

ARTICLE

Received 19 Aug 2013 | Accepted 23 Apr 2014 | Published 13 Jun 2014

DOI: 10.1038/ncomms4934

Integrating sequence and array data to create an improved 1000 Genomes Project haplotype reference panel

Olivier Delaneau¹, Jonathan Marchini^{1,2} & The 1000 Genomes Project Consortium*

A major use of the 1000 Genomes Project (1000GP) data is genotype imputation in genome-wide association studies (GWAS). Here we develop a method to estimate haplotypes from low-coverage sequencing data that can take advantage of single-nucleotide polymorphism (SNP) microarray genotypes on the same samples. First the SNP array data are phased to build a backbone (or 'scaffold') of haplotypes across each chromosome. We then phase the sequence data 'onto' this haplotype scaffold. This approach can take advantage of relatedness between sequenced and non-sequenced samples to improve accuracy. We use this method to create a new 1000GP haplotype reference set for use by the human genetic community. Using a set of validation genotypes at SNP and bi-allelic indels we show that these haplotypes have lower genotype discordance and improved imputation performance into downstream GWAS samples, especially at low-frequency variants.

¹Department of Statistics, University of Oxford, Oxford OX1 3TG, UK. ²Wellcome Trust Centre for Human Genetics, University of Oxford, Oxford OX3 7BN, UK. Correspondence and requests for materials should be addressed to J.M. (email: marchini@stats.ox.ac.uk).

*Lists of participants and their affiliations appear at the end of the paper.

Over the last few years the use of next generation sequencing technologies has led to new insights in both population and disease genetics, by providing a more complete characterization of DNA sequences than is possible using genome-wide micro arrays. However, high coverage sequencing in large cohorts is still prohibitively expensive, and an experimental design involving low-coverage sequencing has become popular. For example, the 1000 Genomes Project (1000GP) is using $4 \times$ coverage sequencing of $\sim 2,500$ samples from a diverse set of worldwide populations¹. A consequence of the low-coverage sequencing is that some genotypes are only partially observed, and directly calling genotypes one site at a time can lead to low-quality call rates².

The current paradigm for detecting, genotyping and phasing polymorphic sites from low-coverage sequence data starts by mapping sequence reads to a reference genome. Mapped reads that overlap a given site in a single individual are then combined together to form genotype likelihoods (GLs). Genotype likelihoods are the probabilities of observing the reads given the underlying (unknown) genotypes at each site.

Improved call rates can be achieved by aggregating information across many samples through the use of phasing methods that estimate the underlying haplotypes of the study samples. Inference of the underlying haplotypes dictates the genotype calls of each sample. This builds on the idea that over small genomic regions, the samples will share haplotypes due to local genealogical relationships, leading to a per-haplotype coverage much higher than the per-individual coverage.

To achieve this haplotype phasing and genotype calling, the hidden Markov model (HMM)-based phasing methods that were primarily designed to estimate haplotypes from single-nucleotide polymorphism (SNP) array data were adapted to deal with sequencing data. For example, the 1000GP phase 1 set of haplotypes from 1,092 individuals was estimated using a combination of Beagle³ and MaCH/Thunder⁴. Such haplotype reference panels are now routinely used to impute unobserved genotypes in genome-wide association studies (GWAS), as this increases power to detect and resolve associated variants and facilitates meta-analysis⁵.

Our recent research suggests that the SHAPEIT2 method is currently the most accurate method for phasing sets of known genotypes. The method uses a similar HMM to approaches such as Impute2 (ref. 6) and MaCH. A key feature of the method is that the HMM calculations are linear in the number of haplotypes being estimated, whereas Impute2 and MaCH scale quadratically. The method uses a unique approach that represents the space of all possible haplotypes consistent with an individual's genotype data in a graphical model. A pair of haplotypes consistent with an individual's genotypes are represented as a pair of paths through this graph, with constraints to ensure consistency that are easy to apply due to the model structure. For this reason SHAPEIT2 is among the most computationally tractable methods^{7,8}.

Here we present a new version of SHAPEIT2 that estimates haplotypes from GLs generated by low-coverage sequencing data. In addition, our new method can also take advantage of SNP microarray genotypes on the same samples. The majority of the $\sim 2,500$ 1000GP sequenced samples have been genotyped on either the Illumina Omni 2.5 or Affymetrix 6.0 microarray, as well as an additional set of 1,198 unsequenced samples, many of whom are close relatives of the $\sim 2,500$ sequenced samples. Our overall approach has two steps: first the SNP array data are phased to build a backbone of haplotypes across each chromosome, which we refer to as the scaffold. Second, we take GL data at sequenced variant sites, and jointly phase this data 'onto' this haplotype scaffold.

The first advantage of this approach is that the relatedness between the extended set of genotyped samples leads to a very accurate phased scaffold. For the analysis in the paper, this set included 392 mother–father–child trios, 30 parent–child duos and 905 nominally unrelated samples. The phasing of trios and duos is expected to be highly accurate due to the Mendelian constraints on the underlying haplotypes. The phasing of the unrelated samples will benefit from being phased together with these trios and duos. The second advantage is that the phasing of the GL data onto the scaffold is carried out in chunks. As the variants in each region are phased 'onto' the scaffold, no further work is needed to combine the regions together. As such, the method is highly parallelizable. This approach generalizes our MVNcall⁹, approach which is designed to phase one variant site at a time onto a haplotype scaffold, and improves upon its accuracy, by phasing multiple sites jointly onto the scaffold and using a more sophisticated underlying model.

Our method is unique in its ability to phase GL data at multiple sites jointly, together with a phased scaffold at a subset of sites. Methods such as Beagle³ and MaCH/Thunder⁴ could be made to accept a scaffold of unphased genotypes, by recoding the genotypes as sequenced variants at very high coverage. However, our two-stage approach allows valuable family information to be used in phasing the scaffold.

Results

To demonstrate the benefits of this new method, we applied it to the 1000GP phase 1 sequence data to produce new haplotypes. We then compared these haplotypes with the existing set of 1000GP phase 1 haplotypes, and also to a set of haplotypes produced by Beagle. In all the experiments, we used the set of GLs available on the FTP website for 1,092 phase 1 samples. These consist of GLs at 36,820,992 SNPs, 1,384,273 bi-allelic indels and 14,017 structural variations (SVs). To create the haplotype scaffold (Omni 2.5 M), we used Illumina Omni 2.5 genotypes available on 2,141 samples and 2,368,234 SNPs. We phased this data set using the existing version of SHAPEIT2 (r644). Supplementary Table 1 shows the number of trios, duos and unrelated samples in each of the 14 populations. To mimic the use of a sparser haplotype scaffold, we also created a new scaffold by thinning the Omni scaffold down to 1,000,000 SNPs (1 M). We then phased the GL data set on chromosome 20 in three different ways using (a) the Omni 2.5 M scaffold, (b) the 1 M scaffold, (c) no scaffold.

We evaluated the quality of the different sets of haplotypes by looking at the concordance of the inferred genotypes to validation sets of SNP and indel genotypes. We used two validation data sets derived from Complete Genomics (CG) sequencing: a set of publicly available genotypes on 69 samples (CG1), and a larger set of 250 individuals sequenced for the purposes of 1000GP validation (CG2). Both of these data sets contain accurate genotypes that were derived from high coverage ($\sim 80 \times$), and show enough overlap in variants and samples with phase 1 for relevant genotype discordance analysis. Supplementary Tables 2 and 3 show the overlap between the CG and 1000GP data sets in terms of samples and variant sites, respectively.

Figure 1a shows the genotype discordance at CG1 SNPs. We measure discordance using just the validation genotypes that contain at least one copy of the non-reference allele (ALT) and all validation genotypes (ALL). These results show that the three haplotype sets produced by SHAPEIT2 (blue bars) have lower levels of discordance compared with Beagle haplotypes (green) and the 1000GP haplotypes (orange). For example, the CG1 ALT discordance of the SHAPEIT2 haplotypes made using the Omni 2.5 scaffold, and the ALT discordance of the 1000GP haplotypes,

are 1.03 and 1.38%, respectively. In addition, we observe that the Omni 2.5 scaffold produced better results than the 1 M scaffold, which is in turn better than using no scaffold. Figure 2a,b shows the genotype discordance at CG2 SNPs and indels, where we observe the same pattern of performance between methods. We also find that this pattern holds across different ancestries (Supplementary Fig. 1). The discordance on indels is worse than on SNPs (Fig. 2c). A reason for this difference may be that it is

more challenging to map sequencing reads that contain indels, so the GLs for indels may be less informative than GLs at SNPs.

We also used the CG samples not included in phase 1 to assess the quality of the estimated haplotypes when used as a reference panel for GWAS imputation^{5,10}. We divided the CG1 sites into those on the Illumina 1 M SNP array, and then used these together with the different haplotype sets to impute the CG1 genotypes not on the array. We then measured the imputation

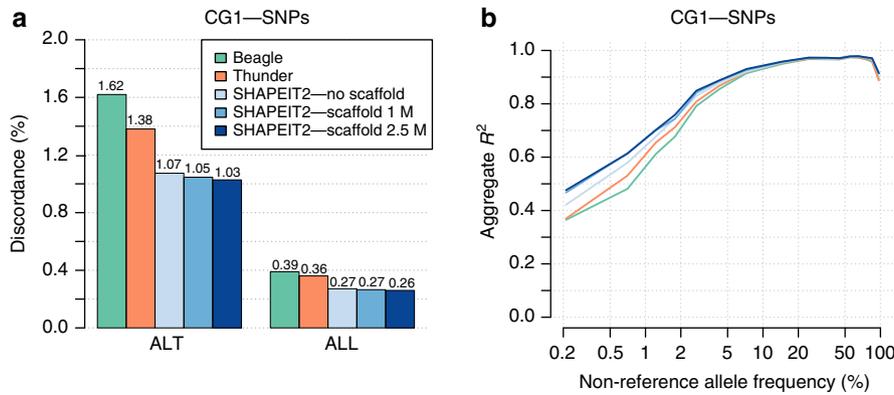


Figure 1 | Methods comparison of genotype discordance and imputation accuracy using the CG1 data. (a) Shows the discordance at chr20 CG1 SNP genotypes of Beagle (green), Thunder (orange) and SHAPEIT2 without using a scaffold (light blue), using a 1 M SNPs haplotype scaffold (medium blue) and using a 2.5 M SNPs haplotype scaffold (dark blue). ALT stands for the discordance at genotypes involving at least one non-reference allele, and ALL for the overall discordance. (b) Shows the performance of the previous call sets when used as a reference panel to impute four CG1 European samples genotyped on Illumina 1 M SNP array. The x axis shows the non-reference allele frequency of the SNP being imputed. The y axis shows imputation accuracy measure by aggregate R^2 .

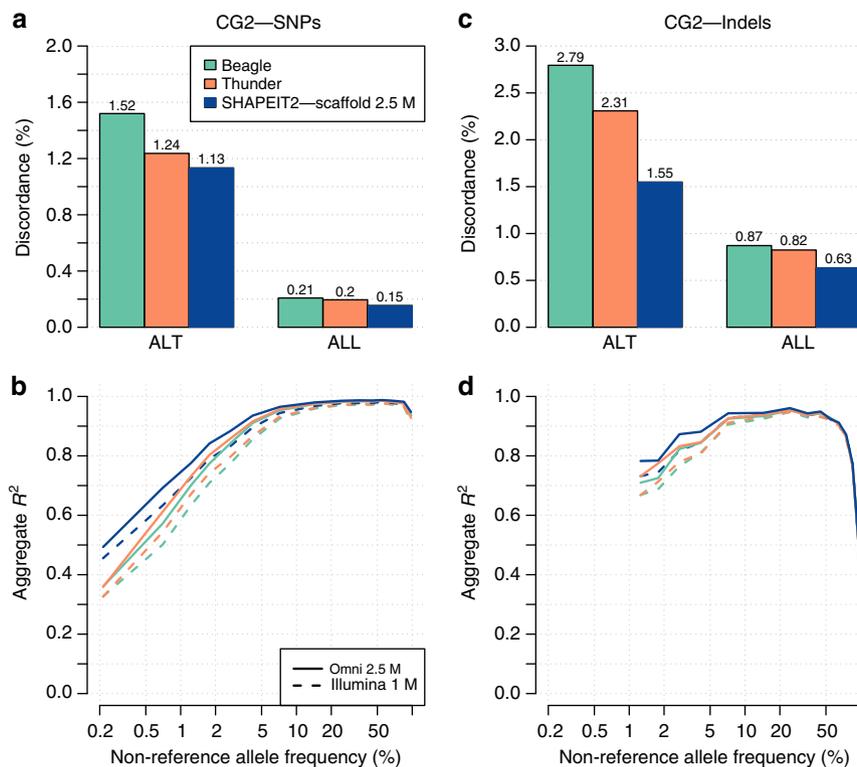


Figure 2 | Methods comparison of genotype discordance and imputation accuracy using the CG2 data. (a) Shows the whole genome genotype discordance of Beagle (green), Thunder (orange) and SHAPEIT2 using a 2.5 M SNPs haplotype scaffold (dark blue) at CG2 SNPs. (b) Shows the performance of the three call sets to impute SNPs on chromosome 10 in 10 CG2 European samples typed on Illumina 1 M and Omni 2.5 M chips. The x axis shows the non-reference allele frequency of the SNP being imputed. The y axis shows imputation accuracy measure by aggregate R^2 . (c) and (d) show similar results than a and b, respectively for short bi-allelic indels instead of SNPs.

accuracy against the CG1 genotypes. In the same way as previous evaluations¹, we stratified SNPs and indels by their non-reference allele frequency in the 1000GP haplotypes so that each site is always assigned to the same frequency bin in the results. For each SNP or indel, we measured the R^2 of the imputed dosage estimates with the validation genotypes. Figure 1b plots the non-reference allele frequency versus R^2 and shows that the use of a haplotype scaffold clearly leads to an increase in R^2 especially at lower frequencies. For example, at 0.5% frequency, the SHAPEIT2 haplotypes made with a 2.5 M scaffold increase R^2 by 0.1 compared with the 1000GP phase 1 set of haplotypes. We also find that using the 1 M scaffold produces almost identical imputation performance to the 2.5 M scaffold. Running SHAPEIT2 without a scaffold produces results intermediate to those of the scaffolded haplotypes and the 1000GP phase 1 set of haplotypes.

Figure 2c,d shows the imputation performance of SNPs and indels, respectively when using the CG2 validation set. For this experiment we carried out imputation using genotypes on the Illumina 1 M and Omni 2.5 M chip. We also observe that SHAPEIT2 haplotypes using the 2.5 M scaffold produce improved imputation performance compared with the 1000GP phase 1 set of haplotypes and the Beagle haplotypes, again independently of the sample ancestry (Supplementary Fig. 2). As expected, using a denser chip the imputation improves the results. At 1% frequency SNPs, we find that the imputation from the SHAPEIT2 scaffold reference haplotypes into genotypes on the Omni 2.5 M chip and the Illumina 1 M chip produce R^2 measures of 0.78 and 0.73, respectively. Interestingly, imputation from the 1000GP phase 1 set of haplotypes into genotypes on the Omni 2.5 M chip produces an $R^2 = 0.73$. This highlights the value of using a scaffolded set of haplotypes. In terms of imputation performance, the value of using a scaffold set of haplotypes is equivalent to the use of a much denser SNP chip in the GWAS samples.

The indel imputation results in Fig. 2d show some differences to the SNP imputation results at high frequencies, but are otherwise broadly similar. We investigated this issue and discovered that indels within 50 bp of another indel had noticeable lower imputation accuracy than more isolated indels. Figure 3 shows the imputation performance of indels stratified by

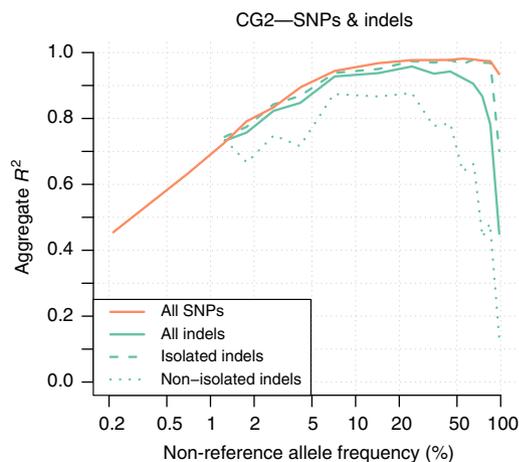


Figure 3 | Imputation accuracy at SNPs and indels using the CG2 data.

The imputation performance at SNPs and indels are shown with the orange and green lines, respectively. Performance at all indels, isolated indels and non-isolated indels are shown using plain, dashed and dotted lines. An indel is isolated when no other indels is in the 50 bp flanking regions. The x axis shows the non-reference allele frequency of the SNP being imputed. The y axis shows imputation accuracy measure by aggregate R^2 .

distance to another indel, together with the SNP imputation results. This figure shows that isolated indels can be imputed with very similar levels of accuracy to SNPs.

Discussion

Over the past year, the 1000 Genomes phase 1 haplotypes have been extensively used in many genetic studies, most of the time as reference panel to carry out GWAS imputation. In this paper, we showed that using the SHAPEIT2 phasing model, and integrating phased SNP array data, produces more accurate genotype and haplotype estimates. Using the resulting haplotypes as reference panel for GWAS imputation provides better prediction of untyped variants at rare SNPs and indels across a range of ancestries and SNP arrays. This highlights the potential of using this new set of haplotypes in future GWAS studies. The new haplotype reference set is available from the website ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/phase1/analysis_results/shapeit2_phased_haplotypes/ and our new methods are available from the website <http://www.stats.ox.ac.uk/~marchini/#software>.

We expect that many other studies may be able to make use of our approach to produce highly accurate haplotypes in their samples. It is likely that many cohorts that undergo sequencing will already have SNP microarray genotypes available. For example, twin studies that have sequenced one individual from each dizygotic twin pair, and also have genotype data on all individuals, may benefit substantially from using our approach. The phasing of the twins genotype data will be highly accurate in regions of shared haplotypes, and this will help in genotype calling and phasing of the sequence data. Studies which have sequenced one individual from parent–child pairs will benefit in a similar manner. The final version of the 1000GP haplotypes on all of the $\sim 2,500$ samples will be phased using our new approach.

We predict that further advances in haplotype accuracy are possible. First, it has recently been shown by ourselves and others that leveraging phase information in sequencing reads can lead to improved genotype calls and haplotype sets with lower switch error. In parallel work¹¹, we have extended SHAPEIT2 to utilize phase informative reads after genotypes have been called, and have shown that this improves phasing accuracy. Other authors^{12,13} have recently shown that joint inference of genotypes and haplotypes can improve both genotype and haplotype calls. However, it is yet to be determined how such improvements translate into downstream imputation accuracy. It is more likely that downstream imputation accuracy can be improved by increasing sample size of the reference panel. Efforts are now under way to create larger sets of haplotypes by combining together many low-coverage sequencing studies <http://www.haplotype-reference-consortium.org/>.

Methods

The phasing model for low-coverage sequence data. We wish to estimate the haplotypes of N unrelated individuals with sequence data at L bi-allelic variants, which could be either SNPs, indels or structural variants. Our new algorithm extends the SHAPEIT2 model and the Markov chain Monte Carlo (MCMC) method used to carry out inference from this model. We use a Gibbs sampling scheme in which each individual's haplotypes are sampled conditional upon the sequence reads of the individual and the current estimates of all the other individuals. Thus it is sufficient for us to consider the details of a single iteration in which we update the haplotypes of the i th individual. We use R to denote the sequence data available for this individual and H to denote the current haplotype estimates of other individuals being used in the iteration. We define the genotype likelihood as the probability of observing the sequence data R at a particular site l given the unobserved genotype G_l ; $P(R|G_l)$, where $G_l = 0, 1, 2$ counts the number of non-reference alleles in the genotype. These GLs can be obtained using specialised software like SAMtools¹⁴, SNPtools¹⁵ or GATK¹⁶ that derive these likelihoods directly from the BAM files containing the sequence reads.

In each iteration we must sample a pair of haplotypes (h_1, h_2) for the i th individuals given both R and H . To do so, we adapted the parsimonious representation of the possible haplotypes of SHAPEIT to deal with GLs. We divide

the region being phased into a number, C , of consecutive non-overlapping segments such that each segment contains eight possible haplotypes consistent with the GLs. In the case of bi-allelic variants, it means that each segment spans three sites, and we will see in the next section how this number can be increased. We use $S_l \in \{1, \dots, C\}$ to denote the segment that contains the l th SNP and b_s and e_s to denote the first site and the last site included in the s th segment, respectively. We use A_{lb} to denote the allele carried at the l th site by the b th consistent haplotype. We can now represent a possible haplotype as a vector of labels $X = \{X_1, \dots, X_L\}$ where X_l denotes the label of the haplotype at the l th site in the S_l th segment. The segmentation implies that the labels are identical within each segment so that we always have $X_l = X_{l-1}$ when $S_l = S_{l-1}$. We use $X_{\{s\}}$ to define the label of the haplotype across all sites residing in the s th segment. Moreover, we represent a pair of haplotypes as a pair of vectors of labels (X^1, X^2) . An illustration of this graph representation of the possible haplotypes can be seen in Supplementary Fig. 3a.

Given the segment representation described above, sampling a diplotype (pair of haplotypes) given a set of known haplotypes H and a set of sequencing reads R involves sampling from the posterior distribution $Pr(X^1, X^2 | H, R)$. By assuming first that the reads for the individual we are updating, R , are conditionally independent of the haplotypes in other individuals, H , given the pair of haplotypes (X^1, X^2) we can write

$$P(X^1, X^2 | H, R) \propto P(X^1, X^2, R, H) \tag{1}$$

$$\propto P(R | X^1, X^2) P(X^1, X^2 | H) \tag{2}$$

This factorization involves a model of the diplotype given the observed haplotypes, $P(X^1, X^2 | H)$ and for this we use the previously described SHAPEIT2 model⁸. The term $P(R | X^1, X^2)$ is constructed from the GLs.

On the basis of the segmentation of the chromosome into C segments, we employ a similar Markov model as the one introduced in the SHAPEIT2 method⁸. It can be written as:

$$P(X^1, X^2 | H, R) = P(X^1_{\{1\}}, X^2_{\{1\}} | H, R) \prod_{s=2}^C P(X^1_{\{s\}}, X^2_{\{s\}} | X^1_{\{s-1\}}, X^2_{\{s-1\}}, H, R) \tag{3}$$

The idea here is to sample first a diplotype for the first segment $s = 1$ from $P(X^1_{\{1\}}, X^2_{\{1\}} | H, R)$ and then for each successive segment from $P(X^1_{\{s\}}, X^2_{\{s\}} | X^1_{\{s-1\}}, X^2_{\{s-1\}}, H, R)$. The scheme we use is described by the following steps:

1. A pair of haplotypes in the first segment with labels (i, j) is sampled with probability proportional to $P(X^1_i = i, X^2_j = j | H, R)$.
2. While $s \leq C$ a pair of haplotypes (d, f) for the s th segment is sampled given the previously sampled pair (i, j) for the $\{s-1\}$ th segment with probability proportional to $P(X^1_d = d, X^2_f = f | X^1_{\{s-1\}} = i, X^2_{\{s-1\}} = j, H, R)$.
3. Set $s = s + 1$.
4. If $s = C + 1$ then stop, else go to step 2.

The result is a pair of vectors of haplotype labels, X^1 and X^2 , across the whole region being phased and these can be turned into new haplotype estimates, (h_1, h_2) , using $h_{il} = A_{lX^i}$ for $i \in \{1, 2\}$. These haplotype estimates can then be added back into the haplotype set H and the next individual's haplotypes can be estimated, although their current haplotype estimates must be removed from H first.

To carry out this Markov-based sampling, we need now to describe how to obtain the two distributions $P(X^1_i = i, X^2_j = j | H, R)$ and $P(X^1_d = d, X^2_f = f | X^1_{\{s-1\}} = i, X^2_{\{s-1\}} = j, H, R)$. To do so, we decompose them by using equations (1) and (2) as follows:

$$P(X^1_{\{1\}}, X^2_{\{1\}} | H, R) = P(R | X^1_{\{1\}}, X^2_{\{1\}}) P(X^1_{\{1\}}, X^2_{\{1\}} | H) \tag{4}$$

$$P(X^1_{\{s\}}, X^2_{\{s\}} | X^1_{\{s-1\}}, X^2_{\{s-1\}}, H, R) \propto P(R | X^1_{\{s\}}, X^2_{\{s\}}, X^1_{\{s-1\}}, X^2_{\{s-1\}} | H, R) \tag{5}$$

$$\propto P(R | X^1_{\{s\}}, X^2_{\{s\}}, X^1_{\{s-1\}}, X^2_{\{s-1\}}) \times P(X^1_{\{s\}}, X^2_{\{s\}}, X^1_{\{s-1\}}, X^2_{\{s-1\}} | H)$$

We use the SHAPEIT2 model for the terms $P(X^1_{\{1\}}, X^2_{\{1\}} | H)$ and $P(X^1_{\{s\}}, X^2_{\{s\}}, X^1_{\{s-1\}}, X^2_{\{s-1\}} | H)$. We do not give more details here since a complete description can be found in the SHAPEIT2 paper⁸. The GLs enter the model in the term $P(R | X^1, X^2)$ as a product over all L sites as

$$P(R | X^1, X^2) = \prod_{l=1}^L P(R | G_l = A_{lX^1} + A_{lX^2}) \tag{6}$$

which implies that

$$P(R | X^1_{\{1\}}, X^2_{\{1\}}) = \prod_{l=b_1}^{e_1} P(R | X^1_l, X^2_l) \tag{7}$$

$$P(R | X^1_{\{s\}}, X^2_{\{s\}}, X^1_{\{s-1\}}, X^2_{\{s-1\}}) = \prod_{l=b_{s-1}}^{e_s} P(R | X^1_l, X^2_l) \tag{8}$$

Initialization and MCMC iterations. The experience of the 1000GP analysis group is that phasing approaches based on HMMs such as Thunder and Impute2

are slow to converge when applied to low-coverage sequence data if the starting haplotype estimates are initialized randomly. It has been observed that the Beagle method does not have this property, and that Thunder and Impute2 benefit from using an initial set of haplotypes estimated via Beagle. The 1000GP phase 1 haplotypes were estimated in this way by first running Beagle and then using these haplotypes as initial estimates in the Thunder model¹.

We initialize some of the genotypes by using the genotype posteriors $P(G_l | H, R)$ provided by the Beagle phasing model. Our approach relies on fixing the genotypes with high posterior probabilities and then use our model to call all the remaining genotypes (Supplementary Fig. 3b). Fixing highly confident genotypes is beneficial as it implies additional constraints on the space of possible haplotypes. In practice, segments then tend to contain more sites than in the default model: 32 sites on average per segment when applied to 1000GP instead of only three sites if no genotypes are fixed.

We empirically determined a threshold on the Beagle posteriors to fix genotypes while maintaining relatively low discordance rates. This approach relies on the Beagle posteriors being well calibrated. To do so, we defined a set of 23 different threshold values ranging from 0.5 to 0.999 and measured for each (1) the discordance between CG1 and genotypes with a posterior above the threshold and (2) the percentage of genotypes with posteriors falling below the threshold (Supplementary Fig. 4a,b). In addition, we also measured the proportion of discordances of the full Beagle call set falling below each threshold value (Supplementary Fig. 4c,d). From this experiment, we empirically determined that a threshold value of 0.995 gives good performance: it implies that around 97% of the genotypes can be directly fixed while maintaining a discordance against CG1 of 0.07% overall (ALL) and of 0.25% at genotypes involving at least one alternative allele (ALT). We find that the 3% of the genotypes that we choose not to fix contain over 80% of the genotypes found to be discordant. Thus it makes sense that these are the genotypes that we try to improve upon using our model.

Our algorithm starts from the haplotype estimates produced by Beagle and then, each MCMC iteration consists of updating the haplotypes of each sample conditional upon a set of other haplotypes using the Markov model described in section A. Our algorithm for GLs follows an iteration scheme quite different than in the SHAPEIT2 algorithm described in Delaneau *et al.* (2012). Specifically, we carry out several stages of pruning and merging iterations, instead of a single set of pruning and merging. In practice, we use 12 stages of four iterations ($= 48$ iterations). We do not use burn-in iterations as we already have an initial estimate provided by Beagle. Each pruning and merging stage is used to remove unlikely states and transitions from the Markov model that describes the space of haplotypes with each individual. When enough transitions are pruned we merge adjacent segments together. This has the effect of simplifying the space of possible haplotypes so that a final set of sampling iterations can be carried out more efficiently. In practice, as we multiply these pruning and merging stages, the size of the model (that is, the graphs) tend to converge as shown by the evolutions of the number of sites per segment (Supplementary Fig. 5a) and the total number of segments (Supplementary Fig. 5b).

Finally, to complete the model, we only use a subset of all available haplotypes when updating each individual as done in SHAPEIT2. We used a carefully chosen subset containing $K_1 = 400$ haplotypes that most closely match the haplotypes of the individual being updated¹⁰. Note that the haplotype matching is carried out on overlapping windows of size $W = 0.1$ Mb. Moreover, we also found useful to use an additional set of $K_2 = 200$ randomly chosen haplotypes to help the mixing of the MCMC. So in total, we used $K = 600$ conditioning haplotypes. Using such a large number of conditioning haplotypes is facilitated as SHAPEIT2 has linear complexity with K .

Using a haplotype scaffold. We denote as F the pair of haplotypes derived from SNP array for the i th individual, now the goal is to sample a pair of haplotypes from $P(X^1, X^2 | H, R, F)$ such that they are fully consistent with F . The scaffold F imposes a set of hard constraints on the space of possible haplotypes generated by the sampling scheme as illustrated in Supplementary Fig. 3c. So in the first segment $s = 1$: $P(X^1_{\{1\}}, X^2_{\{1\}} | H, R, F) = P(X^1_{\{1\}}, X^2_{\{1\}} | H, R)$ when the pair of haplotypes defined by $(X^1_{\{1\}}, X^2_{\{1\}})$ is fully consistent with F over the first segment, and 0 otherwise. Similarly, we define

$$P(X^1_{\{s\}}, X^2_{\{s\}} | X^1_{\{s-1\}}, X^2_{\{s-1\}}, H, R, F) = P(X^1_{\{s\}}, X^2_{\{s\}} | X^1_{\{s-1\}}, X^2_{\{s-1\}}, H, R) \tag{9}$$

when the haplotype pair defined by $(X^1_{\{s\}}, X^2_{\{s\}}, X^1_{\{s-1\}}, X^2_{\{s-1\}})$ is fully consistent with F over the segments s and $s-1$, and 0 otherwise. In practice, setting some of the transition probabilities that are inconsistent with F to 0 between successive segments means that it becomes impossible to sample haplotypes inconsistent with F across the full set of L sites.

1000GP phase 1 low-coverage sequence data. We downloaded the GLs for 1,092 1000GP samples from <ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20110521/>. This data set contains GLs for 36,820,992 SNPs, 1,384,273 short bi-allelic indels and 14,017 SVs. The GLs for SNPs were computed using SNPtools¹⁵, for indels using (ref. 16) and SVs using (ref. 17). We ran Beagle and SHAPEIT2 on

the whole genome in chunks of 1.4 Mb with 0.2 Mb overlaps between flanking chunks.

Beagle was run using 20 iterations instead of the 10 by default, otherwise, all other default settings were used. SHAPEIT2 was run using 78 iterations: 12 stages of 4 pruning iterations plus 30 main iterations. The estimation was carried out in windows of size $W = 0.1$ Mb, using $k = 600$ conditioning haplotypes; 400 chosen by Hamming distance and 200 chosen at random. All these computations were done using an $\sim 1,000$ CPU nodes cluster. SHAPEIT2 and Beagle required ~ 289 and ~ 99 CPU months, respectively to phase the whole genome 1000GP phase 1 data set.

The multi-threading property of SHAPEIT2 proved to be very convenient on clusters with low memory nodes (for example, only 2–3 Gb of RAM per CPU core). For instance, on a single 8 CPU node, it is much more memory efficient to phase with SHAPEIT2 eight chunks of data sequentially each using eight threads than running the eight chunks in parallel. Both strategies need roughly the same running times whereas the second requires sharing of memory between the eight chunks.

1000GP Illumina Omni 2.5 SNP array data. For the haplotype scaffold, we used a set of 2,141 samples genotyped on Illumina Omni 2.5 M. This set of samples includes all the 1000GP phase 1 samples. This data set contains some parent–child duos and mother–father–child trios, and in some cases just a subset of each family has been sequenced. Supplementary Table 1 gives details of sequenced and non-sequenced samples. We found that 380 and 30 phase 1 1000GP sequenced samples are part of trios and duos in this data set. SNPs with a missing data rate above 10% and a Mendel error rate above 5% were removed, leaving a total of 2,368,234 SNPs ready for phasing. We phased this data using SHAPEIT2 (r644) using all default settings ($W = 2$ Mb, $K = 100$ haplotypes, iterations = 45) and using all available family information. We used the resulting haplotypes as a scaffold to call the variant sites in 1000GP. The whole genome overlap between both data sets contains 2,183,314 SNPs.

Complete Genomics (CG) validation data. As validation data, we used two different data sets: the 69 genomes from Complete Genomics (CG1) and an additional set of 250 samples (CG2) also sequenced by CG. All these samples were sequenced using the Complete Genomics sequencing technology at an average of $80\times$. The CG1 can be found at <http://www.completegenomics.com/public-data/69-Genomes/> and the CG2 at ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/working/20130524_cgi_combined_calls/. On these data sets, we filtered out all variants with a call rate below 66% and ignored them in all posterior validation analysis. In both the data sets, we used called SNPs as validations. We found 15,060,295 and 17,399,956 1000GP SNPs overlapping CG1 and CG2, respectively. In addition, we found 554,886 1000GP indels also in CG2.

In terms of sample overlap with 1000GP, CG1 and CG2 contain 34 and 125 samples, respectively. We used genotypes of these samples to measure discordance with the 1000GP call sets. As CG genotypes were derived from an average coverage of $80\times$, we assume that they are accurate and thus can be considered as the truth in the validation process. We define the discordance as being the percentage of these CG genotypes that are miscalled by a software (Beagle, Thunder or SHAPEIT). We measure both the overall (ALL) discordance and the discordance at genotypes with at least one non-reference allele (ALT). In all discordance measures, we systematically exclude all genotypes at SNPs included in the Omni 2.5 M chips.

We also used CG samples that are not in 1000GP nor related with any samples in 1000GP to assess the performance of the various call sets when used as reference panels for imputation. In CG1, we found 20 such samples, and 51 in CG2. To mimic a standard GWAS, we extracted genotypes at subsets of SNPs in both the data sets: for CG1, at all SNPs on chromosome 20 also included in the Illumina 1 M chip for CG1 (set A), and for CG2, at all SNPs on chromosome 10 also included in the Illumina 1 M (set B) and Illumina Omni 2.5 M (set C) chips. We then imputed all remaining CG SNP genotypes available using Impute2 (default parameters) and the various call sets as reference panels. We imputed 315,326 SNPs from set A, 823,570 SNPs and 27,511 indels from set B, and 775,818 SNPs and 27,511 indels from set C. We defined as isolated, an indel with no other indel in the 50 bp flanking regions. We found 23,641 (85.9%) isolated indels and 3,870 (14.1%) non-isolated indels. All these variants were then classified into frequency bins that were derived from the official release of haplotypes on a per continental group basis as defined in Supplementary Table 2. Then, for each continental group and frequency bin separately, we measured the squared Pearson correlation coefficient between the true (CG derived) and the imputed dosages, ranging from 0 in case of completely wrong imputation to 1 in the case of a perfect imputation. Note that a genotype dosage is the expected number of copies of non-reference alleles; being 0, 1 or 2 in the case of a known genotype and ranging from 0 to 2 in the case of an

imputed genotype. Indels in the phase 1 1000GP haplotypes were filtered at 1% which explains why there are no results for very low-frequency indels in Fig. 2d.

References

1. The 1000 Genomes Project Consortium. An integrated map of genetic variation from 1,092 human genomes. *Nature* **491**, 56–65 (2012).
2. Nielsen, R., Paul, J. S., Albrechtsen, A. & Song, Y. S. Genotype and SNP calling from next-generation sequencing data. *Nat. Rev. Genet.* **12**, 443–451 (2011).
3. Browning, B. & Browning, S. A unified approach to genotype imputation and haplotype-phase inference for large data sets of trios and unrelated individuals. *Am. J. Hum. Genet.* **84**, 210–223 (2009).
4. Li, Y., Willer, C. J., Ding, J., Scheet, P. & Abecasis, G. R. MaCH: using sequence and genotype data to estimate haplotypes and unobserved genotypes. *Genet. Epidemiol.* **34**, 816–834 (2010).
5. Marchini, J. & Howie, B. Genotype imputation for genome-wide association studies. *Nat. Rev. Genet.* **11**, 499–511 (2010).
6. Howie, B. N., Donnelly, P. & Marchini, J. A flexible and accurate genotype imputation method for the next generation of genome-wide association studies. *PLoS Genet.* **5**, e1000529 (2009).
7. Delaneau, O., Marchini, J. & Zagury, J.-F. A linear complexity phasing method for thousands of genomes. *Nat. Methods* **9**, 179–181 (2012).
8. Delaneau, O., Zagury, J.-F. & Marchini, J. Improved whole-chromosome phasing for disease and population genetic studies. *Nat. Methods* **10**, 5–6 (2013).
9. Menelaou, A. & Marchini, J. Genotype calling and phasing using next-generation sequencing reads and a haplotype scaffold. *Bioinformatics* **29**, 84–91 (2013).
10. Howie, B., Marchini, J. & Stephens, M. Genotype imputation with thousands of genomes. *G3 (Bethesda)* **1**, 457–470 (2011).
11. Delaneau, O., Howie, B., Cox, A. J., Zagury, J.-F. & Marchini, J. Haplotype estimation using sequencing reads. *Am. J. Hum. Genet.* **93**, 687–696 (2013).
12. Zhang, K. & Zhi, D. Joint haplotype phasing and genotype calling of multiple individuals using haplotype informative reads. *Bioinformatics* **29**, 2427–2434 (2013).
13. Yang, W. *et al.* Leveraging reads that span multiple single nucleotide polymorphisms for haplotype inference from sequencing data. *Bioinformatics* **29**, 2245–2252 (2013).
14. Li, H. A statistical framework for SNP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. *Bioinformatics* **27**, 2987–2993 (2011).
15. Wang, Y., Lu, J., Yu, J., Gibbs, R. A. & Yu, F. An integrative variant analysis pipeline for accurate genotype/haplotype inference in population NGS data. *Genome Res.* **23**, 833–842 (2013).
16. DePristo, M. A. *et al.* A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nat. Genet.* **43**, 491–498 (2011).
17. Handsaker, R. E., Korn, J. M., Nemesh, J. & McCarroll, S. A. Discovery and genotyping of genome structural polymorphism by sequencing on a population scale. *Nat. Genet.* **43**, 269–276 (2011).

Acknowledgements

J.M. and O.D. acknowledge support from the Medical Research Council (G0801823). We thank Androniki Menelaou, Bryan Howie and members of the 1000 Genomes analysis group for their comments.

Author contributions

O.D. and J.M. designed and performed the research. J.M. supervised the research. J.M. and O.D. wrote the paper. The 1000 Genomes Project Consortium provided data.

Additional information

Supplementary Information accompanies this paper at <http://www.nature.com/naturecommunications>

Competing financial interests: The authors declare no competing financial interests.

Reprints and permission information is available online at <http://npg.nature.com/reprintsandpermissions/>

How to cite this article: Delaneau, O. *et al.* Integrating sequence and array data to create an improved 1000 Genomes Project haplotype reference panel. *Nat. Commun.* **5**:3934 doi: 10.1038/ncomms4934 (2014).

Gil A. McVean^{1,2}, Peter Donnelly^{1,2}, Gerton Lunter¹, Jonathan L. Marchini^{1,2}, Simon Myers^{1,2}, Anjali Gupta-Hinch¹, Zamin Iqbal¹, Iain Mathieson¹, Andy Rimmer¹, Dionysia K. Xifara^{1,2}, Angeliki Kerasidou¹, Claire Churchhouse², Olivier Delaneau², David M. Altshuler^{3,4,5}, Stacey B. Gabriel³, Eric S. Lander³, Namrata Gupta³, Mark J. Daly³, Mark A. DePristo³, Eric Banks³, Gaurav Bhatia³, Mauricio O. Carneiro³, Guillermo del Angel³, Giulio Genovese³, Robert E. Handsaker^{3,5}, Chris Hartl³, Steven A. McCarroll³, James C. Nemeshegyi³, Ryan E. Poplin³, Stephen F. Schaffner³, Khalid Shakir³, Pardis C. Sabeti^{3,39}, Sharon R. Grossman^{3,39}, Shervin Tabrizi^{3,39}, Ridhi Tariyal^{3,39}, Heng Li^{3,6}, David Reich⁵, Richard M. Durbin⁶, Matthew E. Hurles⁶, Senduran Balasubramaniam⁶, John Burton⁶, Petr Danecek⁶, Thomas M. Keane⁶, Anja Kolb-Kokocinski⁶, Shane McCarthy⁶, James Stalker⁶, Michael Quail⁶, Qasim Ayub⁶, Yuan Chen⁶, Alison J. Coffey⁶, Vincenza Colonna^{6,86}, Ni Huang⁶, Luke Jostins⁶, Aylwyn Scally⁶, Klaudia Walter⁶, Yali Xue⁶, Yujun Zhang⁶, Ben Blackburne⁶, Sarah J. Lindsay⁶, Zemin Ning⁶, Adam Frankish⁶, Jennifer Harrow⁶, Chris Tyler-Smith⁶, Gonalo R. Abecasis⁷, Hyun Min Kang⁷, Paul Anderson⁷, Tom Blackwell⁷, Fabio Busonero^{7,69,71}, Christian Fuchsberger⁷, Goo Jun⁷, Andrea Maschio^{7,69,71}, Eleonora Porcu^{7,69,71}, Carlo Sidore^{7,69,71}, Adrian Tan⁷, Mary Kate Trost⁷, David R. Bentley⁸, Russell Grocock⁸, Sean Humphray⁸, Terena James⁸, Zoya Kingsbury⁸, Markus Bauer⁸, R. Keira Cheetham⁸, Tony Cox⁸, Michael Eberle⁸, Lisa Murray⁸, Richard Shaw⁸, Aravinda Chakravarti⁹, Andrew G. Clark¹⁰, Alon Keinan¹⁰, Juan L. Rodriguez-Flores¹⁰, Francisco M. De La Vega¹⁰, Jeremiah Degenhardt¹⁰, Evan E. Eichler¹¹, Paul Flicek¹², Laura Clarke¹², Rasko Leinonen¹², Richard E. Smith¹², Xiangqun Zheng-Bradley¹², Kathryn Beal¹², Fiona Cunningham¹², Javier Herrero¹², William M. McLaren¹², Graham R. S. Ritchie¹², Jonathan Barker¹², Gavin Kelman¹², Eugene Kulesha¹², Rajesh Radhakrishnan¹², Asier Roa¹², Dmitriy Smirnov¹², Ian Streeter¹², Iliana Toneva¹², Richard A. Gibbs¹³, Huyen Dinh¹³, Christie Kovar¹³, Sandra Lee¹³, Lora Lewis¹³, Donna Muzny¹³, Jeff Reid¹³, Min Wang¹³, Fuli Yu¹³, Matthew Bainbridge¹³, Danny Challis¹³, Uday S. Evani¹³, James Lu¹³, Uma Nagaswamy¹³, Aniko Sabo¹³, Yi Wang¹³, Jin Yu¹³, Gerald Fowler¹³, Walker Hale¹³, Divya Kalra¹³, Eric D. Green¹⁴, Bartha M. Knoppers¹⁵, Jan O. Korbel¹⁶, Tobias Rausch¹⁶, Adrian M. Stutz¹⁶, Charles Lee¹⁷, Lauren Griffin¹⁷, Chih-Heng Hsieh¹⁷, Ryan E. Mills^{17,33}, Marcin von Grotthuss¹⁷, Chengsheng Zhang¹⁷, Xinghua Shi¹⁸, Hans Lehrach^{19,20}, Ralf Sudbrak¹⁹, Vyacheslav S. Amstislavskiy¹⁹, Matthias Lienhard¹⁹, Florian Mertes¹⁹, Marc Sultan¹⁹, Bernd Timmermann¹⁹, Marie-Laure Yaspo¹⁹, Sudbrak, Ralf Herwig¹⁹, Elaine R. Mardis²¹, Richard K. Wilson²¹, Lucinda Fulton²¹, Robert Fulton²¹, George M. Weinstock²¹, Asif Chinwalla²¹, Li Ding²¹, David Dooling²¹, Daniel C. Koboldt²¹, Michael D. McLellan²¹, John W. Wallis²¹, Michael C. Wendt²¹, Qunyuan Zhang²¹, Gabor T. Marth²², Erik P. Garrison²², Deniz Kural²², Wan-Ping Lee²², Wen Fung Leong²², Alistair N. Ward²², Jiantao Wu²², Mengyao Zhang²², Deborah A. Nickerson²³, Can Alkan^{23,82}, Fereydoun Hormozdizadeh²³, Arthur Ko²³, Peter H. Sudmant²³, Jeanette P. Schmidt²⁴, Christopher J. Davies²⁴, Jeremy Gollub²⁴, Teresa Webster²⁴, Brant Wong²⁴, Yiping Zhan²⁴, Stephen T. Sherry²⁵, Chunlin Xiao²⁵, Deanna Church²⁵, Victor Ananiev²⁵, Zinaida Belaia²⁵, Dimitriy Beloslyudtsev²⁵, Nathan Bouk²⁵, Chao Chen²⁵, Robert Cohen²⁵, Charles Cook²⁵, John Garner²⁵, Timothy Hefferon²⁵, Mikhail Kimelman²⁵, Chunlei Liu²⁵, John Lopez²⁵, Peter Meric²⁵, Yuri Ostapchuk²⁵, Lon Phan²⁵, Sergiy Ponomarev²⁵, Valerie Schneider²⁵, Eugene Shekhtman²⁵, Karl Sirotkin²⁵, Douglas Slotka²⁵, Hua Zhang²⁵, Jun Wang^{26,27,28,29}, Xiaodong Fang²⁶, Xiaosen Guo²⁶, Min Jian²⁶, Hui Jiang²⁶, Xin Jin²⁶, Guoqing Li²⁶, Jingxiang Li²⁶, Yingrui Li²⁶, Xiao Liu²⁶, Yao Lu²⁶, Xuedi Ma²⁶, Shuaishuai Tai²⁶, Meifang Tang²⁶, Bo Wang²⁶, Guangbiao Wang²⁶, Honglong Wu²⁶, Renhua Wu²⁶, Ye Yin²⁶, Wenwei Zhang²⁶, Jiao Zhao²⁶, Meiru Zhao²⁶, Xiaole Zheng²⁶, Lachlan J.M. Coin²⁶, Lin Fang²⁶, Qibin Li²⁶, Zhenyu Li²⁶, Haoxiang Lin²⁶, Binghang Liu²⁶, Ruibang Luo²⁶, Haojing Shao²⁶, Bingqiang Wang²⁶, Yinlong Xie²⁶, Chen Ye²⁶, Chang Yu²⁶, Hancheng Zheng²⁶, Hongmei Zhu²⁶, Hongyu Cai²⁶, Hongzhi Cao²⁶, Yeyang Su²⁶, Zhongming Tian²⁶, Huanming Yang^{26,29,30}, Ling Yang²⁶, Jiayong Zhu²⁶, Zhiming Cai²⁶, Jian Wang²⁶, Marcus W. Albrecht³¹, Tatiana A. Borodina³¹, Adam Auton³², Seungtae C. Yoon³⁴, Jayon Lihm³⁴, Vladimir Makarov³⁵, Hanjun Jin³⁶, Wook Kim³⁷, Ki Cheol Kim³⁷, Srikanth Gottipati³⁸, Danielle Jones³⁸, David N. Cooper⁴⁰

Edward V. Ball⁴⁰, Peter D. Stenson⁴⁰, Bret Barnes⁴¹, Scott Kahn⁴¹, Kai Ye⁴², Mark A. Batzer⁴³, Miriam K. Konkel⁴³, Jerilyn A. Walker⁴³, Daniel G. MacArthur⁴⁴, Monkol Lek⁴⁴, Mark D. Shriver⁴⁵, Carlos D. Bustamante⁴⁶, Simon Gravel⁴⁶, Eimear E. Kenny⁴⁶, Jeffrey M. Kidd⁴⁶, Phil Lacroute⁴⁶, Brian K. Maples⁴⁶, Andres Moreno-Estrada⁴⁶, Fouad Zakharia⁴⁶, Brenna Henn⁴⁶, Karla Sandoval⁴⁶, Jake K. Byrnes⁴⁷, Eran Halperin^{48,49,50}, Yael Baran⁴⁸, David W. Craig⁵¹, Alexis Christoforides⁵¹, Tyler Izatt⁵¹, Ahmet A. Kurdoglu⁵¹, Shripad A. Sinari⁵¹, Nils Homer⁵², Kevin Squire⁵³, Jonathan Sebat^{54,55}, Vineet Bafna⁵⁶, Kenny Ye⁵⁷, Esteban G. Burchard⁵⁸, Ryan D. Hernandez⁵⁸, Christopher R. Gignoux⁵⁸, David Haussler^{59,60}, Sol J. Katzman⁵⁹, W. James Kent⁵⁹, Bryan Howie⁶¹, Andres Ruiz-Linares⁶², Emmanouil T. Dermitzakis^{63,64,65}, Tuuli Lappalainen^{63,64,65}, Scott E. Devine⁶⁶, Xinyue Liu⁶⁶, Ankit Maroo⁶⁶, Luke J. Tallon⁶⁶, Jeffrey A. Rosenfeld^{67,68}, Leslie P. Michelson^{67,68}, Andrea Angius⁶⁹, Francesco Cucca^{69,71}, Serena Sanna⁶⁹, Abigail Bigham⁷⁰, Chris Jones⁷², Fred Reinier⁷², Yun Li⁷³, Robert Lyons⁷⁴, David Schlessinger⁷⁵, Philip Awadalla⁷⁶, Alan Hodgkinson⁷⁶, Taras K. Oleksyk⁷⁷, Juan C. Martinez-Cruzado⁷⁷, Yunxin Fu⁷⁸, Xiaoming Liu⁷⁸, Momiao Xiong⁷⁸, Lynn Jorde⁷⁹, David Witherspoon⁷⁹, Jinchuan Xing⁸⁰, Brian L. Browning⁸¹, Iman Hajirasouliha⁸³, Ken Chen⁸⁴, Cornelis A. Albers⁸⁵, Mark B. Gerstein^{87,88,89}, Alexej Abyzov^{87,89}, Jieming Chen⁸⁷, Yao Fu⁸⁷, Lukas Habegger⁸⁷, Arif O. Harmanci⁸⁷, Xinmeng Jasmine Mu⁸⁷, Cristina Sisu⁸⁷, Suganthi Balasubramanian⁸⁹, Mike Jin⁸⁹, Ekta Khurana⁸⁹, Declan Clarke⁹⁰, Jacob J. Michaelson⁹¹, Chris OSullivan⁹², Kathleen C. Barnes⁹³, Neda Gharani⁹⁴, Lorraine H. Toji⁹⁴, Norman Gerry⁹⁴, Jane S. Kaye⁹⁵, Alastair Kent⁹⁶, Rasika Mathias⁹⁷, Pilar N. Ossorio^{98,99}, Michael Parker¹⁰⁰, Charles N. Rotimi¹⁰¹, Charmaine D. Royal¹⁰², Sarah Tishkoff¹⁰³, Marc Via¹⁰⁴, Walter Bodmer¹⁰⁵, Gabriel Bedoya¹⁰⁶, Gao Yang¹⁰⁷, Chu Jia You¹⁰⁸, Andres Garcia-Montero¹⁰⁹, Alberto Orfao¹¹⁰, Julie Dutil¹¹¹, Lisa D. Brooks¹¹², Adam L. Felsenfeld¹¹², Jean E. McEwen¹¹², Nicholas C. Clegg¹¹², Mark S. Guyer¹¹², Jane L. Peterson¹¹², Audrey Duncanson¹¹³, Michael Dunn¹¹³ and Leena Peltonen[‡]

¹Wellcome Trust Centre for Human Genetics, Oxford University, Oxford OX3 7BN, UK; ²Department of Statistics, Oxford University, Oxford OX1 3TG, UK; ³The Broad Institute of MIT and Harvard, 7 Cambridge Center, Cambridge, Massachusetts 02142, USA; ⁴Center for Human Genetic Research, Massachusetts General Hospital, Boston, Massachusetts 02114, USA; ⁵Department of Genetics, Harvard Medical School, Cambridge, Massachusetts 02142, USA; ⁶Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Cambridge CB10 1SA, UK; ⁷Center for Statistical Genetics, Biostatistics, University of Michigan, Ann Arbor, Michigan 48109, USA; ⁸Illumina United Kingdom, Chesterford Research Park, Little Chesterford, Near Saffron Walden, Essex CB10 1XL, UK; ⁹McKusick-Nathans Institute of Genetic Medicine, Johns Hopkins University School of Medicine, Baltimore, Maryland 21205, USA; ¹⁰Center for Comparative and Population Genomics, Cornell University, Ithaca, New York 14850, USA; ¹¹Department of Genome Sciences, University of Washington School of Medicine and Howard Hughes Medical Institute, Seattle, Washington 98195, USA; ¹²European Bioinformatics Institute, Wellcome Trust Genome Campus, Cambridge CB10 1SD, UK; ¹³Brendan Vaughan Baylor College of Medicine, Human Genome Sequencing Center, Houston, Texas 77030, USA; ¹⁴US National Institutes of Health, National Human Genome Research Institute, 31 Center Drive, Bethesda, Maryland 20892, USA; ¹⁵Centre of Genomics and Policy, McGill University, Montreal, Quebec, Canada H3A 1A4; ¹⁶European Molecular Biology Laboratory, Genome Biology Research Unit, Meyerhofstrae 1, 69117 Heidelberg, Germany; ¹⁷Department of Pathology, Brigham and Women's Hospital and Harvard Medical School, Boston, Massachusetts 02115, USA; ¹⁸Department of Bioinformatics and Genomics, University of North Carolina at Charlotte, Charlotte, North Carolina 28223, USA; ¹⁹Max Planck Institute for Molecular Genetics, Ihnestrae 63-73, 14195 Berlin, Germany; ²⁰Dahlem Centre for Genome Research and Medical Systems Biology, D-14195 Berlin-Dahlem, Germany; ²¹The Genome Center, Washington University School of Medicine, St Louis, Missouri 63108, USA; ²²Department of Biology, Boston College, Chestnut Hill, Massachusetts 02467, USA; ²³Department of Genome Sciences, University of Washington School of Medicine, Seattle, Washington 98195, USA; ²⁴Affymetrix, Inc., Santa Clara, California 95051, USA; ²⁵US National Institutes of Health, National Center for Biotechnology Information, 45 Center Drive, Bethesda, Maryland 20892, USA; ²⁶BGI-Shenzhen, Shenzhen 518083, China; ²⁷The Novo Nordisk Foundation Center for Basic Metabolic Research, University of Copenhagen, DK-2200 Copenhagen, Denmark; ²⁸Department of Biology, University of Copenhagen, DK-2100 Copenhagen, Denmark; ²⁹Prince Aljawhra Center of Excellence in Research of Hereditary Disorders, King Abdulaziz University, Saudi Arabia; ³⁰James D. Watson Institute of Genome Science, Hangzhou 310008, China; ³¹Alacris Theranostics GmbH, D-14195 Berlin-Dahlem, Germany; ³²Department of Genetics, Albert Einstein College of Medicine, Bronx, New York 10461, USA; ³³Department of Computational Medicine and Bioinformatics, University of Michigan, Ann Arbor, Michigan 48109, USA; ³⁴Cold Spring Harbor Laboratory, Cold Spring Harbor, New York 11724, USA; ³⁵Seaver Autism Center and Department of Psychiatry, Mount Sinai School of Medicine, New York, New York 10029, USA; ³⁶Department of Nanobiomedical Science, Dankook University, Cheonan 330-714, South Korea; ³⁷Department of Biological Sciences, Dankook University, Cheonan 330-714, South Korea; ³⁸Department of Biological Statistics and Computational Biology, Cornell University, Ithaca, New York 14853, USA; ³⁹Center for Systems Biology and Department Organismic and Evolutionary Biology, Harvard University, Cambridge, Massachusetts 02138, USA; ⁴⁰Institute of Medical Genetics, School of Medicine, Cardiff University, Heath Park, Cardiff CF14 4XN, UK; ⁴¹Illumina, Inc., San Diego, California 92122, USA; ⁴²Molecular Epidemiology Section, Department of Medical Statistics and Bioinformatics, Leiden University Medical Center, 2333 ZA, The Netherlands; ⁴³Department of Biological Sciences, Louisiana State University, Baton Rouge, Louisiana 70803, USA; ⁴⁴Analytic and Translational Genetics Unit, Massachusetts General Hospital, Boston, Massachusetts 02114, USA; ⁴⁵Department of Anthropology, Penn State University, University Park, Pennsylvania 16802, USA; ⁴⁶Department of Genetics, Stanford University, Stanford, California 94305, USA; ⁴⁷Ancestry.com, San Francisco, California 94107, USA; ⁴⁸Blavatnik School of Computer Science, Tel Aviv University, 69978 Tel Aviv, Israel; ⁴⁹Department of Microbiology, Tel Aviv University, 69978 Tel Aviv, Israel; ⁵⁰International Computer Science Institute, Berkeley, California 94704, USA; ⁵¹The Translational Genomics Research Institute, Phoenix, Arizona 85004, USA; ⁵²Life Technologies, Beverly, Massachusetts 01915, USA; ⁵³Department of Human Genetics, David Geffen School of Medicine, University of California, Los Angeles, California 90024, USA; ⁵⁴Department of Psychiatry, University of California, San Diego, La Jolla, California 92093, USA; ⁵⁵Department of Cellular and Molecular Medicine, University of California, San Diego, La Jolla,

California 92093, USA; ⁵⁶Department of Computer Science, University of California, San Diego, La Jolla, California 92093, USA; ⁵⁷Department of Epidemiology and Population Health, Albert Einstein College of Medicine, Bronx, New York 10461, USA; ⁵⁸Department of Bioengineering and Therapeutic Sciences and Medicine, University of California, San Francisco, California 94158, USA; ⁵⁹Center for Biomolecular Science and Engineering, University of California, Santa Cruz, California 95064, USA; ⁶⁰Howard Hughes Medical Institute, Santa Cruz, California 95064, USA; ⁶¹Department of Human Genetics, University of Chicago, Chicago, Illinois 60637, USA; ⁶²Department of Genetics, Evolution and Environment, University College London, London WC1E 6BT, UK; ⁶³Department of Genetic Medicine and Development, University of Geneva Medical School, 1211 Geneva, Switzerland; ⁶⁴Institute for Genetics and Genomics in Geneva (IGE3), University of Geneva, 1211 Geneva, Switzerland; ⁶⁵Swiss Institute of Bioinformatics, 1211 Geneva, Switzerland; ⁶⁶Institute for Genome Sciences, University of Maryland School of Medicine, Baltimore, Maryland 21201, USA; ⁶⁷IST/High Performance and Research Computing, University of Medicine and Dentistry of New Jersey, Newark, New Jersey 07107, USA; ⁶⁸Department of Invertebrate Zoology, American Museum of Natural History, New York, New York 10024, USA; ⁶⁹Istituto di Ricerca Genetica e Biomedica, CNR, Monserrato, 09042 Cagliari, Italy; ⁷⁰Department of Anthropology, University of Michigan, Ann Arbor, Michigan 48109, USA; ⁷¹Dipartimento di Scienze Biomediche, Universit delgi Studi di Sassari, 07100 Sassari, Italy; ⁷²Center for Advanced Studies, Research, and Development in Sardinia (CRS4), AGCT Program, Parco Scientifico e tecnologico della Sardegna, 09010 Pula, Italy; ⁷³Department of Genetics, University of North Carolina at Chapel Hill, Chapel Hill, North Carolina 27599, USA; ⁷⁴University of Michigan Sequencing Core, University of Michigan, Ann Arbor, Michigan 48109, USA; ⁷⁵National Institute on Aging, Laboratory of Genetics, Baltimore, Maryland 21224, USA; ⁷⁶Department of Pediatrics, University of Montreal, Sainte-Justine Hospital Research Centre, Montreal, Quebec, Canada H3T 1C5; ⁷⁷Department of Biology, University of Puerto Rico, Mayaguez, Puerto Rico 00680, USA; ⁷⁸The University of Texas Health Science Center at Houston, Houston, Texas 77030, USA; ⁷⁹Eccles Institute of Human Genetics, University of Utah School of Medicine, Salt Lake City, Utah 84112, USA; ⁸⁰Department of Genetics, Rutgers University, The State University of New Jersey, Piscataway, New Jersey 08854, USA; ⁸¹Department of Medicine, Division of Medical Genetics, University of Washington, Seattle, Washington 98195, USA; ⁸²Department of Computer Engineering, Bilkent University, TR-06800 Bilkent, Ankara, Turkey; ⁸³Department of Computer Science, Simon Fraser University, Burnaby, British Columbia, Canada V5A 1S6; ⁸⁴Department of Bioinformatics and Computational Biology, The University of Texas MD Anderson Cancer Center, Houston, Texas 77230, USA; ⁸⁵Department of Haematology, University of Cambridge and National Health Service Blood and Transplant, Cambridge CB2 1TN, UK; ⁸⁶Institute of Genetics and Biophysics, National Research Council (CNR), 80125 Naples, Italy; ⁸⁷Program in Computational Biology and Bioinformatics, Yale University, New Haven, Connecticut 06520, USA; ⁸⁸Department of Computer Science, Yale University, New Haven, Connecticut 06520, USA; ⁸⁹Department of Molecular Biophysics and Biochemistry, Yale University, New Haven, Connecticut 06520, USA; ⁹⁰Department of Chemistry, Yale University, New Haven, Connecticut 06520, USA; ⁹¹Beyster Center for Genomics of Psychiatric Diseases, University of California, San Diego, La Jolla, California 92093, USA; ⁹²US National Institutes of Health, National Human Genome Research Institute, 50 South Drive, Bethesda, Maryland 20892, USA; ⁹³Division of Allergy and Clinical Immunology, School of Medicine, Johns Hopkins University, Baltimore, Maryland 21205, USA; ⁹⁴Coriell Institute for Medical Research, Camden, New Jersey 08103, USA; ⁹⁵Centre for Health, Law and Emerging Technologies, University of Oxford, Oxford OX3 7LF, UK; ⁹⁶Genetic Alliance, London N1 3QP, UK; ⁹⁷Johns Hopkins University School of Medicine, Baltimore, Maryland 21205, USA; ⁹⁸Department of Medical History and Bioethics, Morgridge Institute for Research, University of Wisconsin-Madison, Madison, Wisconsin 53706, USA; ⁹⁹University of Wisconsin Law School, Madison, Wisconsin 53706, USA; ¹⁰⁰The Ethox Centre, Department of Public Health, University of Oxford, Old Road Campus, Oxford OX3 7LF, UK; ¹⁰¹US National Institutes of Health, Center for Research on Genomics and Global Health, National Human Genome Research Institute, 12 South Drive, Bethesda, Maryland 20892, USA; ¹⁰²Institute for Genome Sciences and Policy, Duke University, Durham, North Carolina 27708, USA; ¹⁰³Department of Genetics, University of Pennsylvania School of Medicine, Philadelphia, Pennsylvania 19104, USA; ¹⁰⁴Department of Animal Biology, Unit of Anthropology, University of Barcelona, 08028 Barcelona, Spain; ¹⁰⁵Cancer and Immunogenetics Laboratory, University of Oxford, John Radcliffe Hospital, Oxford OX3 9DS, UK; ¹⁰⁶Laboratory of Molecular Genetics, Institute of Biology, University of Antioquia, Medellin, Colombia; ¹⁰⁷Peking University Shenzhen Hospital, Shenzhen 518036, China; ¹⁰⁸Institute of Medical Biology, Chinese Academy of Medical Sciences and Peking Union Medical College, Kunming 650118, China; ¹⁰⁹Instituto de Biología Molecular y Celular del Cancer, Centro de Investigación del Cancer/IBMCC (CSIC-USAL), Institute of Biomedical Research of Salamanca (IBSAL), Banco Nacional de ADN Carlos III, University of Salamanca, 37007 Salamanca, Spain; ¹¹⁰Instituto de Biología Molecular y Celular del Cancer, Centro de Investigación del Cancer/IBMCC (CSIC-USAL), Institute of Biomedical Research of Salamanca (IBSAL), Cytometry Service and Department of Medicine, University of Salamanca, 37007 Salamanca, Spain; ¹¹¹Ponce School of Medicine and Health Sciences, Ponce, Puerto Rico 00716, USA; ¹¹²US National Institutes of Health, National Human Genome Research Institute, 5635 Fishers Lane, Bethesda, Maryland 20892, USA; ¹¹³Wellcome Trust, Gibbs Building, 215 Euston Road, London NW1 2BE, UK. ‡Deceased.