framework for variation discovery and genotyping using next-  
448 generation DNA sequencing data. *Nat. Genet.* **43**, 491–498 (2011).
- 449 16. Byrska-Bishop, M. *et al.* High coverage whole genome sequencing of the expanded 1000  
450 Genomes Project cohort including 602 trios. *bioRxiv* 2021.02.06.430068 (2021)  
451 doi:10.1101/2021.02.06.430068.
- 452 17. Fuchsberger, C., Abecasis, G. R. & Hinds, D. A. minimac2: faster genotype imputation.  
453 *Bioinformatics* **31**, 782–784 (2015).
- 454 18. Rubinacci, S., Delaneau, O. & Marchini, J. Genotype imputation using the Positional  
455 Burrows Wheeler Transform. *PLOS Genet.* **16**, e1009049 (2020).

- 456 19. Tobin, M. D., Sheehan, N. A., Scurrah, K. J. & Burton, P. R. Adjusting for treatment  
457 effects in studies of quantitative traits: antihypertensive therapy and systolic blood pressure.  
458 *Stat. Med.* **24**, 2911–2935 (2005).
- 459 20. Frazer, K. A. *et al.* A second generation human haplotype map of over 3.1 million SNPs.  
460 *Nature* **449**, 851–861 (2007).
- 461 21. Maller, J. B. *et al.* Bayesian refinement of association signals for 14 loci in 3 common  
462 diseases. *Nat. Genet.* **44**, 1294–1301 (2012).
- 463