

AperTO - Archivio Istituzionale Open Access dell'Università di Torino

## A bird's-eye view of Italian genomic variation through whole-genome sequencing

**This is a pre print version of the following article:**

*Original Citation:*

*Availability:*

This version is available <http://hdl.handle.net/2318/1725857> since 2020-02-28T16:15:34Z

*Published version:*

DOI:10.1038/s41431-019-0551-x

*Terms of use:*

Open Access

Anyone can freely access the full text of works made available as "Open Access". Works made available under a Creative Commons license can be used according to the terms and conditions of said license. Use of all other works requires consent of the right holder (author or publisher) if not exempted from copyright protection by the applicable law.

(Article begins on next page)

1 **Title**

2 **A bird's eye view of Italian genomic variation and deleterious variants pattern**

3

4 **Cocca Massimiliano<sup>1</sup>, Barbieri Caterina<sup>2</sup>, Concas Maria Pina<sup>1,3</sup>, Gandin Ilaria<sup>1</sup>,**  
5 **Brumat Marco<sup>1,3</sup>, Robino Antonietta<sup>1</sup>, Trudu Matteo<sup>7</sup>, Sala Cinzia<sup>2</sup>, Vuckovic**  
6 **Dragana<sup>4</sup>, Girotto Giorgia<sup>1,3</sup>, Matullo Giuseppe<sup>5</sup>, Polasek Ozren<sup>6</sup>, Ivana Kolčić<sup>6</sup>,**  
7 **Paolo Gasparini<sup>1,3</sup>, Soranzo Nicole<sup>4</sup>, Toniolo Daniela<sup>2</sup>, Massimo Mezzavilla<sup>1</sup>**

8

9 **Affiliations**

10 **1) Institute for Maternal and Child Health IRCCS Burlo Garofolo, Trieste, Italy.**

11 **2) Division of Genetics and Cell Biology, San Raffaele Scientific Institute, Milan,**  
12 **Italy.**

13 **3) Department of Medical, Surgical and Health Sciences, University of Trieste,**  
14 **Trieste, Italy.**

15 **4) Sanger Institute, Wellcome Trust Genome Campus, Hinxton CB10 1SA, UK.**

16 **5) Department of Medical Sciences, University of Turin, Turin, Italy**

17 **6) Public Health, University of Split, Croatia**

18 **7) Molecular Genetics of Renal Disorders Unit, Division of Genetics and Cell**  
19 **Biology, San Raffaele Scientific Institute, Milan**

20

21

22 **Abstract**

23 The genomic variation in the Italian peninsula populations is currently under  
24 represented: the only Italian whole genome reference are the Tuscans from the 1000  
25 Genome Project. To address this issue, we sequenced a total of 947 Italian genomes  
26 from three different geographical areas that could be representative of a large portion of  
27 the whole country genomic pool. First, we defined a new Italian Genome Reference  
28 Panel (IGRP) for imputation, which showed high-performance, especially for rare  
29 variants imputation, and we subsequently validated it by GWAS analysis. Furthermore,  
30 we widened the catalogue of genetic variation and investigated population structure,  
31 pattern of natural selection, distribution of deleterious variants and human knockouts

32 (HKO). All the results emphasise a high level of genomic differentiation between  
33 populations, diverse signatures of natural selection and a distinctive distribution of  
34 deleterious variants and HKO, confirming the necessity of multiple genome references  
35 for the Italian population.

36

37

## 38 **Introduction**

39 Large sequencing projects have identified the majority of common variants and millions  
40 of rare and low frequency variants (Gudbjartsson et al., 2015; The 1000 Genomes  
41 Project Consortium, 2015; The UK10K Consortium, 2015). Most of the rare variants  
42 were detected in protein coding genes and it was calculated that each individual may  
43 carry more than 20.000 variants per exome (Karczewski et al., 2017; The ENCODE  
44 Project Consortium, 2007), a finding that complicates our understanding of gene  
45 function since only few genes may underline a disorder or be associated with a given  
46 phenotype. The filtering of candidate variants by frequency in unselected individuals is a  
47 key step in any pipeline for the discovery of causal variants. The efficacy of such filtering  
48 depends on both the size and the ancestral diversity of the available reference data.  
49 From this point of view, the catalogue of rare and low frequency variants is still largely  
50 incomplete, and its completion will represent a major challenge.

51 In the available human genome reference sequence data sets (i.e 1000G PH3, ExAC  
52 databases, etc.), Southern European populations, which represent a significant  
53 proportion of the overall European populations, are highly underrepresented (i.e. only a  
54 small group of subjects from Tuscany, Italy, and Spain). To fill this gap, we obtained  
55 whole genomes from founder populations - for which the presence of stratification (Esko  
56 et al., 2013; Sazzini et al., 2016) and the different level of isolation were demonstrated  
57 (Xue et al., 2017) - localized in three different parts of Italy: North-West (Val Borbera),  
58 North-East (Friuli Venezia Giulia) and South-East (Carlantino). In founder populations,  
59 variants that are rare or absent elsewhere can occur at higher frequencies and  
60 overcome the difficulty of identifying rare and low frequency variants. In this respect, our  
61 Italian genomes could be also extremely useful for the genetic analysis of other Italian  
62 and South European populations. An Italian Reference Genome panel for imputation

63 was also developed, tested and validated with GWAS analysis for red blood cells  
64 parameters and results were compared with those previously obtained using the 1000G  
65 data imputation panel only. Our work aims to answer the following questions: 1) *Are we*  
66 *able to increment the catalogue of genotypic variation, possibly in the low frequency*  
67 *spectrum, with new data?* 2) *Do we add useful information in terms of genetic variability,*  
68 *non-redundant with respect to the South European-Italian data already present in the*  
69 *commonly used reference panels?* 3) *Will we be able to identify new loci/variants,*  
70 *characteristic of a South-European subpopulation through GWAS using the new*  
71 *reference panel for imputation?* 4) *How much homogeneous are genomes coming from*  
72 *different regions of Italy in terms of population structure, natural selection signatures,*  
73 *deleterious variants distribution and human knockouts (HKO)? And, as a consequence,*  
74 *how reliable is to use only one reference population for Italians such as Tuscans?*

75

## 76 **Results**

77

### 78 **WGS data generation: variant calling and quality control**

79 A total of 947 DNAs from three cohorts were sequenced at 6 to 10X coverage; 381  
80 individuals from Friuli Venezia Giulia (FVG), 433 from Val Borbera (VBI) and 133 from  
81 Carlantino (CAR) (**Figure 1a**). Genotype calls for autosomal chromosomes were  
82 produced separately for each population. After filtering, 926 samples were retained.  
83 Approximately 27M sites (i.e. >24M SNVs and >2M indels) were detected (**Table 1**) in  
84 the joint dataset. Overall, 7.1 M sites (26%) were common (MAF>5%), 3.1M (12%) were  
85 low frequency (MAF between 1% and 5%) and 16.6M (62%) were rare (MAF <1%) with  
86 a similar partition in all cohorts. Singletons variants (AC=1) were >6M (24%) (**Table 1**  
87 **and Figure 1 b**). For each individual, we identified on average ~3.5M variant sites  
88 including ~0.56M indels and ~7.000 singletons. Considering each cohort separately, we  
89 noticed an excess of singletons in Carlantino cohort (CAR): most of them were shared  
90 with the other INGI cohorts, confirming that this is an artefact due to the lower sample  
91 size (124 samples, after QC). The comparison with outbred references (EUR subset  
92 from 1000G Phase 3, the whole 1000G Phase 3 and UK10K) highlighted that 34% to  
93 45% of the INGI variants are not represented (~12M with EUR, ~10M with 1000G and

94 ~9M with UK10K respectively): 89% of those variants are private to each INGI cohort.  
95 Moreover 8% of the sites shared between two or all three INGI cohorts were not found  
96 either in the whole 1000G or in the EUR subpopulation from 1000G (which includes  
97 Italian samples from the Tuscany region - TSI), suggesting that they may be  
98 characteristic of the general Italian population. The majority of the private variants are  
99 within the range of the low and rare frequencies (MAF < 1%) (**Figure 1c**) while the  
100 proportion of low frequency and common variants are similar in the pool of shared sites  
101 (**figure supplement 1, table supplement 1**).

102

103

#### 104 **IGRP1.0: Reference panel and imputation**

105 To increase the burden of good quality low frequency sites imputed in our isolated  
106 cohorts and possibly in the general Italian population, we generated a custom reference  
107 panel integrating our WGS data with already available resources from the 1000  
108 Genomes project.

109 Variants with read depth (DP) lower than 5 and all singleton variants not overlapping  
110 between all INGI populations or the 1000 Genome project data were excluded. After  
111 filtering, 95.6%, 94.29% and 92.06% variants were retained for CAR, FVG and VBI,  
112 respectively (**table supplement 2**). Merging our data with the 1000G Phase 3 reference  
113 resulted in the addition of 6.9M Italian population specific variants, 7.8% of the merged  
114 INGI+1000G (IGRP1.0, from now on) panel (**table supplement 3**).

115 We tested our resource on the INGI populations and on an outbred Italian cohort of  
116 randomly selected samples. As shown in **Figure 2**, the panel including our data (red  
117 line) always outperforms the 1000G phase 3 reference panel for the INGI cohorts in  
118 terms of genotype concordance ( $r^2$  - right y-axes), while there are not significant  
119 improvements for the outbred population (NW-ITA).

120 We compared our resource performances also in terms of the IMPUTE 'info score'  
121 metric. The proportion of well imputed sites (info score  $\geq 0.4$ ) in the IGRP1.0 reference  
122 panel was always higher compared to the 1000G phase 3 reference panel (red and blue  
123 bars respectively) with an increase from 20% to 36% of the rare sites (MAF<0.5%) with  
124 info score  $\geq 0.4$  (**Figure 2, table supplement 4**). Imputation of an outbred Italian

125 population showed a similar outcome: the variants added by our resource spread evenly  
126 across the info score bins without jeopardizing the imputation results. In particular, for  
127 the lowest frequency bin we could impute 800.721 sites with IGRP1.0 versus 698.140  
128 sites with 1000G phase 3 panel with info scores  $\geq 0.4$  and a 13% increase of good  
129 imputation of the rare sites. We further validated our resource on three Croatian cohorts  
130 (VIS, KORCULA, SPLIT): the IGRP1.0 panel has higher proportion of well imputed sites  
131 with respect to other panels with a result similar to the outbred Italian population (**figure**  
132 **supplement 2 and table supplement 5**). A direct comparison with the recent HRC  
133 reference panel (McCarthy et al., 2016) was not performed since our populations (as  
134 well as the 1000G samples) are included in that reference. However, we checked the  
135 quality of sites belonging to the INGI cohorts that are excluded because of filtering from  
136 the HRC reference: among seven test cohorts, we identified 696.895 to 624.434  
137 polymorphic sites with an average proportion of good quality sites (info score  $\geq 0.4$ ) of  
138 71% (63% - 81.5%). Focusing on rare variants for this subset, we can identify 256.222  
139 to 326.076 polymorphic sites with a proportion of good quality sites between 15 and  
140 63% (**table supplement 6**).

141

#### 142 **IGRP1.0: GWAS studies**

143 To assess the reliability of our new reference panel, we conducted a GWAS study with  
144 the newly imputed data on a series of red blood cells (RBC) traits (MCH, HGB, MCHC,  
145 RBC, HCT and MCV) for each INGI cohort followed by a meta-analysis. A total of 3292  
146 individuals (age  $\geq 18$  years) were included in the analysis. The characteristics of the  
147 samples are summarised in **table supplement 7**. Results from this analysis were  
148 compared with GWAS results for our cohorts with data imputed on the 1000G reference.  
149 Manhattan plots of all the meta-analysis results are given in **figure supplement 3**.

150 Lambda values of GWAS with 1000G showed no stratification (**figure supplement 4**).

151 Meta-analysis of GWAS with 1000G showed significant results ( $P < 6.23E-9$ ) only for  
152 MCH and MCV (**table supplement 8**). MCH analysis identified rs4820268 ( $P = 4.54E-$   
153  $10$ ) in TMPRSS3, a gene already associated to MCH, MCV and MCHC (Ferreira et al.,  
154 2009; Kullo et al., 2010). A locus on chromosome 11 at 3.8 Mb was identified for MCH  
155 and MCV both (rs117802349, MCH  $P = 2.67E-10$ , MCV  $P = 2.33E-11$ ) and two loci on

156 chromosome 11 at 5.2 and 5.4 Mb were identified for MCV, near beta-globin cluster  
157 (rs113853911, p-value = 4.01E-09 and rs80297185, p-value = 1.84E-14 - D' = 1,  
158 r<sup>2</sup>=0.328). Other significant results were obtained for low frequency variants on  
159 chromosome 2 (MCH, P= 9.44E-10), 5 (MCH P= 1.98E-10 and MCV P= 1.16E-09), and  
160 8 (MCV P=3.86E-10). Finally, a significant association was found for rs112483810  
161 which lies 3Mb close to rs7844723, already found in association with HGB (Yang et al.,  
162 2007).

163 As shown in **figure supplement 5**, lambda of meta-analysis of GWAS with IGRP1.0  
164 imputation were higher than lambda of 1000G imputation, due the high number of rare  
165 variants included in the new panel. However, the values ranged from 1.032 (MCHC) to  
166 1.0505 (RBC) indicating adequate control of population stratification. The meta-analysis  
167 of results of IGRP1.0 imputed data showed several GWAS significant results, mainly in  
168 low frequency and rare variants (**table supplement 9**). The best hits for HGB, MCH,  
169 MCV and RBC were found in HBB cluster (chr11p15.4). The IGRP1.0 meta-analysis for  
170 HGB, MCH, MCV and RBC identified the pathogenic SNP rs11549407 located in HBB  
171 gene and responsible of beta-thalassemia (MCV P=1.86E-59, MCH P=4.88E-52, RBC  
172 P=8.30E-14, HGB P=3.67E-10) (Danjou et al., 2015). This SNP was also replicated with  
173 higher p-values for MCHC (P=0.0001) and HCT (P=5.38E-07). This locus was found  
174 only in CAR and VBI and this rare variant (CAR MAF= 0.48% and VBI MAF= 0.28%)  
175 was present neither in FVG nor in 1000G EUR. Furthermore, this variant is at very low  
176 frequency in Exac European (AF=0.07%) and has a r<sup>2</sup> value of 0.328 (D'=1) with the  
177 rs113853911 variant identified in the previous analysis (replicated with IGRP1.0 only in  
178 HGB and HCT with p-values of 1.68E-4 and 2.49E-4 respectively).

179 IGRP1.0 meta-analysis confirmed TMPRSS6 gene for MCH and MCV already found in  
180 1000G analysis (i.e. significant only in MCH).

181 Overall, IGRP1.0 imputation panel allowed us to replicate known loci and loci identified  
182 through the 1000G imputation, increasing also the number of significant variants, as  
183 shown in **Figure 3 a-b**.

184

185

186 **Population structure**

187 Using only European populations for PCA analysis, each INGI population separates  
188 from each other in the first four principal components (**Figure 4 a and Figure 4 b**).  
189 Regarding the FVG cohort, we can appreciate the separation of the six villages included  
190 in the isolate: Erto (ERT), Illegio (ILG), Resia (RSI), Sauris (SAU), San Martino del  
191 Carso (SMC) and Clauzetto (CLZ), underlining the evidence of population structure and  
192 of a high degree of isolation, as shown previously (Xue et al., 2017). Analyses and  
193 clustering using genomic pairwise  $F_{st}$  (**Figure 4 c**) highlight how INGI populations  
194 cluster with Europeans. However, the six villages from FVG show high levels of  $F_{st}$  in  
195 respect to other Italians. A closer look was taken with Treemix (Pickrell and Pritchard,  
196 2012) analyses (**Figure 4 d**). Different lengths in the tree due to both inbreeding and  
197 genetic drift confirmed the peculiar structure of the FVG cohorts but, most importantly, it  
198 showed gene flow between North European population and North Eastern Italians. This  
199 adds more complexity to the Italian genomic pool.

200 Admixture (Alexander et al., 2009) analyses at different cluster solutions from  $K=2$  to  
201  $K=14$  were also performed using worldwide reference populations from 1000G. The  
202 cluster solution with lowest cross validation error was for  $K=9$  (**figure supplement 6**).  
203 VBI showed an admixture pattern similar to the one of Tuscany (TSI) from 1000G. The  
204 more isolated FVG populations showed their own ancestral component, that was  
205 however present at different fractions in all European and Italian populations suggesting  
206 a strong differential isolation of Italian subpopulations.

207 Finally, inbreeding coefficients and total homozygosity (due to ROH) showed high levels  
208 of variance among different Italian subpopulations as shown by the shape of the  
209 beanplots (**Figure 4 e-f**). In particular, VBI shows the lowest mean coefficient  
210 (mean=0.008) CAR, CLZ and SMC had similar distributions (0.0149, 0.0134, 0.0151,  
211 respectively). Inbreeding was particularly high for ERT, SAU, ILG, and RSI (0.0191,  
212 0.0325, 0.0304, 0.0311, respectively), the same pattern is followed by the total  
213 homozygosity due to ROH, which is quite different from the reference Italian population.

214

## 215 **Natural Selection**

216 We tested natural selection using the statistics  $iHS$  (Voight et al., 2006), we grouped the  
217 markers accordingly: markers with  $|iHS| \geq 2$  in only one population and markers with



218 |iHS|>=2 in all Italian populations. We applied the following stringent criteria to select  
219 genes with signature of positive selection: at least 20 markers with |iHS|>=2.

220 A total of 37 other genes was found under putative selection in all Italian populations  
221 (**Supplementary table 10**). Interestingly, six genes (FHIT, CSMD1, CNTNAP2,  
222 MACROD2, RBFOX1 and PTPRD) were found under putative selection in all Italian  
223 populations but with different markers. Some of them had been previously associated  
224 with complex traits such as FHIT associated with BMI (Hoffmann et al., 2018), CSMD1  
225 associated with 79 different phenotypes, including age of menarche (Perry et al., 2014),  
226 schizophrenia (Bergen et al., 2012) and educational attainment (Lee et al., 2018) (data  
227 from GWAS catalogue), CNTNAP2 associated with mathematical ability (Lee et al.,  
228 2018) and DNA methylation variation (Zhang et al., 2018), MACROD2 associated with  
229 several phenotypes, including educational attainment (Lee et al., 2018) and blood  
230 protein levels (Sun et al., 2018), RBFOX1 associated with eyes (Pickrell et al., 2016),  
231 other neurological traits and also educational attainment (Lee et al., 2018) and PTPRD  
232 associated with restless leg syndrome (Schormair et al., 2017) and blood pressure  
233 (Evangelou et al., 2018).

234 As shown in **Figure 5**, the majority of genes with signatures of selection are not found in  
235 the TSI (the available Italian reference population from 1000G project). More in detail  
236 the fraction of private genes under selection ranges from 74% in VBI to 86% in RSI.

237 As a further example of the complex puzzle of signature of selection present in the  
238 Italian peninsula we selected the highest-ranking genes in terms of |iHS| and number of  
239 SNPs with |iHS|>=2 that are found only in one population. We then provided some  
240 examples reporting the genes with the average highest |iHS| and highest number of  
241 SNPs with |iHS|>=2 that are found only in one population.

242 Starting from the current Italian reference TSI, we found a strong signal for TYW1B,  
243 associated with triglycerides (Teslovich et al., 2010) and educational attainment. We  
244 found signature in CYP2C19 in CAR (associated with diastolic blood pressure (Liu et al.,  
245 2016), ABCG8 in VBI associated with lipid traits (Chasman et al., 2009), SLC25A12 in  
246 SMC (Educational attainment (Lee et al., 2018)), ERI3 in CLZ (associated with  
247 Educational attainment (Lee et al., 2018)), in ERT we found strong signatures for  
248 ANKRD30A, which was associated with paediatric autoimmune diseases,metabolite

249 levels and vestibular neuritis (Li et al., 2015), in SAU evidences were found for CLOCK  
250 gene associated with height (Wood et al., 2014), we found SSPN in ILG (associated  
251 with atrial fibrillation,(Nielsen et al., 2018)) and finally PBRM1 in RSI which was linked to  
252 blood protein levels (Sun et al., 2018),schizophrenia,general cognitive ability (Davies et  
253 al., 2018, p. 4). These are few examples of the different genomic patterns that can be  
254 found in the various subpopulations of the Italian peninsula.

255

### 256 **Deleterious variants enrichment**

257 The profound differentiation in all our Italian samples, and subsequent different level of  
258 isolation and selection signature led to the question whether there was any difference in  
259 deleterious or neutral variant distribution among different populations compared to the  
260 Italian reference population.

261 To answer this question, we applied the DVxy statistic (Xue et al., 2017) for DV variants  
262 (Drifted Variants respect to a reference) between 1-2 allele count (AC) and 3-5 AC in  
263 each population using as actual Italian reference (TSI). Variants were grouped  
264 according to CADD score. In our analysis we discovered a significant relative  
265 enrichment in deleterious variants with  $CADD > 20$  in the more isolated/inbred  
266 populations compared to the TSI ( $DV_{xy-sd} > 1$ ), this is true when we are considering  
267 populations such as ILG, RSI, SAU and also SMC, whereas no differences were found  
268 when considering neutral or low deleterious variants ( $CADD 0-5$ ,  $DV_{xy+/-sd} = 1$ ) (**Figure**  
269 **6**).

270 This level of enrichment could be explained with lower effectiveness of purifying  
271 selection due to isolation and small effective population size (Xue et al., 2017), however  
272 the interesting point lies somewhere else: these Italian populations show enrichment for  
273 high deleterious variants ( $CADD > 20$ ) at low frequency (3-5 AC).

274 In order to show the complexity of the Italian catalogue of deleterious variants, we  
275 estimated the ratio of DV variants (3-5 AC and  $CADD > 20$ ) that are different between  
276 pairs of sub-populations with DV variants shared between pairs (**Figure 7**), with the  
277 lowest value being 12 for the pair CLZ/VBI and the highest 31 for RSI/SAU. All values  
278 are highly positive indicating that the majority of DV variants are private of each  
279 subpopulations.

280

## 281 **Human Knockout**

282 Homozygote loss of Function (LoF) variants represent a category of deleterious variants  
283 and examples of human KO (HKO). Considering that the highest enrichment was found  
284 for variants with CADD score >20, we used this value to select LoF variation.

285 In our total cohort, 509 LoF presenting with a CADD score >20 were found at  
286 homozygous state in at least one individual per population (**table supplement**  
287 **11**). Gene ontology analysis revealed an excess of transmembrane signalling receptor  
288 genes including olfactory receptors, as already described (MacArthur et al., 2012).

289 In order to have an high reliable dataset, we used a stringent filtering criteria analysing  
290 only variants that affected all transcripts (considered as TOTAL LoF in opposition to  
291 PARTIAL LoF), resulting in 205 variants affecting 195 different genes (**table**  
292 **supplement 12**). Among these 205 variants, the majority (150, ~73%) was shared  
293 among all 3 populations and more than a half (~60%) had frequency  $\geq 0.05$ . A large  
294 number of HKOs was located in genes involved in hair/skin/epithelium or eye  
295 phenotypes, and many were members of gene families. As a matter of facts, 5 HKO  
296 were found in keratin genes: (KRT37, KRT24, KRT31 and KRT83) and 5 in keratin  
297 associated protein (KRTAP1-5, KRTAP1-1, KRTAP19-6 and KRTAP13-2, KRTAP29-1).  
298 Two different rare stop gain mutations were found in KRT83 (AF in Europeans  
299  $rs146753414=0.0265$ ,  $rs2857667=0.0063$ ). Two missense mutations in this gene were  
300 associated to a mild form of monilethrix (MNLIX; OMIM #158000), a rare autosomal  
301 dominant hair disease that results in fragile, brittle hair that tends to fracture and  
302 produce some degree of alopecia (Steensel et al., 2005). The hair of three carriers of  
303 the KO of KRT83 in the VB population and of nine heterozygotes in three families was  
304 investigated and resulted normal. Lack of KRT83 does not seem to affect hair structure  
305 as much as substitutions of amino acids that are highly conserved and affect the helix  
306 termination motifs, known hotspots for monilethrix mutations. Finally, we found a very  
307 rare stop gain ( $rs11355796$ ) in COL6A5 gene (collagen type VI, alpha 5) identified in  
308 homozygous state in one individual from VBI. The variant is enriched in all three Italian  
309 populations (AF~0.015 compared with reference Europeans of 0.0013. Mutations in the  
310 COL6A5 gene were shown to cause familial neuropathic chronic itch (Martinelli-

311 Boneschi et al., n.d., p. 5). Unfortunately, we do not have so far clinical data on our  
312 homozygous HKO, but for sure a related heterozygous carrier reported to complain of  
313 itching all his life.

314 We then analysed only variants reported in gnomAD, resulting in 133 different genes  
315 which are distributed among the populations as shown in **Figure 8**: we found that the  
316 majority of genes are private of FVG, VBI and CAR (61, 36 and 10 respectively)  
317 whereas only 13 genes are shared among all populations. Among these HKO genes  
318 only few of them show evidence of selection (11 genes out of 133) (see **table**  
319 **supplement 13**), in particular signatures of selection are found in populations where the  
320 LOF variants are not present with the exception of CDH23 (associated with LDL  
321 cholesterol) and KLHL23 (associated with obesity-related traits); however, these HKO  
322 are considered PARTIAL.

323 The novel aspect that we report is that we add new information into the variability and  
324 genomic signatures in HKO genes present in the Italian genomic pool.

325

## 326 **Discussion**

327 The ability to interrogate all kind of genetic variations is critical for the classification of  
328 genetic determinants of complex and monogenic disorders: the whole genome  
329 sequencing of peculiar populations such as isolates has given a significant contribution,  
330 providing denser data and allowing a better mapping of the genomic features under  
331 study (Hatzikotoulas et al., 2014).

332 Here, we report the results of a series of analyses obtained through the investigation of  
333 WGS from 947 subjects coming from different Italian geographic areas (i.e. South, North  
334 West and North East) and their contribution to the identification and description of a  
335 relevant proportion of the Italian population pool of genetic variation.

336 The number of new variants described and discovered, especially in comparison with  
337 1000G data, which include the Italian reference TSI (~1.86 M variants shared between  
338 all INGI cohorts but not TSI) confirms that these genomes are able to increment the  
339 catalogue of Italian genotypic variation, in particular in the low frequency spectrum.

340 This leads us to the inevitable next step: the creation of a reference panel for imputation  
341 using the Italian whole genome data. The “Italian core” of the reference panel was

342 assembled with INGI data only and it proved to be an extremely useful resource when  
343 merged with other larger reference panels: the IGRP1.0 outperformed the 1000G Phase  
344 3 reference panel for imputation of inbred and outbred Italian and other European  
345 populations such as the Croatians cohorts.

346 More in detail, our new reference panel could facilitate the imputation of rare variants for  
347 GWAS studies and help the identification of population specific variants of different  
348 Italian and possibly Southern European populations: a notable point is that we are  
349 incrementing the total number of variants that are valuable for GWAS studies without  
350 adding “noise” neither INGI populations, as expected, nor in other outbred populations in  
351 terms of imputation quality.

352 With this resource at our disposal, another question arises: will we be able to increment  
353 the power to detect genome wide significant loci/variants using this new reference panel  
354 for imputation?

355 In this case, the reliability of IGRP1.0 panel was proven running a series of GWAS tests  
356 on some selected RBC traits, demonstrating that it performs better than the 1000G  
357 panel alone. As a matter of facts, GWAS studies carried out with IGRP1.0 panel  
358 imputed data, not only replicated previous findings with higher statistical significance,  
359 but also demonstrated that several previously found suggestive signals ( $p < 1E-5$ )  
360 became genome wide significant ( $p < 1E-8$ ).

361 One interesting example is the RBFOX1 gene: we showed that it harbours signals of  
362 selection in all Italian populations and, moreover, it carries two variants significantly  
363 associated to MCV and MCH traits, that we were able to pinpoint only through our  
364 custom reference panel (**table supplement 11**). These results need further dissection,  
365 but are, again, a clue of the fact that the genetic features of our cohorts represent an  
366 important resource in the understanding of gene function and association to different  
367 traits.

368 On the other hand, one of the major point of our work consisted in addressing the issue  
369 of the underrepresentation of South European population in whole genome databases.  
370 Recent works based on array data pinpointed the genetic diversity in the Italian  
371 peninsula (Sazzini et al., 2016) along with the presence of isolates (Esko et al., 2013).  
372 This foregoing information prompted the inquiry about the homogeneity of genomes

373 coming from different regions of Italy in terms of diverse genomic aspects (population  
374 structure, natural selection signatures, deleterious variants distribution and HKO) and,  
375 as a consequence, how reliable is to use only one reference population for Italians such  
376 as the Tuscan (TSI).

377 Population structure analysis revealed that our populations fall in the European pole of  
378 variation, but their separation from the North European cohorts is clear.

379 Principal component analysis, tree graph analyses, ancestry coefficient distribution  
380 confirm the non-homogeneous genetic background of the Italian populations from North  
381 to South and highlight the fact that the use of the only TSI as genome reference leads to  
382 an underestimation of the Italian genomic variability, and if that was not enough, ROH  
383 pattern and inbreeding coefficient showed a wide array of values not comparable using  
384 the only South European reference of 1000 Genomes.

385 Discussing Natural Selection, it was demonstrated that environmental differences along  
386 the peninsula might have shaped the genome through mechanisms such as evolution  
387 and selective pressure (Sazzini et al., 2016). Our analyses pinpointed the presence of  
388 shared selective pressure on specific genes in all Italian populations such as HIT,  
389 CSMD1, CNTNAP2, MACROD2, RBFOX1 and PTPRD, however the striking point was  
390 the level of selection signatures that are private of single populations (when substructure  
391 is taken into account) ranging up to 86% of the total genes found for RSI. In addition,  
392 considering the relationship of some populations (RSI, SAU, SMC) with North European  
393 populations (as shown in Treemix analyses), we can suppose that a number of  
394 haplotypes passed in some North East Italian populations but not others: this peculiar  
395 gene flow could be responsible for some unique signals of selection.

396 For what concerns the distribution of deleterious variants it was already demonstrated  
397 how the relative relaxation of purifying selection in presence of isolation (Chheda et al.,  
398 2017; Xue et al., 2017) leads to an increased frequency of specific deleterious variants.

399 This aspect reinforces our thesis about the need of a more broadened reference for the  
400 Italian genomic variation, as we demonstrated that not only have we an enrichment of  
401 low frequency deleterious variants ( $CADD \geq 20$ ) in our genomes, but also most of this  
402 enrichment is population-specific.

403 In our analyses of human knockout we showed that the majority of the loss of function

404 (>70%) was shared among all the three populations and, as expected, they belong to  
405 the category of transmembrane signalling receptor genes, including olfactory receptors.  
406 Nevertheless, while analysing only the genes with at least 1 homozygous individual in  
407 each cohort, we discovered an inverse pattern: the majority of genes harbouring HKO  
408 are private of each cohort. In addition the majority of them (91%) was not found in any  
409 selection scan, suggesting the lack of evolutionary constraints for these genes, even  
410 though more accurate analyses are needed to confirm this consideration. Still, this gives  
411 us another hint of the necessity of multiple genomes to describe the catalogue of HKO  
412 present in Italy: with our data we are starting to scratch the surface, providing some of  
413 them.

414 Furthermore, HKO and pattern of deleterious variants are useful examples to show how  
415 clinical-relevant polymorphisms could be found enriched in frequency in specific  
416 populations within the same country and provide extremely relevant information that  
417 could be used for developing personalized medicine strategies: another great added  
418 value of our cohorts is that a large series of instrumental and clinical phenotypes is  
419 already available. Thus, future efforts should be pointed towards the functional  
420 characterization of putatively enriched variants and deep phenotyping of carriers of such  
421 polymorphisms as well as deep phenotyping of HKO.

422 In conclusion, we showed how our unique dataset of populations and WGS data fill a  
423 gap in publicly available human genome sequence data sets (i.e. 1000G, gnomAD  
424 databases, etc.), in which Southern European populations - a significant proportion of  
425 the overall European populations - are highly underrepresented and it will be able to  
426 produce regionally appropriate reference panels.

427 Furthermore, considering the fact that in Italy, a National Genomic BioBank is not yet  
428 existing the availability of a catalogue of rare and low frequency variants for Italians  
429 populations will facilitate the understanding of these genetic loci improving the accuracy  
430 and efficacy of a series of genetics/genomics studies, and subsequently opening new  
431 perspectives for precise medicine and drug targets identification.

432

## 433 **Materials and Methods**

434

435 **WGS data generation: variant calling and quality control**

436

437 All samples selected for sequencing already had genotype data from other platforms  
438 (SNP array and Exome chip): this data allowed us to assess genotype concordance  
439 against a “trusted” set of variants. For all cohorts, samples were selected randomly. The  
440 sequencing was carried out at different sequencing centres: the Wellcome Trust Sanger  
441 Institute in Hinxton (UK), the BGI, Shenzhen (PRC) and the HSR in Milan. **Table**  
442 **supplement 14** summarises the total number of samples sequenced. All the data were  
443 post processed at the Sanger Institute. A written informed consent for participation was  
444 obtained from all subjects. Regarding the FVG cohort the project was approved by the  
445 Ethical committee of the IRCCS Burlo-Garofolo. Regarding the CAR cohort the project  
446 was approved by the local administration of Carlintino, the Health Service of Foggia  
447 Province, Italy, and ethical committee of the IRCCS Burlo-Garofolo of Trieste. For VBI  
448 cohort, data collection and genotyping were approved by the institutional ethical  
449 committee of the San Raffaele Hospital in Milan and by the Regione Piemonte.

450 The raw data were checked first at lane level to remove any sample with bad quality  
451 data. 54 samples were realigned to the hs37d5 reference sequence because they were  
452 aligned to a previous version of the GRCh37 build: this process has been carried out  
453 using the ‘Bridgebuilder system’ developed by the Human Genetics Informatic group at  
454 the Wellcome Trust Sanger Institute (*BridgeBuilder*, 2015) .

455 After the alignment, performed with bwa software (Li and Durbin, 2010), each bam file  
456 was improved through the implementation of the following steps: 1) Realignment  
457 around known and discovered INDELS using GATK (McKenna et al., 2010)  
458 RealignerTargetCreator and IndelRealigner; 2) Base Quality Recalibration by GATK  
459 BQSR using the BaseRecalibrator and PrintReads tools; 3) Recalculation of the MD tag  
460 by samtools (Li, 2011) calmd; 4) Bam indexing.

461 For the Carlintino cohort, sequencing was carried out using Illumina technology  
462 (Genome Analyzer and HiSeq 2000) at the Wellcome Trust Sanger Institute for 115  
463 samples with an average coverage of 4x, an additional batch of 40 samples was  
464 sequenced at Beijing Genomics Institute (BGI) with an average coverage of 10x. Among  
465 the 115 samples sequenced at the Sanger Institute 27 failed the quality check at the



466 lane level: 5 were re-processed while 22 were excluded from further analyses.  
467 The most common cause of failure was the high percentage of adapter contamination  
468 and a bimodal insert size distribution.  
469 For the Friuli Venezia Giulia cohort, 200 samples were sequenced at the Wellcome  
470 Trust Sanger Institute with a mean coverage of 4x and 192 samples at BGI with a mean  
471 coverage of 10x. Among the 200 Sanger samples, only 4 failed the quality check at the  
472 lane level, thus were excluded from further analyses. Among the BGI set 6 samples  
473 were duplicated from the Sanger pool: we merged the two sets of data to increase the  
474 coverage of each sample. We removed 1 additional sample from this set because of  
475 data corruption.  
476 The data for the Val Borbera Cohort were generated at a mean coverage of 6x for all  
477 selected samples: 210 were sequenced at the Wellcome Trust Sanger Institute, 209  
478 were sequenced at BGI and a small batch of 29 was processed at the San Raffaele  
479 Hospital. After the first step of quality check we removed 2 samples from the Sanger  
480 Institute set for contamination and bad quality DNA respectively, 12 samples from the  
481 OSR dataset for bad quality and 1 sample from the BGI set for data corruption.  
482 Finally a total set of 947 samples was sent forward for the Variant Calling step.  
483 We produced genotype calls for autosomal chromosomes separately for each  
484 population using the following pipeline.  
485 Samtools mpileup (v.1.2) (Li, 2011) was used for multisample genotype calling  
486 (parameter set: -E -t DP,DV,SP -C50 -pm3 -F0.2 -d 10000). The generated BCF files  
487 were converted to VCF format with bcftools call (v.1.2) (parameter set: -Nvm) and  
488 filtered with bcftools filter (v.1.2) (parameter set: -m+ -sLowQual -e"%QUAL<=0"-g3 -  
489 G10 -Ov - ). Variant Quality Score Recalibrator (VQSR) filtering was applied to the raw  
490 call data with GATK v.3.3 (DePristo et al., 2011). Raw calls from samtools were used  
491 with the UnifiedGenotyper module in "Given allele mode" to generate all the annotation  
492 needed to calculate the VQSLOD scores through the VariantRecalibrator module,  
493 separately for SNVs and INDELs. For SNVs we selected the following parameters: i)  
494 Annotations: QD, DP, FS, HaplotypeScore, MQRankSum, ReadPosRankSum,  
495 InbreedingCoeff; ii) Training set: HapMap 3.3, Omni 2.5M chip, 1000 Genomes Phase I;  
496 iii) Truth set: HapMap 3.3, Omni 2.5M chip; iv) Known set: dbSNP build 138. For

497 INDELS we selected: i) Annotations: DP, FS, ReadPosRankSum, MQRankSum; ii)  
498 Training set: Mills-Devine, 1000 Genomes Phase I, dbSNP v138; iii) Truth set: Mills-  
499 Devine; iv) Known set: Mills-Devine, dbSNP build 138. For each population the lowest  
500 VQSLOD threshold was chosen according to the output produced by  
501 VariantRecalibrator to select the best cut-off in terms of specificity and sensitivity of the  
502 trained model. The Transition/Transversion (Ti/Tv) ratio was used as a parameter to  
503 select the best threshold, taking as a reference the empirical value of  $\sim 2$  calculated by  
504 (1000 Genomes Project Consortium et al., 2012) . For SNPs the minimum VQSLOD  
505 values selected were -59.1994 (99.94% truth sensitivity threshold), -15.0283 (99.80%  
506 truth sensitivity threshold), -22.6034 (99.9% truth sensitivity threshold) for VBI, FVG and  
507 CAR cohort respectively. For INDELS we used a more conservative approach, selecting  
508 a sensitivity threshold of 95% for each population. The filter was applied to each call set  
509 with GATK ApplyRecalibration module.

510 We performed several genotype refinement steps on the filtered data: 1.  
511 BEAGLEv4.r1230 (Browning and Browning, 2007) was used to assign posterior  
512 probabilities to all remaining genotypes. 2. SHAPEITv2 (Delaneau et al., 2013) to phase  
513 all genotypes calls and 3. IMPUTEv2 (Howie et al., 2009) to perform internal imputation  
514 in order to correct genotyping errors.

515 Finally, bcftools annotate (v.1.2) was used to add information about Ancestral Allele and  
516 allele frequencies from 1000G phase 3 (Sudmant et al., 2015) populations and rsIDs  
517 from dbSNP v.141 (Sherry et al., 2001). The Variant Effect Predictor v.90 (McLaren et  
518 al., 2010) provided all consequence annotation as well as Polyphen and Sift information.  
519 CADD score (Kircher et al., 2014) information was also added.

520 Samples and sites were again investigated for outliers or artefacts after the variant  
521 calling.

522 First, we looked for batch effects due to the different sequencing centres: we conducted  
523 an MDS analysis on each cohort testing the first PCA component for correlation with the  
524 sequencing centre variable with a Pearson's correlation test and obtaining a significant  
525 outcome only for the FVG cohort ( $p=0.001728$ ). We compared the analysis for the FVG  
526 cohort with data available from a previous work (Esko et al., 2013) showing that the  
527 pattern is consistent with the underlying population structure. We then generated a sites

528 exclusion list, focusing on: a) Hardy-Weinberg equilibrium (sites removed if exact test p-  
529 value was below the threshold of  $1e-8$ ); b) Heterozygosity rate distribution (removed  
530 sites with values greater than 3 standard deviations of the mean); c) MAF mismatch  
531 when compared with SNP array data; d) Non Reference Discordance rate (NRDR),  
532 defined as the ratio between the sum of concordant calls of the alternative allele in WGS  
533 and array data and the sum of all discordant calls of the alternative allele in WGS and  
534 array data (cut off value for removal of 3 standard deviations of the mean).

535 We removed 5,552 , 2,577 and 2,502 sites from CAR, FVG and VBI respectively.

536 We excluded samples using the following parameters: a) Singleton number, b)  
537 Heterozygosity rate and c) Non Reference Discordance rate. We removed one sample  
538 from the FVG cohort for an excess of singletons (~100000 singletons counted). We  
539 calculated also the heterozygosity rate for each sample and removed all samples with  
540 values exceeding a threshold of 3 SD from the average value for each population: one  
541 sample was removed from the CAR cohort, one sample from the FVG cohort and 4  
542 samples from the VBI cohort. Finally, we calculate the samples' non-reference  
543 discordance rate and removed all individuals with an NRDR greater than 5%: 8 samples  
544 from the CAR cohort, 1 sample from FVG cohort and 5 samples from the VBI cohort.

545

#### 546 **Reference imputation panel**

547 We selected a 'highly reliable' subset of variants to include in our reference panel.

548 In order to avoid mismatches between the INGI datasets, we split all multi-allelic variant  
549 sites in different vcf records and performed INDELs normalization with bcftools norm to  
550 prepare the data. Data from 1000G Project phase 3 and UK10K project (The UK10K  
551 Consortium, 2015) were processed in the same way.

552 SNPs and INDELs from the INGI WGS data to be included in the reference panel were  
553 selected with the following criteria: a) all sites with Alternative Allele count (AC)  $\geq 2$  and  
554 Read depth (DP)  $\geq 5$ ; b ) all the singleton sites (AC = 1) either shared at least between  
555 two INGI cohorts or which were known sites or present at least in one of the external  
556 resources selected (UK10K and 1000G Project Phase 3).

557 To build the Italian reference dataset, a 'core' INGI panel was created merging data  
558 from the different INGI cohorts, using the method implemented by the IMPUTE2

559 software (Howie et al., 2011). The data were then added to the 1000G phase 3  
560 reference panel to obtain a final reference (INGI+1000G also called IGRP1.0).

561 The imputation test was performed on chromosome 2 genotypes in different cohorts: a)  
562 INGI cohorts; b) a cohort of 567 unselected outbred samples from North Western Italy  
563 (NW-ITALY); c) three cohorts from Croatia (VIS - 960 samples, KORCULA - 1812  
564 samples and SPLIT - 466 samples)

565 Imputation metrics across the different panels were compared for each population. We  
566 assessed the  $r^2$  metric, which estimates the correlation between the true genotype and  
567 the imputed genotype and the IMPUTE **info score** parameter, which provides a  
568 measure of the observed statistical information associated with the allele frequency  
569 estimate for each variant (Marchini and Howie, 2010). We removed from each INGI  
570 cohort all the samples represented in the reference panel.

571

## 572 **Genome Wide Association Studies (GWAS)**

573 GWA studies on Red Blood Cells indexes (MCH, HGB, MCHC, RBC - normalized with  
574 natural logarithm , HCT and MCV) were performed in each population separately, using  
575 age and gender as covariate in an additive model, once using 1000G imputation and  
576 once IGRP1.0. The analyses were carried out using the mixed linear models as  
577 implemented in R ABEL packages (Aulchenko et al., 2007). Genomic kinship was used  
578 to take into account the relatedness. Variants with info score  $\leq 0.4$  were excluded if the  
579 MAF was  $\geq 1\%$ . For rare variants (MAF 0.1%-1%), a more stringent Info Score cut-off  
580 ( $\geq 0.8$ ) was used (Pistis et al., 2015). Meta-analysis was performed using the software  
581 METAL (Willer et al., 2010) and heterogeneity Cochran Q test was performed. After  
582 meta-analysis, the variants that were not present with the same direction in at least two  
583 of the three cohorts were excluded. Variants with significant p-value ( $< 0.05$ ) for  
584 heterogeneity test were also excluded. Bonferroni correction was applied: the thresholds  
585 were  $P=6.23e-9$  for 1000G and  $4.69e-9$  for IGRP1.0. The positions are referred to the  
586 build 37. Manhattan plots were generated with the R library qqman (Turner, 2017) and  
587 hudson package (Lucas, 2018) .

588

## 589 **Population Structure**

590 Principal component analysis (PCA) was carried out to define the genetic structure of  
591 our population using PLINK (Daly et al., 2007). PCA was carried out after removing  
592 markers in high LD ( $r^2 > 0.4$ ), using the function `--indep-pairwise 200 50 0.4` and with  
593 MAF  $< 0.02$ . Runs of homozygosity (ROH) and inbreeding coefficient we estimated as  
594 well using PLINK using the command `--homozyg` and `--het`. Pairwise  $F_{st}$  between  
595 worldwide populations was calculated using the software 4p (Benazzo et al., 2015).  
596 The same dataset was used for tree graph analyses implemented in Treemix (Pickrell  
597 and Pritchard, 2012). The analysis of ancestral component was performed using  
598 ADMIXTURE v 1.2 (Alexander et al., 2009) using the European population plus one  
599 African reference (YRI) one East Asian (CHB) and one South Asian (GIH). Cross  
600 validation error procedure was implemented to select the best cluster solution.

601

### 602 **Natural Selection**

603 Evidence of positive selection was estimated for each population using iHS statistic  
604 (Voight et al., 2006) implemented in selscan program (Szpiech and Hernandez, 2014),  
605 we used only markers with MAF  $> 0.05$ , furthermore we adopted a conservative approach  
606 for genes under putative positive selection: we selected only genes with at least 20  
607 markers with standardized  $|iHS| \geq 2$ .

608

### 609 **Deleterious Variants**

610 After the exclusion of multiallelic variants, we subdivided all variant in bins according to  
611 their CADD score and frequency. The following minor allele frequency classes were  
612 created: between 1-2 allele count, 3-5-allele count, 5-10AC and more 10 AC, thus the  
613 variants were binned in the following CADD categories 0-5, 5-15, 15-20  $> 20$ . We then  
614 applied the DVxy statistic as described in Xue et al., using as reference the TSI  
615 population from 1000 Genomes. In addition, we estimated the ratio of of private and  
616 shared DV variants (variants enriched).

617

### 618 **Human Knockouts**

619 To identify HKO, we considered only deleterious variants in protein coding genes: we  
620 first selected variants with high impact as defined by VEP (i.e. frameshift, splice

621 acceptor variant, splice donor variant, stop gained, stop lost, start lost, transcript  
622 ablation, transcript amplification) and among those we further selected for CADD  
623 score  $\geq 20$ . A total of 12,231 variants (8,832 SNV and 3,399 indels) were selected and  
624 5,916 had a CADD score  $\geq 20$ . Among this subset of variants, those presenting at least  
625 one homozygous individual in one population were defined putative HKO. After filtration  
626 for total KO, the average number of HKO per individual was 20 (12-31), in agreement  
627 with previous determinations (Narasimhan et al., 2016). HKO's were classified as  
628 TOTAL when the variant was predicted as LOF in all Ensembl database transcript,  
629 otherwise they were classified as PARTIAL, even though this approach is highly  
630 conservative, as some PARTIAL loci could still affect the functional transcripts. Overlaps  
631 of HKOs between populations were analysed using the R package "VennDiagram"  
632 (Chen, 2018) (<https://cran.r-project.org/package=VennDiagram> ).

633

634

### 635 **Acknowledgements**

636 We would like to thank the people of the Friuli Venezia Giulia Region and of Carlantino  
637 for the everlasting support. We thank the inhabitants of the Val Borbera that made this  
638 study possible, the local administrations, the Tortona and Genova archdiocese and the  
639 ASL-22, Novi Ligure (AI) for support. We also thank Clara Camaschella for data  
640 collection supervision and organization of the clinical data collection, Fiammetta Viganò  
641 for technical help, Corrado Masciullo for building the analysis platform.

642

643

### 644 **Funding**

645 Fort FVG and CAR cohorts: Project co-financed by the European Regional  
646 Development Fund under the Regional Operational Programme of Friuli Venezia Giulia -  
647 Objective "Regional Competitiveness and Employment" 2007/2013, Telethon  
648 Foundation (GGP09037), Fondo Trieste (2008), Regione FVG (L.26.2008), and Italian  
649 Ministry of Health (RC16/06, ART. 13 D.LGS 297/99) (to PG). For VBI cohort: The  
650 research was supported by funds from Compagnia di San Paolo, Torino, Italy;  
651 Fondazione Cariplo, Italy and Ministry of Health, Ricerca Finalizzata 2008 and CCM

652 2010, and Telethon, Italy to DT. The funders had no role in study design, data collection  
653 and analysis, decision to publish, or preparation of the manuscript.

654

655

#### 656 **Competing interests**

657 No competing interests declared

658

#### 659 **References**

660

661 1000 Genomes Project Consortium, Abecasis GR, Auton A, Brooks LD, DePristo MA,  
662 Durbin RM, Handsaker RE, Kang HM, Marth GT, McVean GA. 2012. An  
663 integrated map of genetic variation from 1,092 human genomes. *Nature* **491**:56–  
664 65. doi:10.1038/nature11632

665 Alexander DH, Novembre J, Lange K. 2009. Fast model-based estimation of ancestry in  
666 unrelated individuals. *Genome Res* **19**:1655–1664. doi:10.1101/gr.094052.109

667 Aulchenko YS, Ripke S, Isaacs A, van Duijn CM. 2007. GenABEL: an R library for  
668 genome-wide association analysis. *Bioinforma Oxf Engl* **23**:1294–1296.  
669 doi:10.1093/bioinformatics/btm108

670 Benazzo A, Panziera A, Bertorelle G. 2015. 4P: fast computing of population genetics  
671 statistics from large DNA polymorphism panels. *Ecol Evol* **5**:172–175.  
672 doi:10.1002/ece3.1261

673 Bergen SE, O'Dushlaine CT, Ripke S, Lee PH, Ruderfer DM, Akterin S, Moran JL,  
674 Chambert KD, Handsaker RE, Backlund L, Ösby U, McCarroll S, Landen M,  
675 Scolnick EM, Magnusson PKE, Lichtenstein P, Hultman CM, Purcell SM, Sklar P,  
676 Sullivan PF. 2012. Genome-wide association study in a Swedish population  
677 yields support for greater CNV and MHC involvement in schizophrenia compared  
678 with bipolar disorder. *Mol Psychiatry* **17**:880–886. doi:10.1038/mp.2012.73

679 BridgeBuilder efficiently remaps BAM/SAM reads to a new reference by first building a  
680 "bridge" reference, first mapping to that bridge, and then remapping  
681 only a subset of reads to the fu.. 2015. . Wellcome Trust Sanger Institute -  
682 Human Genetics Informatics.

683 Browning SR, Browning BL. 2007. Rapid and Accurate Haplotype Phasing and Missing-  
684 Data Inference for Whole-Genome Association Studies By Use of Localized  
685 Haplotype Clustering. *Am J Hum Genet* **81**:1084–1097. doi:10.1086/521987

686 Chasman DI, Paré G, Mora S, Hopewell JC, Peloso G, Clarke R, Cupples LA, Hamsten  
687 A, Kathiresan S, Mälarstig A, Ordovas JM, Ripatti S, Parker AN, Miletich JP,  
688 Ridker PM. 2009. Forty-Three Loci Associated with Plasma Lipoprotein Size,  
689 Concentration, and Cholesterol Content in Genome-Wide Analysis. *PLOS Genet*  
690 **5**:e1000730. doi:10.1371/journal.pgen.1000730

691 Chen H. 2018. VennDiagram: Generate High-Resolution Venn and Euler Plots.

692 Chheda H, Palta P, Pirinen M, McCarthy S, Walter K, Koskinen S, Salomaa V, Daly M,  
693 Durbin R, Palotie A, Aittokallio T, Ripatti S. 2017. Whole-genome view of the  
694 consequences of a population bottleneck using 2926 genome sequences from  
695 Finland and United Kingdom. *Eur J Hum Genet*. doi:10.1038/ejhg.2016.205

696 Daly M, Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira M, Bender D, Maller J,  
697 Sklar P, Debakker P. 2007. PLINK: A Tool Set for Whole-Genome Association  
698 and Population-Based Linkage Analyses. *Am J Hum Genet* **81**:559–575.  
699 doi:10.1086/519795

700 Danjou F, Zoledziwska M, Sidore C, Steri M, Busonero F, Maschio A, Mulas A, Perseu  
701 L, Barella S, Porcu E, Pistis G, Pitzalis M, Pala M, Menzel S, Metrustry S,  
702 Spector TD, Leoni L, Angius A, Uda M, Moi P, Thein SL, Galanello R, Abecasis  
703 GR, Schlessinger D, Sanna S, Cucca F. 2015. Genome-wide association  
704 analyses based on whole-genome sequencing in Sardinia provide insights into  
705 regulation of hemoglobin levels. *Nat Genet* **advance online publication**.  
706 doi:10.1038/ng.3307

707 Davies G, Lam M, Harris SE, Trampush JW, Luciano M, Hill WD, Hagenaars SP, Ritchie  
708 SJ, Marioni RE, Fawns-Ritchie C, Liewald DCM, Okely JA, Ahola-Olli AV, Barnes  
709 CLK, Bertram L, Bis JC, Burdick KE, Christoforou A, DeRosse P, Djurovic S,  
710 Espeseth T, Giakoumaki S, Giddaluru S, Gustavson DE, Hayward C, Hofer E,  
711 Ikram MA, Karlsson R, Knowles E, Lahti J, Leber M, Li S, Mather KA, Melle I,  
712 Morris D, Oldmeadow C, Palviainen T, Payton A, Pazoki R, Petrovic K, Reynolds  
713 CA, Sargurupremraj M, Scholz M, Smith JA, Smith AV, Terzikhan N, Thalamuthu



714 A, Trompet S, Lee SJ van der, Ware EB, Windham BG, Wright MJ, Yang J, Yu J,  
715 Ames D, Amin N, Amouyel P, Andreassen OA, Armstrong NJ, Assareh AA, Attia  
716 JR, Attix D, Avramopoulos D, Bennett DA, Böhmer AC, Boyle PA, Brodaty H,  
717 Campbell H, Cannon TD, Cirulli ET, Congdon E, Conley ED, Corley J, Cox SR,  
718 Dale AM, Dehghan A, Dick D, Dickinson D, Eriksson JG, Evangelou E, Faul JD,  
719 Ford I, Freimer NA, Gao H, Giegling I, Gillespie NA, Gordon SD, Gottesman RF,  
720 Griswold ME, Gudnason V, Harris TB, Hartmann AM, Hatzimanolis A, Heiss G,  
721 Holliday EG, Joshi PK, Kähönen M, Kardia SLR, Karlsson I, Kleineidam L,  
722 Knopman DS, Kochan NA, Konte B, Kwok JB, Hellard SL, Lee T, Lehtimäki T, Li  
723 S-C, Liu T, Koini M, London E, Longstreth WT, Lopez OL, Loukola A, Luck T,  
724 Lundervold AJ, Lundquist A, Lyytikäinen L-P, Martin NG, Montgomery GW,  
725 Murray AD, Need AC, Noordam R, Nyberg L, Ollier W, Papenberg G, Pattie A,  
726 Polasek O, Poldrack RA, Psaty BM, Reppermund S, Riedel-Heller SG, Rose RJ,  
727 Rotter JI, Roussos P, Rovio SP, Saba Y, Sabb FW, Sachdev PS, Satizabal CL,  
728 Schmid M, Scott RJ, Scult MA, Simino J, Slagboom PE, Smyrnis N, Soumaré A,  
729 Stefanis NC, Stott DJ, Straub RE, Sundet K, Taylor AM, Taylor KD, Tzoulaki I,  
730 Tzourio C, Uitterlinden A, Vitart V, Voineskos AN, Kaprio J, Wagner M, Wagner  
731 H, Weinhold L, Wen KH, Widen E, Yang Q, Zhao W, Adams HHH, Arking DE,  
732 Bilder RM, Bitsios P, Boerwinkle E, Chiba-Falek O, Corvin A, Jager PLD, DeBette  
733 S, Donohoe G, Elliott P, Fitzpatrick AL, Gill M, Glahn DC, Hägg S, Hansell NK,  
734 Hariri AR, Ikram MK, Jukema JW, Vuoksimaa E, Keller MC, Kremen WS, Launer  
735 L, Lindenberger U, Palotie A, Pedersen NL, Pendleton N, Porteous DJ,  
736 Rääkkönen K, Raitakari OT, Ramirez A, Reinvang I, Rudan I, Rujescu D, Schmidt  
737 R, Schmidt H, Schofield PW, Schofield PR, Starr JM, Steen VM, Trollor JN,  
738 Turner ST, Duijn CMV, Villringer A, Weinberger DR, Weir DR, Wilson JF,  
739 Malhotra A, McIntosh AM, Gale CR, Seshadri S, Mosley TH, Bressler J, Lencz T,  
740 Deary IJ. 2018. Study of 300,486 individuals identifies 148 independent genetic  
741 loci influencing general cognitive function. *Nat Commun* **9**:2098.  
742 doi:10.1038/s41467-018-04362-x

743 Delaneau O, Howie B, Cox AJ, Zagury J-F, Marchini J. 2013. Haplotype Estimation  
744 Using Sequencing Reads. *Am J Hum Genet* **93**:687–696.

745 doi:10.1016/j.ajhg.2013.09.002

746 DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, Philippakis AA, del  
747 Angel G, Rivas MA, Hanna M, McKenna A, Fennell TJ, Kernytsky AM,  
748 Sivachenko AY, Cibulskis K, Gabriel SB, Altshuler D, Daly MJ. 2011. A  
749 framework for variation discovery and genotyping using next-generation DNA  
750 sequencing data. *Nat Genet* **43**:491–498. doi:10.1038/ng.806

751 Esko T, Mezzavilla M, Nelis M, Borel C, Debniak T, Jakkula E, Julia A, Karachanak S,  
752 Khrunin A, Kisfali P, Krulisova V, Aušrelė Kučinskienė Z, Rehnström K, Traglia M,  
753 Nikitina-Zake L, Zimprich F, Antonarakis SE, Estivill X, Glavač D, Gut I, Klovins J,  
754 Krawczak M, Kučinskas V, Lathrop M, Macek M, Marsal S, Meitinger T, Melegh  
755 B, Limborska S, Lubinski J, Paolotie A, Schreiber S, Toncheva D, Toniolo D,  
756 Wichmann H-E, Zimprich A, Metspalu M, Gasparini P, Metspalu A, D’Adamo P.  
757 2013. Genetic characterization of northeastern Italian population isolates in the  
758 context of broader European genetic diversity. *Eur J Hum Genet* **21**:659–665.  
759 doi:10.1038/ejhg.2012.229

760 Evangelou E, Warren HR, Mosen-Ansorena D, Mifsud B, Pazoki R, Gao H, Ntritsos G,  
761 Dimou N, Cabrera CP, Karaman I, Ng FL, Evangelou M, Witkowska K, Tzanis E,  
762 Hellwege JN, Giri A, Edwards DRV, Sun YV, Cho K, Gaziano JM, Wilson PWF,  
763 Tsao PS, Kovesdy CP, Esko T, Mägi R, Milani L, Almgren P, Boutin T, Debette S,  
764 Ding J, Giulianini F, Holliday EG, Jackson AU, Li-Gao R, Lin W-Y, Luan J,  
765 Mangino M, Oldmeadow C, Prins BP, Qian Y, Sargurupremraj M, Shah N,  
766 Surendran P, Thériault S, Verweij N, Willems SM, Zhao J-H, Amouyel P, Connell  
767 J, Mutsert R de, Doney ASF, Farrall M, Menni C, Morris AD, Noordam R, Paré G,  
768 Poulter NR, Shields DC, Stanton A, Thom S, Abecasis G, Amin N, Arking DE,  
769 Ayers KL, Barbieri CM, Batini C, Bis JC, Blake T, Bochud M, Boehnke M,  
770 Boerwinkle E, Boomsma DI, Bottinger EP, Braund PS, Brumat M, Campbell A,  
771 Campbell H, Chakravarti A, Chambers JC, Chauhan G, Ciullo M, Cocca M,  
772 Collins F, Cordell HJ, Davies G, Borst MH de, Geus EJ de, Deary IJ, Deelen J, M  
773 FDG, Demirkale CY, Dörr M, Ehret GB, Elosua R, Enroth S, Erzurumluoglu AM,  
774 Ferreira T, Frånberg M, Franco OH, Gandin I, Gasparini P, Giedraitis V, Gieger  
775 C, Girotto G, Goel A, Gow AJ, Gudnason V, Guo X, Gyllensten U, Hamsten A,

776 Harris TB, Harris SE, Hartman CA, Havulinna AS, Hicks AA, Hofer E, Hofman A,  
777 Hottenga J-J, Huffman JE, Hwang S-J, Ingelsson E, James A, Jansen R, Jarvelin  
778 M-R, Joehanes R, Johansson Å, Johnson AD, Joshi PK, Jousilahti P, Jukema  
779 JW, Jula A, Kähönen M, Kathiresan S, Keavney BD, Khaw K-T, Knekt P, Knight  
780 J, Kolcic I, Kooner JS, Koskinen S, Kristiansson K, Kutalik Z, Laan M, Larson M,  
781 Launer LJ, Lehne B, Lehtimäki T, Liewald DCM, Lin L, Lind L, Lindgren CM, Liu  
782 Y, Loos RJF, Lopez LM, Lu Y, Lyttikäinen L-P, Mahajan A, Mamasoula C,  
783 Marrugat J, Marten J, Milaneschi Y, Morgan A, Morris AP, Morrison AC, Munson  
784 PJ, Nalls MA, Nandakumar P, Nelson CP, Niiranen T, Nolte IM, Nutile T,  
785 Oldehinkel AJ, Oostra BA, O'Reilly PF, Org E, Padmanabhan S, Palmas W,  
786 Palotie A, Pattie A, Penninx BWJH, Perola M, Peters A, Polasek O, Pramstaller  
787 PP, Nguyen QT, Raitakari OT, Ren M, Rettig R, Rice K, Ridker PM, Ried JS,  
788 Riese H, Ripatti S, Robino A, Rose LM, Rotter JI, Rudan I, Ruggiero D, Saba Y,  
789 Sala CF, Salomaa V, Samani NJ, Sarin A-P, Schmidt R, Schmidt H, Shrine N,  
790 Siscovick D, Smith AV, Snieder H, Söber S, Sorice R, Starr JM, Stott DJ,  
791 Strachan DP, Strawbridge RJ, Sundström J, Swertz MA, Taylor KD, Teumer A,  
792 Tobin MD, Tomaszewski M, Toniolo D, Traglia M, Trompet S, Tuomilehto J,  
793 Tzourio C, Uitterlinden AG, Vaez A, Most PJ van der, Duijn CM van, Vergnaud A-  
794 C, Verwoert GC, Vitart V, Völker U, Vollenweider P, Vuckovic D, Watkins H, Wild  
795 SH, Willemsen G, Wilson JF, Wright AF, Yao J, Zemunik T, Zhang W, Attia JR,  
796 Butterworth AS, Chasman DI, Conen D, Cucca F, Danesh J, Hayward C, Howson  
797 JMM, Laakso M, Lakatta EG, Langenberg C, Melander O, Mook-Kanamori DO,  
798 Palmer CNA, Risch L, Scott RA, Scott RJ, Sever P, Spector TD, Harst P van der,  
799 Wareham NJ, Zeggini E, Levy D, Munroe PB, Newton-Cheh C, Brown MJ,  
800 Metspalu A, Hung AM, O'Donnell CJ, Edwards TL, Psaty BM, Tzoulaki I, Barnes  
801 MR, Wain LV, Elliott P, Caulfield MJ. 2018. Genetic analysis of over 1 million  
802 people identifies 535 new loci associated with blood pressure traits. *Nat Genet*  
803 **50**:1412. doi:10.1038/s41588-018-0205-x

804 Ferreira MAR, Hottenga J-J, Warrington NM, Medland SE, Willemsen G, Lawrence RW,  
805 Gordon S, de Geus EJC, Henders AK, Smit JH, Campbell MJ, Wallace L, Evans  
806 DM, Wright MJ, Nyholt DR, James AL, Beilby JP, Penninx BW, Palmer LJ, Frazer

807 IH, Montgomery GW, Martin NG, Boomsma DI. 2009. Sequence Variants in  
808 Three Loci Influence Monocyte Counts and Erythrocyte Volume. *Am J Hum*  
809 *Genet* **85**:745–749. doi:10.1016/j.ajhg.2009.10.005

810 Gudbjartsson DF, Helgason H, Gudjonsson SA, Zink F, Oddson A, Gylfason A,  
811 Besenbacher S, Magnusson G, Halldorsson BV, Hjartarson E, Sigurdsson GT,  
812 Stacey SN, Frigge ML, Holm H, Saemundsdottir J, Helgadóttir HT, Johannsdóttir  
813 H, Sigfusson G, Thorgeirsson G, Sverrisson JT, Gretarsdóttir S, Walters GB,  
814 Rafnar T, Thjodleifsson B, Bjornsson ES, Olafsson S, Thorarinsdóttir H,  
815 Steingrimsdóttir T, Gudmundsdóttir TS, Theodors A, Jonasson JG, Sigurdsson A,  
816 Bjornsdóttir G, Jonsson JJ, Thorarensen O, Ludvigsson P, Gudbjartsson H,  
817 Eyjolfsson GI, Sigurdardóttir O, Olafsson I, Arnar DO, Magnusson OT, Kong A,  
818 Masson G, Thorsteinsdóttir U, Helgason A, Sulem P, Stefansson K. 2015. Large-  
819 scale whole-genome sequencing of the Icelandic population. *Nat Genet* **47**:435–  
820 444. doi:10.1038/ng.3247

821 Hatzikotoulas K, Gilly A, Zeggini E. 2014. Using population isolates in genetic  
822 association studies. *Brief Funct Genomics* **13**:371–377. doi:10.1093/bfpg/elu022

823 Hoffmann TJ, Choquet H, Yin J, Banda Y, Kvale MN, Glymour M, Schaefer C, Risch N,  
824 Jorgenson E. 2018. A Large Multiethnic Genome-Wide Association Study of  
825 Adult Body Mass Index Identifies Novel Loci. *Genetics* **210**:499–515.  
826 doi:10.1534/genetics.118.301479

827 Howie B, Marchini J, Stephens M. 2011. Genotype Imputation with Thousands of  
828 Genomes. *G3 Genes Genomes Genet* **1**:457–470. doi:10.1534/g3.111.001198

829 Howie BN, Donnelly P, Marchini J. 2009. A Flexible and Accurate Genotype Imputation  
830 Method for the Next Generation of Genome-Wide Association Studies. *PLoS*  
831 *Genet* **5**:e1000529. doi:10.1371/journal.pgen.1000529

832 Karczewski KJ, Weisburd B, Thomas B, Solomonson M, Ruderfer DM, Kavanagh D,  
833 Hamamsy T, Lek M, Samocha KE, Cummings BB, Birnbaum D, Daly MJ,  
834 MacArthur DG. 2017. The ExAC browser: displaying reference data information  
835 from over 60 000 exomes. *Nucleic Acids Res* **45**:D840–D845.  
836 doi:10.1093/nar/gkw971

837 Kircher M, Witten DM, Jain P, O’Roak BJ, Cooper GM, Shendure J. 2014. A general

838 framework for estimating the relative pathogenicity of human genetic variants.  
839 *Nat Genet* **46**:310–315. doi:10.1038/ng.2892

840 Kullo IJ, Ding K, Jouni H, Smith CY, Chute CG. 2010. A Genome-Wide Association  
841 Study of Red Blood Cell Traits Using the Electronic Medical Record. *PLOS ONE*  
842 **5**:e13011. doi:10.1371/journal.pone.0013011

843 Lee JJ, Wedow R, Okbay A, Kong E, Maghazian O, Zacher M, Nguyen-Viet TA, Bowers  
844 P, Sidorenko J, Linnér RK, Fontana MA, Kundu T, Lee C, Li H, Li R, Royer R,  
845 Timshel PN, Walters RK, Willoughby EA, Yengo L, Alver M, Bao Y, Clark DW,  
846 Day FR, Furlotte NA, Joshi PK, Kemper KE, Kleinman A, Langenberg C, Mägi R,  
847 Trampush JW, Verma SS, Wu Y, Lam M, Zhao JH, Zheng Z, Boardman JD,  
848 Campbell H, Freese J, Harris KM, Hayward C, Herd P, Kumari M, Lencz T, Luan  
849 J, Malhotra AK, Metspalu A, Milani L, Ong KK, Perry JRB, Porteous DJ, Ritchie  
850 MD, Smart MC, Smith BH, Tung JY, Wareham NJ, Wilson JF, Beauchamp JP,  
851 Conley DC, Esko T, Lehrer SF, Magnusson PKE, Oskarsson S, Pers TH,  
852 Robinson MR, Thom K, Watson C, Chabris CF, Meyer MN, Laibson DI, Yang J,  
853 Johannesson M, Koellinger PD, Turley P, Visscher PM, Benjamin DJ, Cesarini D.  
854 2018. Gene discovery and polygenic prediction from a genome-wide association  
855 study of educational attainment in 1.1 million individuals. *Nat Genet* **50**:1112.  
856 doi:10.1038/s41588-018-0147-3

857 Li H. 2011. A statistical framework for SNP calling, mutation discovery, association  
858 mapping and population genetical parameter estimation from sequencing data.  
859 *Bioinforma Oxf Engl* **27**:2987–2993. doi:10.1093/bioinformatics/btr509

860 Li H, Durbin R. 2010. Fast and accurate long-read alignment with Burrows–Wheeler  
861 transform. *Bioinformatics* **26**:589–595. doi:10.1093/bioinformatics/btp698

862 Li YR, Li J, Zhao SD, Bradfield JP, Mentch FD, Maggadottir SM, Hou C, Abrams DJ,  
863 Chang D, Gao F, Guo Y, Wei Z, Connolly JJ, Cardinale CJ, Bakay M, Glessner  
864 JT, Li D, Kao C, Thomas KA, Qiu H, Chiavacci RM, Kim CE, Wang F, Snyder J,  
865 Richie MD, Flatø B, Førre Ø, Denson LA, Thompson SD, Becker ML, Guthery  
866 SL, Latiano A, Perez E, Resnick E, Russell RK, Wilson DC, Silverberg MS,  
867 Annese V, Lie BA, Punaro M, Dubinsky MC, Monos DS, Strisciuglio C, Staiano A,  
868 Miele E, Kugathasan S, Ellis JA, Munro JE, Sullivan KE, Wise CA, Chapel H,

869 Cunningham-Rundles C, Grant SFA, Orange JS, Sleiman PMA, Behrens EM,  
870 Griffiths AM, Satsangi J, Finkel TH, Keinan A, Prak ETL, Polychronakos C,  
871 Baldassano RN, Li H, Keating BJ, Hakonarson H. 2015. Meta-analysis of shared  
872 genetic architecture across ten pediatric autoimmune diseases. *Nat Med*  
873 **21**:1018–1027. doi:10.1038/nm.3933

874 Liu C, Kraja AT, Smith JA, Brody JA, Franceschini N, Bis JC, Rice K, Morrison AC, Lu  
875 Y, Weiss S, Guo X, Palmas W, Martin LW, Chen Y-DI, Surendran P, Drenos F,  
876 Cook JP, Auer PL, Chu AY, Giri A, Zhao W, Jakobsdottir J, Lin L-A, Stafford JM,  
877 Amin N, Mei H, Yao J, Voorman A, CHD Exome+ Consortium, ExomeBP  
878 Consortium, GoT2DGenes Consortium, T2d-Genes Consortium, Larson MG,  
879 Grove ML, Smith AV, Hwang S-J, Chen H, Huan T, Kosova G, Stitzel NO,  
880 Kathiresan S, Samani N, Schunkert H, Deloukas P, Myocardial Infarction  
881 Genetics and CARDIoGRAM Exome Consortia, Li M, Fuchsberger C, Pattaro C,  
882 Gorski M, CKDGen Consortium, Kooperberg C, Papanicolaou GJ, Rossouw JE,  
883 Faul JD, Kardina SLR, Bouchard C, Raffel LJ, Uitterlinden AG, Franco OH, Vasan  
884 RS, O'Donnell CJ, Taylor KD, Liu K, Bottinger EP, Gottesman O, Daw EW,  
885 Giulianini F, Ganesh S, Salfati E, Harris TB, Launer LJ, Dörr M, Felix SB, Rettig  
886 R, Völzke H, Kim E, Lee W-J, Lee I-T, Sheu WH-H, Tsosie KS, Edwards DRV,  
887 Liu Y, Correa A, Weir DR, Völker U, Ridker PM, Boerwinkle E, Gudnason V,  
888 Reiner AP, van Duijn CM, Borecki IB, Edwards TL, Chakravarti A, Rotter JI,  
889 Psaty BM, Loos RJF, Fornage M, Ehret GB, Newton-Cheh C, Levy D, Chasman  
890 DI. 2016. Meta-analysis identifies common and rare variants influencing blood  
891 pressure and overlapping with metabolic trait loci. *Nat Genet* **48**:1162–1170.  
892 doi:10.1038/ng.3660

893 Lucas A. 2018. An R package for creating mirrored Manhattan plots: anastasia-  
894 lucas/hudson.

895 MacArthur DG, Balasubramanian S, Frankish A, Huang N, Morris J, Walter K, Jostins L,  
896 Habegger L, Pickrell JK, Montgomery SB, Albers CA, Zhang ZD, Conrad DF,  
897 Lunter G, Zheng H, Ayub Q, DePristo MA, Banks E, Hu M, Handsaker RE,  
898 Rosenfeld JA, Fromer M, Jin M, Mu XJ, Khurana E, Ye K, Kay M, Saunders GI,  
899 Suner M-M, Hunt T, Barnes IHA, Amid C, Carvalho-Silva DR, Bignell AH, Snow

900 C, Yngvadottir B, Bumpstead S, Cooper DN, Xue Y, Romero IG, 1000 Genomes  
901 Project Consortium, Wang J, Li Y, Gibbs RA, McCarroll SA, Dermitzakis ET,  
902 Pritchard JK, Barrett JC, Harrow J, Hurler ME, Gerstein MB, Tyler-Smith C.  
903 2012. A systematic survey of loss-of-function variants in human protein-coding  
904 genes. *Science* **335**:823–828. doi:10.1126/science.1215040

905 Marchini J, Howie B. 2010. Genotype imputation for genome-wide association studies.  
906 *Nat Rev Genet* **11**:499–511. doi:10.1038/nrg2796

907 Martinelli-Boneschi F, Colombi M, Castori M, Devigili G, Eleopra R, Malik RA, Ritelli M,  
908 Zoppi N, Dordoni C, Sorosina M, Grammatico P, Fadavi H, Gerrits MM,  
909 Almomani R, Faber CG, Merkies ISJ, Toniolo D, Cocca M, Doglioni C, Waxman  
910 SG, Dib-Hajj SD, Taiana MM, Sassone J, Lombardi R, Cazzato D, Zauli A,  
911 Santoro S, Marchi M, Lauria G. n.d. COL6A5 variants in familial neuropathic  
912 chronic itch. *Brain*. doi:10.1093/brain/aww343

913 McCarthy S, Das S, Kretzschmar W, Delaneau O, Wood AR, Teumer A, Kang HM,  
914 Fuchsberger C, Danecek P, Sharp K, Luo Y, Sidore C, Kwong A, Timpson N,  
915 Koskinen S, Vrieze S, Scott LJ, Zhang H, Mahajan A, Veldink J, Peters U, Pato  
916 C, van Duijn CM, Gillies CE, Gandin I, Mezzavilla M, Gilly A, Cocca M, Traglia M,  
917 Angius A, Barrett JC, Boomsma D, Branham K, Breen G, Brummett CM,  
918 Busonero F, Campbell H, Chan A, Chen S, Chew E, Collins FS, Corbin LJ, Smith  
919 GD, Dedoussis G, Dorr M, Farmaki A-E, Ferrucci L, Forer L, Fraser RM, Gabriel  
920 S, Levy S, Groop L, Harrison T, Hattersley A, Holmen OL, Hveem K, Kretzler M,  
921 Lee JC, McGue M, Meitinger T, Melzer D, Min JL, Mohlke KL, Vincent JB, Nauck  
922 M, Nickerson D, Palotie A, Pato M, Pirastu N, McInnis M, Richards JB, Sala C,  
923 Salomaa V, Schlessinger D, Schoenherr S, Slagboom PE, Small K, Spector T,  
924 Stambolian D, Tuke M, Tuomilehto J, Van den Berg LH, Van Rheenen W, Volker  
925 U, Wijmenga C, Toniolo D, Zeggini E, Gasparini P, Sampson MG, Wilson JF,  
926 Frayling T, de Bakker PIW, Swertz MA, McCarroll S, Kooperberg C, Dekker A,  
927 Altshuler D, Willer C, Iacono W, Ripatti S, Soranzo N, Walter K, Swaroop A,  
928 Cucca F, Anderson CA, Myers RM, Boehnke M, McCarthy MI, Durbin R, the  
929 Haplotype Reference Consortium. 2016. A reference panel of 64,976 haplotypes  
930 for genotype imputation. *Nat Genet* **48**:1279–1283. doi:10.1038/ng.3643

931 McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K,  
932 Altshuler D, Gabriel S, Daly M, DePristo MA. 2010. The Genome Analysis  
933 Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing  
934 data. *Genome Res* **20**:1297–1303. doi:10.1101/gr.107524.110

935 McLaren W, Pritchard B, Rios D, Chen Y, Flicek P, Cunningham F. 2010. Deriving the  
936 consequences of genomic variants with the Ensembl API and SNP Effect  
937 Predictor. *Bioinforma Oxf Engl* **26**:2069–2070. doi:10.1093/bioinformatics/btq330

938 Narasimhan VM, Hunt KA, Mason D, Baker CL, Karczewski KJ, Barnes MR, Barnett AH,  
939 Bates C, Bellary S, Bockett NA, Giorda K, Griffiths CJ, Hemingway H, Jia Z, Kelly  
940 MA, Khawaja HA, Lek M, McCarthy S, McEachan R, O'Donnell-Luria A, Paigen  
941 K, Parisinos CA, Sheridan E, Southgate L, Tee L, Thomas M, Xue Y, Schnall-  
942 Levin M, Petkov PM, Tyler-Smith C, Maher ER, Trembath RC, MacArthur DG,  
943 Wright J, Durbin R, Heel DA van. 2016. Health and population effects of rare  
944 gene knockouts in adult humans with related parents. *Science* **352**:474–477.  
945 doi:10.1126/science.aac8624

946 Nielsen JB, Thorolfsdottir RB, Fritsche LG, Zhou W, Skov MW, Graham SE, Herron TJ,  
947 McCarthy S, Schmidt EM, Sveinbjornsson G, Surakka I, Mathis MR, Yamazaki M,  
948 Crawford RD, Gabrielsen ME, Skogholt AH, Holmen OL, Lin M, Wolford BN, Dey  
949 R, Dalen H, Sulem P, Chung JH, Backman JD, Arnar DO, Thorsteinsdottir U,  
950 Baras A, O'Dushlaine C, Holst AG, Wen X, Hornsby W, Dewey FE, Boehnke M,  
951 Kheterpal S, Mukherjee B, Lee S, Kang HM, Holm H, Kitzman J, Shavit JA, Jalife  
952 J, Brummett CM, Teslovich TM, Carey DJ, Gudbjartsson DF, Stefansson K,  
953 Abecasis GR, Hveem K, Willer CJ. 2018. Biobank-driven genomic discovery  
954 yields new insight into atrial fibrillation biology. *Nat Genet* **50**:1234.  
955 doi:10.1038/s41588-018-0171-3

956 Perry JRB, Day F, Elks CE, Sulem P, Thompson DJ, Ferreira T, He C, Chasman DI,  
957 Esko T, Thorleifsson G, Albrecht E, Ang WQ, Corre T, Cousminer DL, Feenstra  
958 B, Franceschini N, Ganna A, Johnson AD, Kjellqvist S, Lunetta KL, McMahon G,  
959 Nolte IM, Paternoster L, Porcu E, Smith AV, Stolk L, Teumer A, Tšernikova N,  
960 Tikkanen E, Ulivi S, Wagner EK, Amin N, Bierut LJ, Byrne EM, Hottenga J-J,  
961 Koller DL, Mangino M, Pers TH, Yerges-Armstrong LM, Hua Zhao J, Andrusis IL,



962 Anton-Culver H, Atsma F, Bandinelli S, Beckmann MW, Benitez J, Blomqvist C,  
963 Bojesen SE, Bolla MK, Bonanni B, Brauch H, Brenner H, Buring JE, Chang-  
964 Claude J, Chanock S, Chen J, Chenevix-Trench G, Collée JM, Couch FJ, Couper  
965 D, Coviello AD, Cox A, Czene K, D'adamo AP, Davey Smith G, De Vivo I,  
966 Demerath EW, Dennis J, Devilee P, Dieffenbach AK, Dunning AM, Eiriksdottir G,  
967 Eriksson JG, Fasching PA, Ferrucci L, Flesch-Janys D, Flyger H, Foroud T,  
968 Franke L, Garcia ME, García-Closas M, Geller F, de Geus EEJ, Giles GG,  
969 Gudbjartsson DF, Gudnason V, Guénel P, Guo S, Hall P, Hamann U, Haring R,  
970 Hartman CA, Heath AC, Hofman A, Hooning MJ, Hopper JL, Hu FB, Hunter DJ,  
971 Karasik D, Kiel DP, Knight JA, Kosma V-M, Kutalik Z, Lai S, Lambrechts D,  
972 Lindblom A, Mägi R, Magnusson PK, Mannermaa A, Martin NG, Masson G,  
973 McArdle PF, McArdle WL, Melbye M, Michailidou K, Mihailov E, Milani L, Milne  
974 RL, Nevanlinna H, Neven P, Nohr EA, Oldehinkel AJ, Oostra BA, Palotie A,  
975 Peacock M, Pedersen NL, Peterlongo P, Peto J, Pharoah PDP, Postma DS,  
976 Pouta A, Pylkäs K, Radice P, Ring S, Rivadeneira F, Robino A, Rose LM,  
977 Rudolph A, Salomaa V, Sanna S, Schlessinger D, Schmidt MK, Southey MC,  
978 Sovio U, Stampfer MJ, Stöckl D, Storniolo AM, Timpson NJ, Tyrer J, Visser JA,  
979 Vollenweider P, Völzke H, Waeber G, Waldenberger M, Wallaschofski H, Wang  
980 Q, Willemsen G, Winqvist R, Wolffenbuttel BHR, Wright MJ, Australian Ovarian  
981 Cancer Study, The GENICA Network, kConFab, The LifeLines Cohort Study, The  
982 InterAct Consortium, Early Growth Genetics (EGG) Consortium, Boomsma DI,  
983 Econs MJ, Khaw K-T, Loos RJF, McCarthy MI, Montgomery GW, Rice JP,  
984 Streeten EA, Thorsteinsdottir U, van Duijn CM, Alizadeh BZ, Bergmann S,  
985 Boerwinkle E, Boyd HA, Crisponi L, Gasparini P, Gieger C, Harris TB, Ingelsson  
986 E, Järvelin M-R, Kraft P, Lawlor D, Metspalu A, Pennell CE, Ridker PM, Snieder  
987 H, Sørensen TIA, Spector TD, Strachan DP, Uitterlinden AG, Wareham NJ,  
988 Widen E, Zygumt M, Murray A, Easton DF, Stefansson K, Murabito JM, Ong  
989 KK. 2014. Parent-of-origin-specific allelic associations among 106 genomic loci  
990 for age at menarche. *Nature* **514**:92–97. doi:10.1038/nature13545  
991 Pickrell JK, Berisa T, Liu JZ, Séguérel L, Tung JY, Hinds DA. 2016. Detection and  
992 interpretation of shared genetic influences on 42 human traits. *Nat Genet*

993           **48**:709–717. doi:10.1038/ng.3570

994 Pickrell JK, Pritchard JK. 2012. Inference of Population Splits and Mixtures from  
995           Genome-Wide Allele Frequency Data. *PLoS Genet* **8**:e1002967.  
996           doi:10.1371/journal.pgen.1002967

997 Pistis G, Porcu E, Vrieze SI, Sidore C, Steri M, Danjou F, Busonero F, Mulas A,  
998           Zoledziewska M, Maschio A, Brennan C, Lai S, Miller MB, Marcelli M, Urru MF,  
999           Pitzalis M, Lyons RH, Kang HM, Jones CM, Angius A, Iacono WG, Schlessinger  
1000           D, McGue M, Cucca F, Abecasis GR, Sanna S. 2015. Rare variant genotype  
1001           imputation with thousands of study-specific whole-genome sequences:  
1002           implications for cost-effective study designs. *Eur J Hum Genet* **23**:975–983.  
1003           doi:10.1038/ejhg.2014.216

1004 Sazzini M, Gnechi Ruscone GA, Giuliani C, Sarno S, Quagliariello A, De Fanti S,  
1005           Boattini A, Gentilini D, Fiorito G, Catanoso M, Boiardi L, Croci S, Macchioni P,  
1006           Mantovani V, Di Blasio AM, Matullo G, Salvarani C, Franceschi C, Pettener D,  
1007           Garagnani P, Luiselli D. 2016. Complex interplay between neutral and adaptive  
1008           evolution shaped differential genomic background and disease susceptibility  
1009           along the Italian peninsula. *Sci Rep* **6**:32513. doi:10.1038/srep32513

1010 Schormair B, Zhao C, Bell S, Tilch E, Salminen AV, Pütz B, Dauvilliers Y, Stefani A,  
1011           Högl B, Poewe W, Kemlink D, Sonka K, Bachmann CG, Paulus W, Trenkwalder  
1012           C, Oertel WH, Hornyak M, Teder-Laving M, Metspalu A, Hadjigeorgiou GM, Polo  
1013           O, Fietze I, Ross OA, Wszolek Z, Butterworth AS, Soranzo N, Ouwehand WH,  
1014           Roberts DJ, Danesh J, Allen RP, Earley CJ, Ondo WG, Xiong L, Montplaisir J,  
1015           Gan-Or Z, Perola M, Vodicka P, Dina C, Franke A, Tittmann L, Stewart AFR,  
1016           Shah SH, Gieger C, Peters A, Rouleau GA, Berger K, Oexle K, Di Angelantonio  
1017           E, Hinds DA, Müller-Myhsok B, Winkelmann J, Balkau B, Ducimetière P,  
1018           Eschwège E, Rancièrè F, Alhenc-Gelas F, Gallois Y, Girault A, Fumeron F, Marre  
1019           M, Roussel R, Bonnet F, Bonnefond A, Cauchi S, Froguel P, Cogneau J, Born C,  
1020           Caces E, Cailleau M, Lantieri O, Moreau J, Rakotozafy F, Tichet J, Vol S, Agee  
1021           M, Alipanahi B, Auton A, Bell RK, Bryc K, Elson SL, Fontanillas P, Furlotte NA,  
1022           Hinds DA, Hromatka BS, Huber KE, Kleinman A, Litterman NK, McIntyre MH,  
1023           Mountain JL, Northover CA, Pitts SJ, Sathirapongsasuti JF, Sazonova OV,

1024 Shelton JF, Shringarpure S, Tian C, Tung JY, Vacic V, Wilson CH. 2017.  
1025 Identification of novel risk loci for restless legs syndrome in genome-wide  
1026 association studies in individuals of European ancestry: a meta-analysis. *Lancet*  
1027 *Neurol* **16**:898–907. doi:10.1016/S1474-4422(17)30327-7

1028 Sherry ST, Ward MH, Kholodov M, Baker J, Phan L, Smigielski EM, Sirotkin K. 2001.  
1029 dbSNP: the NCBI database of genetic variation. *Nucleic Acids Res* **29**:308–311.

1030 Steensel MAM van, Steijlen PM, Bladergroen RS, Vermeer M, Geel M van. 2005. A  
1031 missense mutation in the type II hair keratin hHb3 is associated with monilethrix.  
1032 *J Med Genet* **42**:e19–e19. doi:10.1136/jmg.2004.021030

1033 Sudmant PH, Rausch T, Gardner EJ, Handsaker RE, Abyzov A, Huddleston J, Zhang Y,  
1034 Ye K, Jun G, Hsi-Yang Fritz M, Konkol MK, Malhotra A, Stütz AM, Shi X, Paolo  
1035 Casale F, Chen J, Hormozdiari F, Dayama G, Chen K, Malig M, Chaisson MJP,  
1036 Walter K, Meiers S, Kashin S, Garrison E, Auton A, Lam HYK, Jasmine Mu X,  
1037 Alkan C, Antaki D, Bae T, Cerveira E, Chines P, Chong Z, Clarke L, Dal E, Ding  
1038 L, Emery S, Fan X, Gujral M, Kahveci F, Kidd JM, Kong Y, Lammeijer E-W,  
1039 McCarthy S, Flicek P, Gibbs RA, Marth G, Mason CE, Menelaou A, Muzny DM,  
1040 Nelson BJ, Noor A, Parrish NF, Pendleton M, Quitadamo A, Raeder B, Schadt  
1041 EE, Romanovitch M, Schlattl A, Sebra R, Shabalina AA, Untergasser A, Walker  
1042 JA, Wang M, Yu F, Zhang C, Zhang J, Zheng-Bradley X, Zhou W, Zichner T,  
1043 Sebat J, Batzer MA, McCarroll SA, The 1000 Genomes Project Consortium, Mills  
1044 RE, Gerstein MB, Bashir A, Stegle O, Devine SE, Lee C, Eichler EE, Korbel JO.  
1045 2015. An integrated map of structural variation in 2,504 human genomes. *Nature*  
1046 **526**:75–81. doi:10.1038/nature15394

1047 Sun BB, Maranville JC, Peters JE, Stacey D, Staley JR, Blackshaw J, Burgess S, Jiang  
1048 T, Paige E, Surendran P, Oliver-Williams C, Kamat MA, Prins BP, Wilcox SK,  
1049 Zimmerman ES, Chi A, Bansal N, Spain SL, Wood AM, Morrell NW, Bradley JR,  
1050 Janjic N, Roberts DJ, Ouwehand WH, Todd JA, Soranzo N, Suhre K, Paul DS,  
1051 Fox CS, Plenge RM, Danesh J, Runz H, Butterworth AS. 2018. Genomic atlas of  
1052 the human plasma proteome. *Nature* **558**:73. doi:10.1038/s41586-018-0175-2

1053 Szpiech ZA, Hernandez RD. 2014. selscan: An Efficient Multithreaded Program to  
1054 Perform EHH-Based Scans for Positive Selection. *Mol Biol Evol* **31**:2824–2827.

1055 doi:10.1093/molbev/msu211

1056 Teslovich TM, Musunuru K, Smith AV, Edmondson AC, Stylianou IM, Koseki M,  
1057 Pirruccello JP, Ripatti S, Chasman DI, Willer CJ, Johansen CT, Fouchier SW,  
1058 Isaacs A, Peloso GM, Barbalic M, Ricketts SL, Bis JC, Aulchenko YS,  
1059 Thorleifsson G, Feitosa MF, Chambers J, Orho-Melander M, Melander O,  
1060 Johnson T, Li X, Guo X, Li M, Shin Cho Y, Jin Go M, Jin Kim Y, Lee J-Y, Park T,  
1061 Kim K, Sim X, Twee-Hee Ong R, Croteau-Chonka DC, Lange LA, Smith JD,  
1062 Song K, Hua Zhao J, Yuan X, Luan J, Lamina C, Ziegler A, Zhang W, Zee RYL,  
1063 Wright AF, Witteman JCM, Wilson JF, Willemsen G, Wichmann H-E, Whitfield  
1064 JB, Waterworth DM, Wareham NJ, Waeber G, Vollenweider P, Voight BF, Vitart  
1065 V, Uitterlinden AG, Uda M, Tuomilehto J, Thompson JR, Tanaka T, Surakka I,  
1066 Stringham HM, Spector TD, Soranzo N, Smit JH, Sinisalo J, Silander K, Sijbrands  
1067 EJG, Scuteri A, Scott J, Schlessinger D, Sanna S, Salomaa V, Saharinen J,  
1068 Sabatti C, Ruukonen A, Rudan I, Rose LM, Roberts R, Rieder M, Psaty BM,  
1069 Pramstaller PP, Pichler I, Perola M, Penninx BWJH, Pedersen NL, Pattaro C,  
1070 Parker AN, Pare G, Oostra BA, O'Donnell CJ, Nieminen MS, Nickerson DA,  
1071 Montgomery GW, Meitinger T, McPherson R, McCarthy MI, McArdle W, Masson  
1072 D, Martin NG, Marroni F, Mangino M, Magnusson PKE, Lucas G, Luben R, Loos  
1073 RJF, Lokki M-L, Lettre G, Langenberg C, Launer LJ, Lakatta EG, Laaksonen R,  
1074 Kyvik KO, Kronenberg F, König IR, Khaw K-T, Kaprio J, Kaplan LM, Johansson  
1075 Å, Jarvelin M-R, Cecile J. W. Janssens A, Ingelsson E, Igl W, Kees Hovingh G,  
1076 Hottenga J-J, Hofman A, Hicks AA, Hengstenberg C, Heid IM, Hayward C,  
1077 Havulinna AS, Hastie ND, Harris TB, Haritunians T, Hall AS, Gyllensten U,  
1078 Guiducci C, Groop LC, Gonzalez E, Gieger C, Freimer NB, Ferrucci L, Erdmann  
1079 J, Elliott P, Ejebe KG, Döring A, Dominiczak AF, Demissie S, Deloukas P, de  
1080 Geus EJC, de Faire U, Crawford G, Collins FS, Chen YI, Caulfield MJ, Campbell  
1081 H, Burtt NP, Bonnycastle LL, Boomsma DI, Boekholdt SM, Bergman RN, Barroso  
1082 I, Bandinelli S, Ballantyne CM, Assimes TL, Quertermous T, Altshuler D,  
1083 Seielstad M, Wong TY, Tai E-S, Feranil AB, Kuzawa CW, Adair LS, Taylor Jr HA,  
1084 Borecki IB, Gabriel SB, Wilson JG, Holm H, Thorsteinsdottir U, Gudnason V,  
1085 Krauss RM, Mohlke KL, Ordovas JM, Munroe PB, Kooner JS, Tall AR, Hegele

1086 RA, Kastelein JJP, Schadt EE, Rotter JI, Boerwinkle E, Strachan DP, Mooser V,  
1087 Stefansson K, Reilly MP, Samani NJ, Schunkert H, Cupples LA, Sandhu MS,  
1088 Ridker PM, Rader DJ, van Duijn CM, Peltonen L, Abecasis GR, Boehnke M,  
1089 Kathiresan S. 2010. Biological, clinical and population relevance of 95 loci for  
1090 blood lipids. *Nature* **466**:707–713. doi:10.1038/nature09270

1091 The 1000 Genomes Project Consortium. 2015. A global reference for human genetic  
1092 variation. *Nature* **526**:68–74. doi:10.1038/nature15393

1093 The ENCODE Project Consortium. 2007. Identification and analysis of functional  
1094 elements in 1% of the human genome by the ENCODE pilot project. *Nature*  
1095 **447**:799–816. doi:10.1038/nature05874

1096 The UK10K Consortium. 2015. The UK10K project identifies rare variants in health and  
1097 disease. *Nature* **526**:82–90. doi:10.1038/nature14962

1098 Turner S. 2017. qqman: Q-Q and Manhattan Plots for GWAS Data.

1099 Voight BF, Kudaravalli S, Wen X, Pritchard JK. 2006. A Map of Recent Positive  
1100 Selection in the Human Genome. *PLOS Biol* **4**:e72.  
1101 doi:10.1371/journal.pbio.0040072

1102 Willer CJ, Li Y, Abecasis GR. 2010. METAL: fast and efficient meta-analysis of  
1103 genomewide association scans. *Bioinforma Oxf Engl* **26**:2190–2191.  
1104 doi:10.1093/bioinformatics/btq340

1105 Wood AR, Esko T, Yang J, Vedantam S, Pers TH, Gustafsson S, Chu AY, Estrada K,  
1106 Luan J, Kutalik Z, Amin N, Buchkovich ML, Croteau-Chonka DC, Day FR, Duan  
1107 Y, Fall T, Fehrmann R, Ferreira T, Jackson AU, Karjalainen J, Lo KS, Locke AE,  
1108 Mägi R, Mihailov E, Porcu E, Randall JC, Scherag A, Vinkhuyzen AAE, Westra  
1109 H-J, Winkler TW, Workalemahu T, Zhao JH, Absher D, Albrecht E, Anderson D,  
1110 Baron J, Beekman M, Demirkan A, Ehret GB, Feenstra B, Feitosa MF, Fischer K,  
1111 Fraser RM, Goel A, Gong J, Justice AE, Kanoni S, Kleber ME, Kristiansson K,  
1112 Lim U, Lotay V, Lui JC, Mangino M, Leach IM, Medina-Gomez C, Nalls MA,  
1113 Nyholt DR, Palmer CD, Pasko D, Pechlivanis S, Prokopenko I, Ried JS, Ripke S,  
1114 Shungin D, Stancáková A, Strawbridge RJ, Sung YJ, Tanaka T, Teumer A,  
1115 Trompet S, van der Laan SW, van Setten J, Van Vliet-Ostaptchouk JV, Wang Z,  
1116 Yengo L, Zhang W, Afzal U, Ärnlöv J, Arscott GM, Bandinelli S, Barrett A, Bellis

1117 C, Bennett AJ, Berne C, Blüher M, Bolton JL, Böttcher Y, Boyd HA, Bruinenberg  
1118 M, Buckley BM, Buyske S, Caspersen IH, Chines PS, Clarke R, Claudi-Boehm S,  
1119 Cooper M, Daw EW, De Jong PA, Deelen J, Delgado G, Denny JC, Dhonukshe-  
1120 Rutten R, Dimitriou M, Doney ASF, Dörr M, Eklund N, Eury E, Folkersen L,  
1121 Garcia ME, Geller F, Giedraitis V, Go AS, Grallert H, Grammer TB, Gräßler J,  
1122 Grönberg H, de Groot LCPGM, Groves CJ, Haessler J, Hall P, Haller T, Hallmans  
1123 G, Hannemann A, Hartman CA, Hassinen M, Hayward C, Heard-Costa NL,  
1124 Helmer Q, Hemani G, Henders AK, Hillege HL, Hlatky MA, Hoffmann W,  
1125 Hoffmann P, Holmen O, Houwing-Duistermaat JJ, Illig T, Isaacs A, James AL,  
1126 Jeff J, Johansen B, Johansson Å, Jolley J, Juliusdottir T, Junttila J, Kho AN,  
1127 Kinnunen L, Klopp N, Kocher T, Kratzer W, Lichtner P, Lind L, Lindström J,  
1128 Lobbens S, Lorentzon M, Lu Y, Lyssenko V, Magnusson PKE, Mahajan A,  
1129 Maillard M, McArdle WL, McKenzie CA, McLachlan S, McLaren PJ, Menni C,  
1130 Merger S, Milani L, Moayyeri A, Monda KL, Morken MA, Müller G, Müller-  
1131 Nurasyid M, Musk AW, Narisu N, Nauck M, Nolte IM, Nöthen MM, Oozageer L,  
1132 Pilz S, Rayner NW, Renstrom F, Robertson NR, Rose LM, Roussel R, Sanna S,  
1133 Scharnagl H, Scholtens S, Schumacher FR, Schunkert H, Scott RA, Sehmi J,  
1134 Seufferlein T, Shi J, Silventoinen K, Smit JH, Smith AV, Smolonska J, Stanton  
1135 AV, Stirrups K, Stott DJ, Stringham HM, Sundström J, Swertz MA, Syvänen A-C,  
1136 Tayo BO, Thorleifsson G, Tyrer JP, van Dijk S, van Schoor NM, van der Velde N,  
1137 van Heemst D, van Oort FVA, Vermeulen SH, Verweij N, Vonk JM, Waite LL,  
1138 Waldenberger M, Wennauer R, Wilkens LR, Willenborg C, Wilsgaard T,  
1139 Wojczynski MK, Wong A, Wright AF, Zhang Q, Arveiler D, Bakker SJL, Beilby J,  
1140 Bergman RN, Bergmann S, Biffar R, Blangero J, Boomsma DI, Bornstein SR,  
1141 Bovet P, Brambilla P, Brown MJ, Campbell H, Caulfield MJ, Chakravarti A,  
1142 Collins R, Collins FS, Crawford DC, Cupples LA, Danesh J, de Faire U, den  
1143 Ruijter HM, Erbel R, Erdmann J, Eriksson JG, Farrall M, Ferrannini E, Ferrières J,  
1144 Ford I, Forouhi NG, Forrester T, Gansevoort RT, Gejman PV, Gieger C, Golay A,  
1145 Gottesman O, Gudnason V, Gyllensten U, Haas DW, Hall AS, Harris TB,  
1146 Hattersley AT, Heath AC, Hengstenberg C, Hicks AA, Hindorff LA, Hingorani AD,  
1147 Hofman A, Hovingh GK, Humphries SE, Hunt SC, Hyponen E, Jacobs KB,

1148 Jarvelin M-R, Jousilahti P, Jula AM, Kaprio J, Kastelein JJP, Kayser M, Kee F,  
1149 Keinanen-Kiukaanniemi SM, Kiemeny LA, Kooner JS, Kooperberg C, Koskinen  
1150 S, Kovacs P, Kraja AT, Kumari M, Kuusisto J, Lakka TA, Langenberg C, Le  
1151 Marchand L, Lehtimäki T, Lupoli S, Madden PAF, Männistö S, Manunta P,  
1152 Marette A, Matise TC, McKnight B, Meitinger T, Moll FL, Montgomery GW, Morris  
1153 AD, Morris AP, Murray JC, Nelis M, Ohlsson C, Oldehinkel AJ, Ong KK,  
1154 Ouwehand WH, Pasterkamp G, Peters A, Pramstaller PP, Price JF, Qi L,  
1155 Raitakari OT, Rankinen T, Rao DC, Rice TK, Ritchie M, Rudan I, Salomaa V,  
1156 Samani NJ, Saramies J, Sarzynski MA, Schwarz PEH, Sebert S, Sever P,  
1157 Shuldiner AR, Sinisalo J, Steinthorsdottir V, Stolk RP, Tardif J-C, Tönjes A,  
1158 Tremblay A, Tremoli E, Virtamo J, Vohl M-C, The Electronic Medical Records and  
1159 Genomics (eMERGE) Consortium, The MIGen Consortium, The PAGE  
1160 Consortium, The LifeLines Cohort Study, Amouyel P, Asselbergs FW, Assimes  
1161 TL, Bochud M, Boehm BO, Boerwinkle E, Bottinger EP, Bouchard C, Cauchi S,  
1162 Chambers JC, Chanock SJ, Cooper RS, de Bakker PIW, Dedoussis G, Ferrucci  
1163 L, Franks PW, Froguel P, Groop LC, Haiman CA, Hamsten A, Hayes MG, Hui J,  
1164 Hunter DJ, Hveem K, Jukema JW, Kaplan RC, Kivimaki M, Kuh D, Laakso M, Liu  
1165 Y, Martin NG, März W, Melbye M, Moebus S, Munroe PB, Njølstad I, Oostra BA,  
1166 Palmer CNA, Pedersen NL, Perola M, Pérusse L, Peters U, Powell JE, Power C,  
1167 Quertermous T, Rauramaa R, Reinmaa E, Ridker PM, Rivadeneira F, Rotter JI,  
1168 Saaristo TE, Saleheen D, Schlessinger D, Slagboom PE, Snieder H, Spector TD,  
1169 Strauch K, Stumvoll M, Tuomilehto J, Uusitupa M, van der Harst P, Völzke H,  
1170 Walker M, Wareham NJ, Watkins H, Wichmann H-E, Wilson JF, Zanen P,  
1171 Deloukas P, Heid IM, Lindgren CM, Mohlke KL, Speliotes EK, Thorsteinsdottir U,  
1172 Barroso I, Fox CS, North KE, Strachan DP, Beckmann JS, Berndt SI, Boehnke  
1173 M, Borecki IB, McCarthy MI, Metspalu A, Stefansson K, Uitterlinden AG, van  
1174 Duijn CM, Franke L, Willer CJ, Price AL, Lettre G, Loos RJJ, Weedon MN,  
1175 Ingelsson E, O'Connell JR, Abecasis GR, Chasman DI, Goddard ME, Visscher  
1176 PM, Hirschhorn JN, Frayling TM. 2014. Defining the role of common variation in  
1177 the genomic and biological architecture of adult human height. *Nat Genet*  
1178 **46**:1173–1186. doi:10.1038/ng.3097

1179 Xue Y, Mezzavilla M, Haber M, McCarthy S, Chen Y, Narasimhan V, Gilly A, Ayub Q,  
1180 Colonna V, Southam L, Finan C, Massaia A, Chheda H, Palta P, Ritchie G,  
1181 Asimit J, Dedoussis G, Gasparini P, Palotie A, Ripatti S, Soranzo N, Toniolo D,  
1182 Wilson JF, Durbin R, Tyler-Smith C, Zeggini E. 2017. Enrichment of low-  
1183 frequency functional variants revealed by whole-genome sequencing of multiple  
1184 isolated European populations. *Nat Commun* **8**. doi:10.1038/ncomms15927

1185 Yang Q, Kathiresan S, Lin J-P, Tofler GH, O'Donnell CJ. 2007. Genome-wide  
1186 association and linkage analyses of hemostatic factors and hematological  
1187 phenotypes in the Framingham Heart Study. *BMC Med Genet* **8**:S12.  
1188 doi:10.1186/1471-2350-8-S1-S12

1189 Zhang Q, Marioni RE, Robinson MR, Higham J, Sproul D, Wray NR, Deary IJ, McRae  
1190 AF, Visscher PM. 2018. Genotype effects contribute to variation in longitudinal  
1191 methylome patterns in older people. *Genome Med* **10**:75. doi:10.1186/s13073-  
1192 018-0585-7

1193

1194



1195 **Table 1. Final data release of WGS data for all the INGI cohorts.**

1196

<b>INGI All samples</b>				
	<b>CAR</b>	<b>FVG</b>	<b>VBI</b>	<b>INGI</b>
<b>Samples</b>	124	378	424	926
<b>Females</b>	66	220	249	535
<b>Males</b>	58	158	175	391
<b>Average coverage</b>	6.31	7.23	6.12	6.55
<b>Sites</b>	13,370,262	17,002,010	19,361,094	26,619,091
<b>Multiallelic Sites</b>	248,638	356,599	393,328	560,918
<b>SNPs</b>	12,208,629	15,521,313	17,830,208	24,557,366
<b>INDELs</b>	1,161,633	1,480,697	1,530,886	2,061,725
<b>Sites MAF &lt;= 1%</b>	3,627,622	7,283,720	9,416,028	16,685,951
<b>Sites 1% &lt; MAF &lt;= 5%</b>	3,007,162	3,069,534	3,121,545	3,125,971
<b>Sites MAF &gt; 5%</b>	6,735,478	6,648,756	6,823,521	7,123,064
<b>Singletons SNPs</b>	2,061,824	2,784,746	3,554,744	6,193,486
<b>Singletons INDELs</b>	92,372	131,275	133,156	273,679
<b>Average Heterozygosity rate per sample</b>	17.57%	13.27%	12.16%	13.34%
<b>Average Derived allele count per sample</b>	4,703,290	4,741,910	4,844,980	4,763,393
<b>Average variations per sample</b>	3,518,020	3,421,910	3,541,760	3,493,897
<b>Average INDELs per sample</b>	531,151	586,740	590,109	569,333
<b>Average singleton per sample</b>	17,285	7,671	8,646	6,925

1197