

ARTICLE

Rare variant genotype imputation with thousands of study-specific whole-genome sequences: implications for cost-effective study designs

Giorgio Pistis^{1,2,3,8}, Eleonora Porcu^{1,2,3,8}, Scott I Vrieze², Carlo Sidore^{1,2,3}, Maristella Steri¹, Fabrice Danjou¹, Fabio Busonero^{1,2}, Antonella Mulas^{1,3}, Magdalena Zoledziewska¹, Andrea Maschio^{1,2}, Christine Brennan⁴, Sandra Lai¹, Michael B Miller⁵, Marco Marcelli⁶, Maria Francesca Urru⁶, Maristella Pitzalis¹, Robert H Lyons⁴, Hyun M Kang², Chris M Jones⁶, Andrea Angius^{1,4}, William G Iacono⁵, David Schlessinger⁷, Matt McGue⁵, Francesco Cucca^{1,3,9}, Gonçalo R Abecasis^{2,9} and Serena Sanna^{*1,9}

The utility of genotype imputation in genome-wide association studies is increasing as progressively larger reference panels are improved and expanded through whole-genome sequencing. Developing general guidelines for optimally cost-effective imputation, however, requires evaluation of performance issues that include the relative utility of study-specific compared with general/multipopulation reference panels; genotyping with various array scaffolds; effects of different ethnic backgrounds; and assessment of ranges of allele frequencies. Here we compared the effectiveness of study-specific reference panels to the commonly used 1000 Genomes Project (1000G) reference panels in the isolated Sardinian population and in cohorts of European ancestry including samples from Minnesota (USA). We also examined different combinations of genome-wide and custom arrays for baseline genotypes. In Sardinians, the study-specific reference panel provided better coverage and genotype imputation accuracy than the 1000G panels and other large European panels. In fact, even gene-centered custom arrays (interrogating ~200 000 variants) provided highly informative content across the entire genome. Gain in accuracy was also observed for Minnesotans using the study-specific reference panel, although the increase was smaller than in Sardinians, especially for rare variants. Notably, a combined panel including both study-specific and 1000G reference panels improved imputation accuracy only in the Minnesota sample, and only at rare sites. Finally, we found that when imputation is performed with a study-specific reference panel, cutoffs different from the standard thresholds of MACH-Rsq and IMPUTE-INFO metrics should be used to efficiently filter badly imputed rare variants. This study thus provides general guidelines for researchers planning large-scale genetic studies.

European Journal of Human Genetics (2015) 23, 975–983; doi:10.1038/ejhg.2014.216; published online 8 October 2014

INTRODUCTION

Genome-wide association studies (GWASs) have successfully identified thousands of common, single-nucleotide polymorphisms (SNPs) associated with complex traits. However, existing genotyping arrays used in GWASs survey only a limited repertoire of sequence variation, and underrepresent rare and population-specific variants. Much more complete extraction of genetic variation is now accessible using next-generation sequencing (NGS) technologies, but efficient detection of rare and low-frequency variants requires sequencing hundreds to thousands of individuals.¹

An alternative cost-effective approach to enlarge the frequency spectrum of variants assessed in GWASs capitalizes on publicly available sequencing reference panels, especially the 1000 Genomes Project (1000G) reference panels. Indeed, ‘probabilistic’ sequenced genomes can be reconstructed by means of imputation methods,

inferring untyped variants by combining partial haplotypes found in a study sample with the full haplotypes available in a more densely characterized reference set. It has, however, been unclear how well general reference panels represent variation in populations that were poorly or not at all represented in projects like 1000 Genomes. Furthermore, even for well-represented populations, a complete evaluation is needed to assess the benefits of sequencing more study samples for successfully imputing rare or low-frequency variants.

How can imputation be further improved? Imputation works very well for common variants, but rapid performance degradation is seen for lower minor allele frequencies. The performance depends on multiple factors, including: choice of baseline array, quality of input genotypes/haplotypes and limited representation of reference haplotypes carrying rare alleles. Also and very importantly, differences in

¹Istituto di Ricerca Genetica e Biomedica (IRGB), CNR, Monserrato, Italy; ²Center for Statistical Genetics, University of Michigan, Ann Arbor, MI, USA; ³Dipartimento di Scienze Biomediche, Università di Sassari, Sassari, Italy; ⁴University of Michigan Sequencing Core, University of Michigan Medical School, Ann Arbor, MI, USA; ⁵Department of Psychology, University of Minnesota, Minneapolis, MN, USA; ⁶CRS4, Parco tecnologico della Sardegna, Pula, Cagliari, Italy; ⁷Laboratory of Genetics, NIA, Baltimore, MD, USA

*Correspondence: Dr S Sanna, Istituto di Ricerca Genetica e Biomedica (IRGB), Consiglio Nazionale delle Ricerche (CNR), c/o Cittadella Universitaria di Monserrato, SS 554 Km 4500, Monserrato, 09042 Cagliari, Italy. Tel: +39 070 6754593; Fax: +39 070 6754652; E-mail: serena.sanna@irgb.cnr.it

⁸These authors should be regarded as joint first authors.

⁹These authors should be regarded as joint senior authors.

Received 3 June 2014; revised 14 August 2014; accepted 9 September 2014; published online 8 October 2014

linkage disequilibrium (LD) patterns and allele frequency spectrum significantly decrease the quality of imputation overall, especially when using public reference panels for ancestral or geographically isolated populations.^{2,3}

To investigate these factors, we compared imputation quality using three complementary sets of reference panels: 1488 Sardinians from Sardinia, Italy; 1325 individuals of Northern European ancestry from Minnesota, USA; and 1092 individuals from the 1000 Genomes project. These reference panels permit comparison of the relative efficiency of study-specific imputation in founder (ie, Sardinia) and continental (ie, Northern European) populations that have also been genotyped, and contrast those results with the current standard approach (ie, 1000 Genomes). Finally, we evaluated the efficiency of the conventional quality thresholds to discard poorly imputed rare and low-frequency variants, focusing on metrics defined by the two most commonly used imputation software, MACH⁴ and IMPUTE.⁵ Figure 1 shows a schematic representation of the study.

MATERIALS AND METHODS

Sample description and genotyping

The study sample consists of the SardiNIA and the MCTFR cohorts. Both studies were approved by the corresponding institutional review boards and a signed informed consent was obtained from every volunteer. The SardiNIA cohort comprises 6921 individuals, representing > 60% of the adult population of four villages in the Lanusei Valley in Sardinia. Details on the study have been previously described.⁶ The Minnesota Center for Twin and Family Research (MCTFR^{7,8}) at the University of Minnesota specializes in the use of genetically informative family cohorts to investigate the etiology of behavioral and psychiatric phenotypes. The MCTFR consists of two complementary cohorts. One is a population-based cohort of twins and their parents, and the other is a family adoption study.

The entire SardiNIA cohort was genotyped using the HumanOmniExpress GWAS array, containing ~ 750K markers, and three different Illumina custom arrays: the Cardio-MetaboChip, the ImmunoChip and the HumanExome, each containing ~ 200 000 markers.^{9,10} Genotype calling was performed using the Illumina GenCall algorithm (Illumina, San Diego, CA, USA), and an additional 2968 rare variants were called for HumanExome using Zcall.¹¹ A subset of 1072 samples was also previously genotyped with Affymetrix 6.0 (Affymetrix, Santa Clara, CA, USA).¹²

After performing quality control checks (see Supplementary Information and Supplementary Table S1 for details), we used the quality-checked (QCed) autosomal markers from the HumanOmniExpress, ImmunoChip and Cardio-MetaboChip arrays as baseline genotypes to impute variants detected through sequencing, as described below. In order to have fully comparable data sets for all analyses described here, we considered only the 6602 samples for which all four Illumina arrays were successfully genotyped. Data from the Affymetrix 6.0 array were instead not combined with the Illumina arrays, given the smaller number of samples available (1072 vs 6602); for this set, quality control filters have been already described.¹³

From the QCed set of markers, we extracted a subset of 227 745 SNPs representing most of the content of the Illumina HumanCore array (78.9% prior QC), a low-density genome-wide array. Given the extensive overlap, and considering that after quality control filtering the effective content of an array is always reduced, we treated this subset of markers as an approximation of the genomic content accessible with the HumanCore array that we refer to here as 'pseudo-HumanCore'.

Genotyping protocols and quality control for the MCTFR study have been described previously.^{8,14} In short, the full MCTFR study sample was genotyped with the Illumina 660W-quad array, with 7278 (97.8%) samples and 527 829 (94.3%) markers passing quality control filters. The full sample was also genotyped with the Illumina HumanExome array, with 7244 (97.4%) samples and 144 075 (58.1%) markers passing quality control filters. We initially used 6610 individuals of European ancestry, and noticed that the inclusion of the 1181 individuals who were also in the reference panel biased accuracy estimates at rare variants because of perfect match of haplotypes

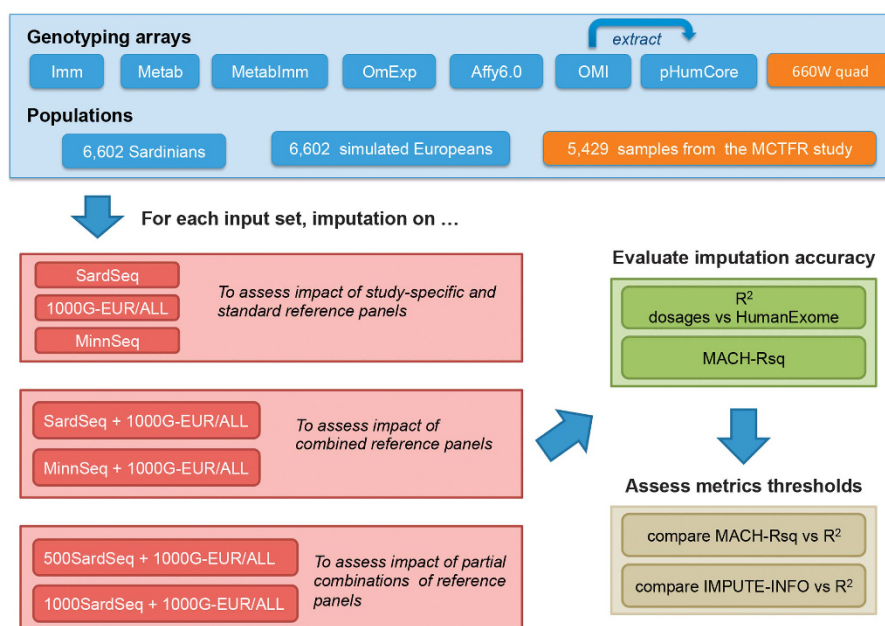


Figure 1 Graphical representation of the analyses and study aims. The figure shows a scheme of analyses carried out. For each genotype input set we carried out several imputation runs (genome-wide for SardiNIA, and on chromosome 20 for other European populations) with different reference panels. We assessed imputation quality of each genotyping array/reference panel combination by looking at the mean imputation quality (MACH-Rsq) and by comparing imputed markers with those directly typed with the HumanExome array (R^2). Finally, we assessed the efficiency of standard thresholds at the commonly used accuracy metrics (MACH-Rsq/IMPUTE-INFO) in filtering badly imputed markers.

(Supplementary Table S2). We therefore restricted the analyses to the 5429 samples not overlapping with the reference panel.

Sequencing and variant calling

Samples to be sequenced were selected in trios, taking advantage of their highly informative content for haplotypes reconstruction. Trios (or parent–offspring pairs for incomplete trios) were selected starting from the founders of all available families to assure the representation of all haplotypes that have been propagated within families (using ExomePicks, <http://genome.sph.umich.edu/wiki/ExomePicks>). For the Sardinians, 2120 samples from 695 nuclear families were sequenced to an average coverage of 4.16-fold. Of these, 1122 samples were part of the SardiNIA project,⁶ whereas the other 998 were individuals enrolled in case–control studies of multiple sclerosis and type I diabetes.^{15,16} The sequencing effort has been described in part previously,¹⁷ and updated details are provided in Supplementary Information.

In the MCTFR study, 1328 individuals from 602 families were sequenced to an average coverage of 10.4-fold. Three samples gave unacceptable sequence quality, leaving 1325 total sequenced samples for analysis.

Variant calling was performed in both studies using GotCloud.¹⁸ Sequencing yielded 17.6 and 27.1 million autosomal bi-allelic SNPs in Sardinian and Minnesota samples, respectively, of which 30.6 and 48.4% were not described in dbSNP135.

Genotype imputation

Genotype imputation for all scenarios were performed on haploid data using Minimac (see URLs in Supplementary Information), a modified version of the MACH⁴ software. For SardiNIA, phased haplotypes were generated using MACH (*–phase* option) with 400 states and 30 rounds by subdividing the variants in 344 groups of 2500 with an overlap of 500, and imputation was subsequently performed independently on each phased chunk (for a description of the code, see the ‘1000G imputation cookbook’ URL in Supplementary Information). Imputation performance was evaluated on seven different input genotype data sets: (1) HumanOmniExpress (OmExp), (2) Cardio-MetaboChip (Metab), (3) ImmunoChip (Imm), (4) Cardio-MetaboChip and ImmunoChip (MetabImm), (5) HumanOmniExpress, Cardio-MetaboChip and ImmunoChip (OMI), (6) pseudo-HumanCore (pHumCore) and (7) Affymetrix 6.0 (Affy 6.0).

For simplicity, we phased the Cardio-MetaboChip, ImmunoChip and HumanOmniExpress arrays jointly, and then extracted haplotypes at relevant SNPs to perform imputation for each particular genotyping set. In actual practice, Cardio-MetaboChip and ImmunoChip will be phased without the additional support of a genome-wide array, and hence we assessed the impact of our procedure by phasing separately each SNP set, for chromosome 20. We noticed that only imputations performed with the SardSeq panel or its combination with 1000G were slightly overestimated (see Supplementary Information and Supplementary Table S3).

In the MCTFR study, haplotypes were phased using SHAPEIT2 (v2.644)¹⁹ with the following model options: *–thread 8 –burn 10 –prune 8 –main 20 –states 200*. Imputation was performed using Minimac and the Illumina 660W-quad array as baseline genotypes.

We used as reference panels the 1000G-ALL (1092 samples) and 1000G-EUR (379 samples) data sets from the 1000 Genomes March 2012 release; the full MCTFR sequencing data (1325 samples, named MinnSeq in the text); a subset of the Sardinian sequencing data (1488 samples, named SardSeq in the text); and combinations of those (see ‘Combination of reference panels’ below). Considering the overall high inbreeding in Sardinia, the SardSeq reference panel was created by selecting only haplotypes of parents at each sequenced trios to avoid overrepresentation of rare variants.

We also performed imputation with IMPUTE2 (newest release of IMPUTE), to test a different approach for reference panels combination (see ‘Combination of reference panels’ paragraph) and to assess the efficiency of its imputation accuracy metric INFO (see ‘Evaluation of imputation accuracy’ paragraph).

Simulation of European haplotypes

Because the Minnesota samples were genotyped with different arrays from those used for Sardinians, they could not be used to assess relative efficiency of

arrays in genotype imputation. We therefore generated, by simulation with the HAPGEN²⁰ software and 1000G-EUR as reference, 6602 unrelated individuals of European ancestry for SNPs present in each different genotyping array considered in the SardiNIA study. For simplicity, we focused only on chromosome 20. Haplotypes were phased using MACH (*–phase* option) with 400 states and 30 rounds, and imputation performed using Minimac, as in the SardiNIA and MCTFR data sets. This simulated data set was only used for assessing the efficiency of different genotyping arrays and reference panels in genotype imputation.

Combination of reference panels

We used VCFtools²¹ to combine the SardSeq and the MinnSeq panels with 1000G-EUR and 1000G-ALL reference panels for chromosome 20. The variants in each set were 331 799, 602 317, 851 702 and 377 494 for SardSeq, MinnSeq, 1000G-ALL and 1000G-EUR, respectively. During the merging procedure, we removed the variants present only in one panel, leading to SardSeq+1000G-ALL, SardSeq+1000G-EUR, MinnSeq+1000G-ALL and MinnSeq+1000G-EUR reference panels containing 249 624, 227 405, 304 899 and 267 550 variants, respectively. Imputation was then performed using Minimac, as for single reference panels. For combinations with 1000G and SardSeq panels, we also performed imputation with IMPUTE2 using the *–merge_ref_panels* option that imputes variants unique to one panel into the other, prior imputation. We observed no difference in imputation accuracy at all frequency ranges when using this approach, which should be preferable for research studies, allowing imputation of all available variants, including those that are study specific, in the same run (Supplementary Table S4).

In addition, to assess the impact of adding a smaller number of population-specific haplotypes, we created two additional reference panels using 500 and 1000 randomly chosen samples from the SardSeq reference panel and merging them with 1000G reference panels (500SardSeq+1000G and 1000SardSeq+1000G, respectively). This analysis was restricted to the SardSeq panel and the SardiNIA cohort, because the advantage in accuracy was substantial for this population.

Evaluation of imputation accuracy

Imputation accuracy was assessed using both the MACH Rsq metric and the squared Pearson’s correlation (R^2)⁴ between dosages and the real genotypes (considered as allele count) available for the same individuals, extracted from the HumanExome array. The Rsq metric is also known as variance ratio, being calculated as the proportion of the empirically observed variance (based on the imputation) to the expected binomial variance $p(1-p)$, where p is the minor allele frequency. In SardiNIA we tested 21 398 SNPs across autosomes for genome-wide evaluation of imputation accuracy and tested a subset of 558 SNPs for comparisons restricted to chromosome 20. For the MCTFR study, as the baseline array was different, we used a subset of 541 SNPs. The number of SNPs tested for comparing imputation with SardSeq *versus* 500SardSeq+1000G and 1000SardSeq+1000G was reduced to 517 because 41 SNPs (MAF range 0.0008–0.0072%) were not detected in the selected subset of sequenced samples.

We also assessed efficiency in discriminating between well and poorly imputed markers of the imputation accuracy metrics estimated by MACH (Rsq) and IMPUTE (INFO).⁴ The INFO metric, also known as imputed information score (INFO), is a measure of the relative statistical information about the SNP allele frequency from the imputed data. We defined good- and bad-quality imputed SNPs as in the original MACH paper, that is, those with $R^2 > 0.5$ and with $R^2 < 0.2$, respectively, and stratified imputed SNPs based on their Rsq and INFO scores. This analysis was restricted to chromosome 20, and performed using as baseline genotypes the OmExp for the SardiNIA study and the Illumina 660W-quad for the MCTFR cohort.

RESULTS

Effect of baseline genotyping array

This subsection is restricted to the SardiNIA study and the simulated European haplotypes, because the MCTFR study used only one array. We found clear differences in imputation performance depending on the baseline genotyping set. Comparable differences were seen when assessments were done with either the Rsq metric—the imputation

quality metric from MACH⁴—or the R^2 metric, the squared Pearson correlation, between dosages and real genotypes⁴ (Table 1).

When using the 1000G reference panels for Sardinians, the two custom arrays (Cardio-MetaboChip and ImmunoChip) provided very limited information for imputation and far less accuracy than the genome-wide arrays, reflecting their low marker density. However, the Cardio-MetaboChip array performed very well when imputing with the SardSeq panel, allowing accurate inference of the rest of the genome (mean $R_{sq}=0.62$, and mean $R^2=0.70$ at HumanExome SNPs). The relative efficiency was similar when considering all autosomes (Table 1 and Supplementary Table S5) or focusing only on chromosome 20 (Figure 2 and Supplementary Table S6). The extended LD in the population and the increased genetic similarity of the reference panel aid in haplotype reconstruction when using a relatively small set of markers.

The addition of the two custom arrays to the OmExp genome-wide array (OmExp+Metab+Imm, called OMI here) did not improve quality for common or low-frequency variants compared with that reached using OmExp alone. Thus, such arrays provide direct genotyping of low-frequency and rare variants in genes of interest but do not contribute to an overall improvement in imputation accuracy. We also observed negligible differences in imputation accuracy between the two tested Illumina genome-wide arrays, OmExp and pHumCore (Table 1 and Supplementary Tables S5 and S6), when imputing the SardSeq panel. In particular, we noticed that the low-density genome-wide array pHumCore provided only slightly less accuracy than the denser OmExp array when the SardSeq sequencing panel was used for imputation (mean $R^2=0.85$ and 0.87 , for pHumCore and OmExp, respectively, at HumanExome SNPs, Supplementary Table S5) and a very similar genomic coverage (92.6 and 91.8% of markers imputed with $R_{sq} > 0.3$, Table 1).

Of note, performance was patently lower for both arrays and more significantly for pHumCore when imputation was performed with the 1000G panels (mean $R^2=0.54$ and 0.64 , for pHumCore and OmExp, respectively, imputing with the 1000G-ALL; Table 1, Supplementary Tables S5 and S6, and Figure 2). In contrast, in the simulated European data, the Cardio-MetaboChip performed poorly, with insufficient genomic coverage. Contrarily to previous observations,²² the pHumCore was fairly comparable in efficiency to the OmExp array (Figure 3 and Supplementary Table S7), but we expect performance to be overestimated (because the genotypes were simulated based on 1000 Genomes). In fact, when we extracted subset of SNPs that are present in HumanOmniExpress and HumanCore from the MCTFR genotypes, the difference between the two arrays was clearly evident (Supplementary Table S8). This difference has also been observed for another European population.²³

Thus, in founder populations it appears that highly accurate imputation can be achieved with cost-effective sparse genotyping arrays when a population-specific reference panel is available.

Effect of study-specific reference panels

Study-specific reference panels increased the accuracy and completeness of coverage in both Sardinian and Minnesota samples, but the gain in accuracy was greater for the Sardinia founder population.

In Sardinians, the 1000G-ALL reference panel provided the highest number of imputed variants—~37 million including both indels and SNPs *vs* ~15 million SNPs for the SardSeq panel—but the majority were of poor quality and were subsequently discarded. For example, for the Metab/SardSeq combination, 11.5 million imputed SNPs passed the standard $R_{sq} > 0.3$ filter, but only 2.7 million and 3.0 million reached that threshold for Metab/1000G-ALL and Metab/1000G-EUR, respectively. The gap was less striking but still marked

Table 1 Basic imputation statistics on the SardiNIA samples for different panels/genotyping arrays

Array	Reference panel	Whole imputed SNP set		$R_{sq} > 0.3$		Shared imputed SNPs
		No. of SNPs	Mean (SD) R_{sq}	% SNPs	Mean (SD) R_{sq}	Mean (SD) R_{sq}
Imm	SardSeq	15 071 719	0.258 (0.312)	33.33	0.652 (0.213)	0.299 (0.321)
	1000G-ALL	37 798 002	0.037 (0.134)	3.90	0.638 (0.232)	0.099 (0.213)
	1000G-EUR	16 873 087	0.085 (0.203)	9.68	0.647 (0.231)	0.115 (0.232)
Metab	SardSeq	15 069 660	0.617 (0.335)	76.91	0.777 (0.181)	0.685 (0.301)
	1000G-ALL	37 782 741	0.064 (0.170)	7.20	0.614 (0.217)	0.175 (0.260)
	1000G-EUR	16 878 099	0.149 (0.253)	18.05	0.634 (0.219)	0.201 (0.282)
MetabImm	SardSeq	14 977 409	0.734 (0.300)	86.51	0.835 (0.163)	0.808 (0.239)
	1000G-ALL	37 721 853	0.100 (0.218)	11.71	0.644 (0.221)	0.272 (0.311)
	1000G-EUR	16 781 983	0.219 (0.303)	27.12	0.667 (0.222)	0.297 (0.328)
OmExp	SardSeq	14 580 754	0.861 (0.256)	92.61	0.924 (0.131)	0.935 (0.161)
	1000G-ALL	37 424 729	0.297 (0.382)	33.61	0.796 (0.224)	0.742 (0.322)
	1000G-EUR	16 453 325	0.543 (0.406)	60.89	0.840 (0.206)	0.729 (0.341)
OMI	SardSeq	14 319 695	0.862 (0.256)	92.57	0.925 (0.131)	0.937 (0.159)
	1000G-ALL	37 211 511	0.300 (0.385)	34.00	0.799 (0.131)	0.753 (0.318)
	1000G-EUR	16 255 689	0.549 (0.406)	61.50	0.842 (0.206)	0.739 (0.337)
pHumCore	SardSeq	15 020 615	0.840 (0.264)	91.81	0.908 (0.139)	0.913 (0.179)
	1000G-ALL	37 793 052	0.234 (0.341)	26.66	0.759 (0.221)	0.614 (0.354)
	1000G-EUR	16 825 817	0.455 (0.398)	52.64	0.802 (0.207)	0.615 (0.367)
Affy6.0	SardSeq	14 550 658	0.798 (0.342)	84.51	0.937 (0.116)	0.905 (0.232)
	1000G-ALL	37 328 716	0.263 (0.379)	29.55	0.814 (0.217)	0.721 (0.341)
	1000G-EUR	16 350 040	0.515 (0.416)	57.63	0.843 (0.205)	0.708 (0.357)

The table shows, for each genotyping array/reference panel combination, the number of imputed SNPs and the corresponding mean R_{sq} and SD, the percentage of SNPs with $R_{sq} > 0.3$, with the corresponding mean R_{sq} and SD, mean R_{sq} and SD evaluated for 8 842 944 SNPs that were imputed in all genotyping array/reference panel combinations (called 'Shared imputed SNPs').

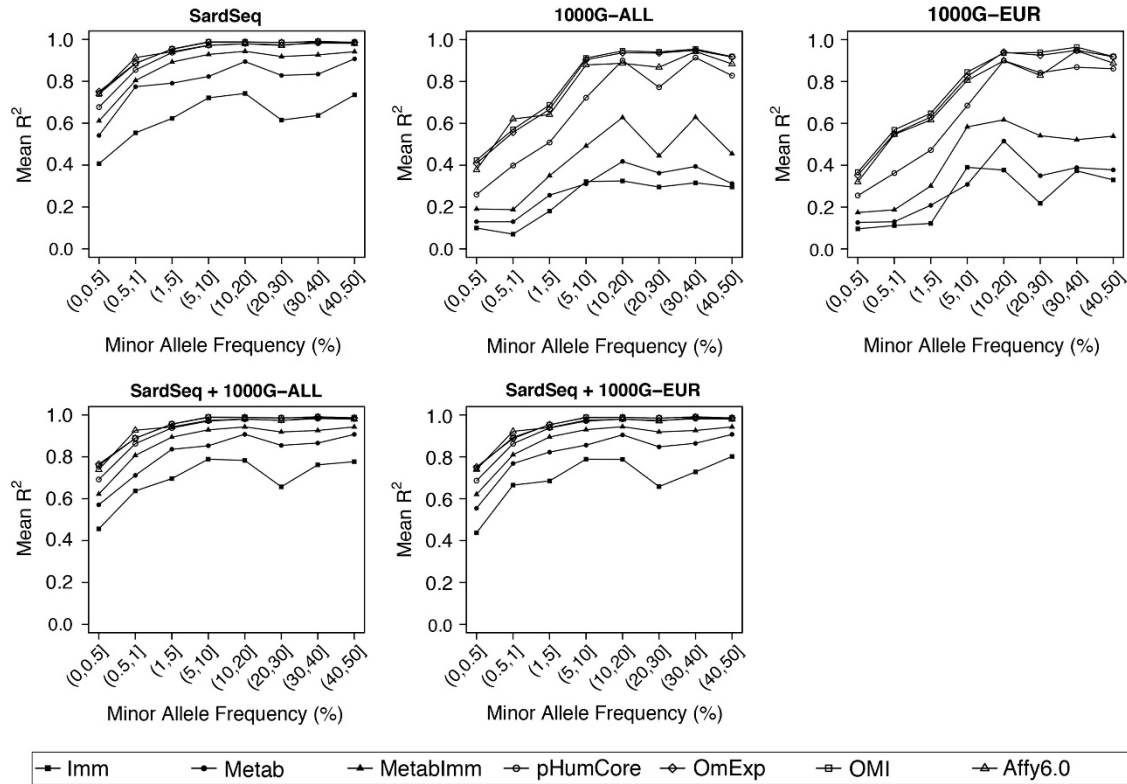


Figure 2 Mean R^2 for each particular genotyping array/reference panel in the Sardinia cohort. The figure shows the mean R^2 at different allele frequency ranges for each particular genotyping array/reference panel combination, including the combination of SardSeq and 1000G panels. Results are restricted to chromosome 20.

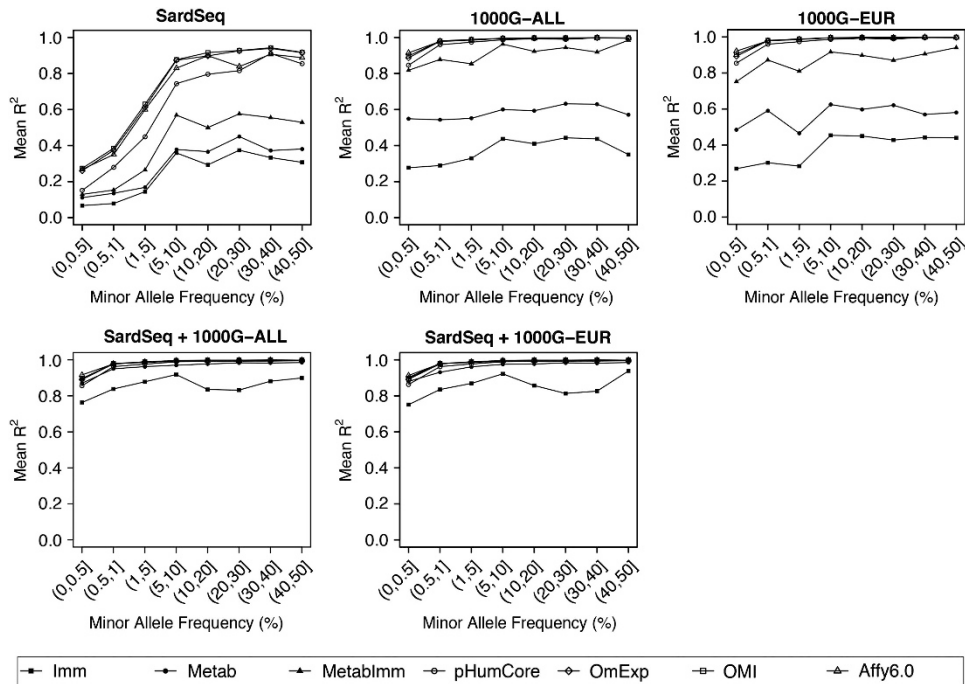


Figure 3 Mean R^2 for each combination of genotyping array/reference panel in the European simulated data set. The figure shows the mean R^2 at different allele frequency ranges for each particular genotyping array/reference panel, including the combination of SardSeq and 1000G panels. Results are restricted to chromosome 20.

when denser genotype data sets were considered, and was still noticeable even considering only SNPs present in all reference panels (which are enriched for high-frequency variants; Table 1). Consistent results were seen for the OmExp, OMI, pHumCore and Affy6.0 data sets, with accuracy consistently better when using SardSeq (Figure 2).

The benefit in overall accuracy was clear at all frequency ranges and even greater for low-frequency and rare variants. For example, using the OMI data set, the average R^2 for SNPs with MAF ranging from 0.5 to 1% is 0.91, 0.57 and 0.52 when using SardSeq, 1000G-ALL and 1000G-EUR reference panels, respectively (Supplementary Table S5). This reinforces the finding that on average, low-frequency variants are hard to impute in founder populations when using external reference panels because these variants appear in fewer haplotypes.² Of note, the results remained the same after removing 646 Sardinian samples that appear in both the genotyping set and the SardSeq reference panel (Supplementary Table S2).

To assess whether the advantage with the SardSeq panel was attributable to the lower number of European haplotypes present in the 1000 Genomes reference, we performed imputation using the MinnSeq panel. There was no appreciable gain in accuracy within Sardinians compared with 1000G-based imputations (Figure 4a, Supplementary Tables S9A and S10).

Similar to results with Sardinians, the MinnSeq panel outperformed the 1000G panels in the MCTFR study at all frequency ranges (Figure 4b and Supplementary Table S9B). However, the gain in accuracy was far less than that observed in Sardinians with the SardSeq panels. For example, for variants with MAF ranging from 1 to 5%, we observed 11% and 42% additional gain in mean R^2 for Minnesotans and Sardinians, respectively. Of note, in both cohorts the study-specific panel also yielded a higher number of SNPs useful for analyses (considering an $R_{sq} > 0.3$) even when the other reference sets contain more SNPs (Supplementary Table S10).

Effect of combined reference panels

We also evaluated the impact on imputation accuracy of extended panels created by combining the two study-specific panels and 1000G haplotypes.

The combined SardSeq+1000G panels provided only marginally higher accuracy at rarer shared SNPs in Sardinians (Figure 2 and Supplementary Tables S5 and S6). Slight increase in accuracy was also observable for more frequent variants (except for the two custom

arrays (Metab and Imm), for which the improvement was substantial across all frequency ranges (Figure 2 and Supplementary Tables S5 and Table S6)). Thus, for Sardinians, the inclusion of 1000G haplotypes would only be beneficial for very rare variants if a genome-wide array was used for baseline imputation.

In the simulated European set, the addition of SardSeq haplotypes to the 1000G panels remarkably increased imputation accuracy for custom genotyping arrays (Metab and Imm) for both common and rare variants (Figure 3 and Supplementary Table S7). For example, for variants with $MAF > 40\%$ and $MAF \leq 50\%$ the mean R^2 is 0.57 and 0.98, when imputing with 1000G-ALL and SardSeq+1000G-ALL and using the Metab data set (Figure 3 and Supplementary Table S7). The impact of a combined panel was instead negligible for the more comprehensive genotype data (OmExp, OMI, pHumCore and Affy6.0). However, imputation on simulated data could give slight overestimations, and this could mask the advantage of adding SardSeq to 1000G panels. Indeed, when considering the MCTFR study, the combined SardSeq+1000G-ALL panel provided benefit at all frequency ranges compared with 1000G-ALL imputation, and for $MAF \leq 0.5\%$ variants accuracy becomes fairly similar to that observed when using the MCTFR-specific panel (Figure 4 and Supplementary Table S9). Thus, the Sardinian panel could be generally useful to increase the overall accuracy in population cohorts other than Sardinians, especially where only custom array genotyping is available or when a study-specific reference is not available.

Compared with imputation with MinnSeq alone, the addition of the 1000G haplotypes to the MinnSeq reference panel was useful only for rare variants in Minnesotans. The difference in accuracy was >4-fold higher than that seen in Sardinians comparing imputations with SardSeq and SardSeq+1000G panels. Thus, for Europeans, the inclusion of 1000G haplotypes in a study-specific panel is sensitively beneficial for very rare variants. Of note, for the Minnesotans, genotype imputation at the full spectrum of frequency ranges never reaches the same accuracy as in Sardinians with the SardSeq panel, even when using the combined MinnSeq+1000G with almost twice as many individuals as there are in the SardSeq panel.

Given the great utility of the Sardinian haplotypes, we further examined whether the advantage achieved by imputing with the SardSeq panel could have been reached sequencing a smaller number of samples and merging their haplotypes with the 1000G panels. For simplicity, we again focused on chromosome 20 and the OmExp

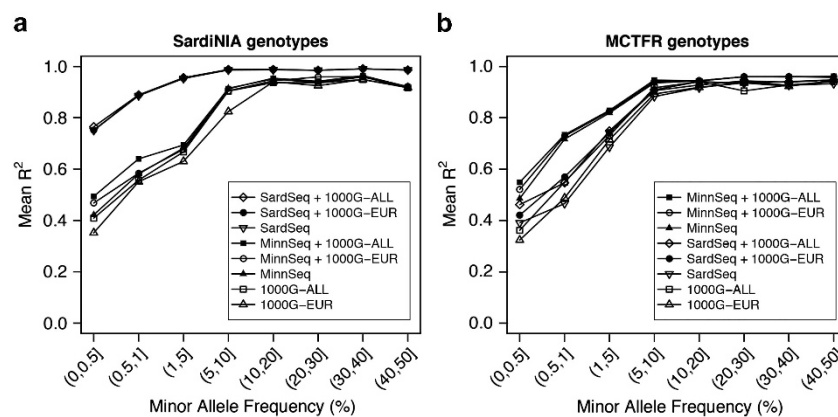


Figure 4 Impact of cross-studies reference panels. The figure shows the mean R^2 at different allele frequency ranges for the chromosome 20 of OmExp genotyping array for SardiNIA (a) and the Illumina 660W-quad array for the MCTFR (b) study, when using different reference panels, including combination of SardSeq/MinnSeq and 1000G panels and cross-studies references.

array. Only for variants with MAF >5% does adding 500 Sardinian samples to the 1000G panels provide the same accuracy as the SardSeq panel alone. Instead, adding 1000 Sardinians to the 1000G panels provides the same accuracy given by the SardSeq panel for all frequency bins, with only a modest difference in accuracy for the very rare variants (MAF<0.5%) (Supplementary Figure S1 and Supplementary Table S11).

Thus, sequencing a smaller number of individuals and combining their haplotypes with the 1000G panels could give imputation accuracy that is highly comparable to a panel comprising a large number of samples. However, the caveat remains that the genotype accuracy and variant discovery in low-pass sequencing is highly dependent on the number of sequenced samples. Consequently, sequencing only 500 samples would not provide genotypes as precise as those obtained by randomly selecting 500 samples from a set of 2000 sequenced genomes. For example, when we performed variant calling on a subset of 508 samples, the heterozygous error rate increased from 2.6 to 11.3% at rare sites (Supplementary Table S12).

Performance of imputation quality metrics

To determine whether the commonly used MACH-Rsq threshold >0.3 and IMPUTE-INFO >0.4 can be applied to all frequency ranges (and if not, to infer appropriate cutoffs), we investigated how well imputation quality metrics can predict true imputation accuracy, especially for rare and less common variants. We found that for MAF ≥1%, imputation accuracy and therefore concordance between real genotypes and dosages using study-specific panels was almost

perfect in both Sardinians and Minnesotans (Tables 2 and 3 and Supplementary Figure S2). At these frequency ranges, high but clearly less concordance was also seen when imputing with the 1000G panels. Whatever the reference panel used and the population under study, the standard Rsq cutoff of >0.3 efficiently discarded most badly imputed markers while keeping most of those imputed well (see Materials and Methods). In particular, imputation was so accurate overall that even an Rsq cutoff of >0 would leave no badly imputed markers on chromosome 20 (Tables 2a and 3a) (and only 8 over the entire genome in Sardinians, Supplementary Table S13). Similarly for the INFO metrics, the standard >0.4 threshold was efficient to discriminate between well and poorly inferred genotypes at this range of frequency (Tables 2b and 3b).

In contrast, for MAF <1%, we noticed that both metrics were slightly overestimated when using the study-specific panels, possibly because of the inclusion of relatives with similar haplotypes in the target data set; but overall concordance was better than 1000G imputation for this range of frequency as well. Specifically, in this range and when imputation was performed with the 1000G panels, the threshold of Rsq >0.3 was less efficient, aggressively discarding some well-imputed variants (eliminating 7–18% and 7–25% of the well-imputed markers for ALL and EUR panels) and retaining an excess of the badly imputed ones (Tables 2a and 3a and Supplementary Tables S14 and S15). The INFO >0.4 threshold instead worked efficiently on selecting well-imputed variants, but was too lenient on discarding those of poor quality (Tables 2b and 3b and Supplementary Tables S14

Table 2 Efficiency of imputation quality metrics in the SardiNIA cohort

	MAF < 1%				MAF ≥ 1%			
	SardSeq		1000G-ALL		SardSeq		1000G-ALL	
	% Bad (n)	% Good (n)	% Bad (n)	% Good (n)	% Bad (n)	% Good (n)	% Bad (n)	% Good (n)
(a) Rsq>0	100 (14)	100 (222)	100 (98)	100 (124)	0 (0)	100 (301)	100 (20)	100 (255)
>0.1	92.86 (13)	100 (222)	44.90 (44)	92.74 (115)	0 (0)	100 (301)	90 (18)	99.61 (254)
>0.2	85.71 (12)	100 (222)	19.39 (19)	86.29 (107)	0 (0)	100 (301)	75 (15)	99.61 (254)
>0.3	78.57 (11)	100 (222)	11.22 (11)	81.45 (101)	0 (0)	100 (301)	65 (13)	98.43 (251)
>0.4	71.43 (10)	100 (222)	5.10 (5)	70.16 (87)	0 (0)	100 (301)	45 (9)	97.25 (248)
>0.5	64.29 (9)	99.55 (221)	3.06 (3)	62.90 (78)	0 (0)	100 (301)	30 (6)	94.90 (242)
>0.6	42.86 (6)	95.95 (213)	2.04 (2)	50.81 (63)	0 (0)	100 (301)	20 (4)	89.41 (228)
>0.7	28.57 (4)	91.89 (204)	0 (0)	43.55 (54)	0 (0)	100 (301)	15 (3)	82.35 (210)
>0.8	14.29 (2)	83.78 (186)	0 (0)	33.87 (42)	0 (0)	100 (301)	0 (0)	72.16 (184)
>0.9	7.14 (1)	58.11 (129)	0 (0)	23.39 (29)	0 (0)	98.01 (295)	0 (0)	59.22 (151)
>1	0 (0)	0.90 (2)	0 (0)	0 (0)	0 (0)	3.65 (11)	0 (0)	2.75 (7)
(b) INFO								
>0	100 (7)	100 (189)	100 (81)	100 (83)	0 (0)	100 (307)	100 (32)	100 (251)
>0.1	100 (7)	100 (189)	100 (81)	98.80 (82)	0 (0)	100 (307)	100 (32)	100 (251)
>0.2	100 (7)	99.47 (188)	90.12 (73)	97.59 (81)	0 (0)	100 (307)	100 (32)	100 (251)
>0.3	100 (7)	99.47 (188)	62.96 (51)	96.39 (80)	0 (0)	100 (307)	96.88 (31)	100 (251)
>0.4	100 (7)	99.47 (188)	48.15 (39)	93.98 (78)	0 (0)	100 (307)	96.88 (31)	100 (251)
>0.5	100 (7)	99.47 (188)	27.16 (22)	89.16 (74)	0 (0)	100 (307)	84.38 (27)	99.20 (249)
>0.6	100 (7)	98.94 (187)	17.28 (14)	85.54 (71)	0 (0)	100 (307)	59.38 (19)	98.01 (246)
>0.7	71.43 (5)	97.35 (184)	11.11 (9)	73.49 (61)	0 (0)	100 (307)	37.50 (12)	95.62 (240)
>0.8	42.86 (3)	92.59 (175)	3.70 (3)	60.24 (50)	0 (0)	100 (307)	15.62 (5)	88.84 (223)
>0.9	14.29 (1)	76.19 (144)	0 (0)	44.58 (37)	0 (0)	99.35 (305)	6.25 (2)	72.91 (183)
>1	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)

The table shows the number and the percentage of poorly imputed and well-imputed SNPs (see Materials and Methods) that are captured for each Rsq (a) and INFO (b) threshold. Imputation was performed on chromosome 20 HumanOmniExpress SNPs, using the SardSeq and 1000G-ALL panels. Statistics are reported separately for common and rare variants.

Table 3 Efficiency of imputation quality metrics in the MCTFR cohort

	MAF < 1%				MAF ≥ 1%			
	MinnSeq		1000G-ALL		SardSeq		1000G-ALL	
	% Bad (n)	% Good (n)	% Bad (n)	% Good (n)	% Bad (n)	% Good (n)	% Bad (n)	% Good (n)
(a) <i>Rsq</i> > 0	100 (38)	100 (129)	100 (80)	100 (92)	0 (0)	100 (284)	100 (4)	100 (258)
> 0.1	81.58 (31)	100 (129)	72.50 (58)	96.74 (89)	0 (0)	100 (284)	100 (4)	100 (258)
> 0.2	73.68 (28)	100 (129)	41.25 (33)	95.65 (88)	0 (0)	100 (284)	25 (1)	100 (258)
> 0.3	57.89 (22)	100 (129)	26.25 (21)	92.39 (85)	0 (0)	100 (284)	25 (1)	99.61 (257)
> 0.4	47.37 (18)	100 (129)	17.50 (14)	83.70 (77)	0 (0)	100 (284)	25 (1)	98.84 (255)
> 0.5	28.95 (11)	96.90 (125)	10 (8)	72.83 (67)	0 (0)	100 (284)	0 (0)	96.12 (248)
> 0.6	21.05 (8)	92.25 (119)	3.75 (3)	59.78 (55)	0 (0)	95.07 (270)	0 (0)	87.21 (225)
> 0.7	2.63 (1)	72.09 (93)	1.25 (1)	48.91 (45)	0 (0)	89.79 (255)	0 (0)	77.13 (199)
> 0.8	0 (0)	51.94 (67)	0 (0)	31.52 (29)	0 (0)	79.58 (226)	0 (0)	62.02 (160)
> 0.9	0 (0)	28.68 (37)	0 (0)	17.39 (16)	0 (0)	59.51 (169)	0 (0)	48.45 (125)
> 1	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)
(b) <i>INFO</i>								
> 0	100 (38)	100 (96)	100 (82)	100 (67)	100 (1)	100 (277)	100 (9)	100 (241)
> 0.1	100 (38)	100 (96)	100 (82)	100 (67)	100 (1)	100 (277)	100 (9)	100 (241)
> 0.2	100 (38)	100 (96)	97.56 (80)	100 (67)	100 (1)	100 (277)	100 (9)	100 (241)
> 0.3	97.37 (37)	100 (96)	95.12 (78)	100 (67)	100 (1)	100 (277)	100 (9)	100 (241)
> 0.4	94.74 (36)	100 (96)	69.51 (57)	100 (67)	100 (1)	100 (277)	100 (9)	100 (241)
> 0.5	73.68 (28)	100 (96)	41.46 (34)	100 (67)	100 (1)	100 (277)	100 (9)	100 (241)
> 0.6	50 (19)	98.96 (95)	14.63 (12)	100 (67)	0 (0)	100 (277)	44.44 (4)	100 (241)
> 0.7	23.68 (9)	95.83 (92)	7.32 (6)	92.54 (62)	0 (0)	99.28 (275)	0 (0)	99.59 (240)
> 0.8	5.26 (2)	77.08 (74)	0 (0)	71.64 (48)	0 (0)	91.70 (254)	0 (0)	87.97 (212)
> 0.9	2.63 (1)	40.62 (39)	0 (0)	46.27 (31)	0 (0)	72.20 (200)	0 (0)	63.49 (153)
> 1	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)

The table shows the number and the percentage of poorly imputed and well-imputed SNPs (see Materials and Methods) that are captured for each *Rsq* (a) and *INFO* (b) threshold. Imputation was performed on chromosome 20 Illumina 660W-quadrant array SNPs, using MinnSeq and 1000G-ALL as reference panels. Statistics are reported separately for common and rare variants.

and S15). Nevertheless, *Rsq* > 0.3 and *INFO* > 0.4 still remain the optimal thresholds.

When imputation was performed with the study-specific panels, both the *Rsq* and *INFO* thresholds were more efficient in capturing all well-imputed markers, but less efficient in discarding the poorly imputed.

In such cases, ie for MAF < 1% and when imputation is performed with a reference panel that is genetically close to the study population, an *Rsq* threshold of > 0.6 and *INFO* > 0.7 should be preferred in lieu of the standard thresholds of 0.3 and 0.4, respectively.

DISCUSSION

We used different reference panels and genotype input sets to investigate the effects on imputation in founder and nonfounder populations of European ancestry. We found that a study-specific reference panel considerably improved imputation accuracy and genomic coverage compared with external equally large reference panels, regardless of the genotyping array, especially for rare variants. However, the benefit was strikingly higher in the founder population of Sardinians, with a precision that was not obtainable in Europeans even with a reference panel twice the size. In fact, in such homogenous populations each sequenced genome provides information that can be extended to distant relatives as well, whereas in continental Europeans, haplotypes carrying rare variants can only inform closely related samples.

We also observed that in Sardinians a study-specific panel boosts imputation even for low-coverage genotyping array(s) like the Cardio-

MetaboChip that are barely informative when imputing with the 1000G panels alone, or for the HumanCore that becomes highly comparable for all frequency ranges to the wider HumanOmniExpress. Given the low cost of the sparser arrays, accurate population-scale imputation is more feasible in the Sardinian founder population than in nonfounder populations when combined with large-scale sequencing. For example, at current cost schedules, with an investment of 500 000 dollars one could genotype ~ 8300 Sardinian samples with the HumanCore array instead of ~ 4500 with the HumanOmniExpress. The power to detect association for variants accounting for 0.5% of the trait variance thereby rises from 24 to 84%.

Finally, we observed that standard thresholds on metrics for evaluating accuracy, estimated by two commonly used imputation software, are somewhat imprecise for rare variants. We propose that all cohorts using study-specific reference panels for imputation consider adopting different thresholds for common and rare variants to filter inaccurate genotypes.

Taken together, these imputation-based analyses can guide genetic studies, and complement recent reports^{22,24} with several novel aspects that can improve performance:

- They exploit imputation accuracy with the two larger study-specific reference panels so far published, including one that is population specific.
- They also provide the first evaluation of imputation performance of the 1000 Genomes Project haplotypes in an isolated population.

- They include analyses of large cohorts coupled with the use of HumanExome array, allowing appropriate assessment of results for less frequent and rare variants.
- Using real data sets, they based analyses on a subset of quality-controlled SNPs instead of the full list of markers present on an array (excluding many that are likely to be imperfectly genotyped in a case study).
- They evaluate two widely used custom genotyping arrays, Cardio-MetaboChip and ImmunoChip, providing information for cohorts that are limited to that source of genotypes.
- They also evaluate for rare variants the efficiency of accuracy metric thresholds that were previously suggested for common variants.

Ultimately, full genome sequencing could make imputation methods superfluous, but the timescale remains indeterminate. It should be considered that increasing sample size can augment genome-wide power to assess rare variants more than increasing array density – even up to full genotyping of the complete 1000 Genomes Project variant set.^{22,24} Thus, aids to imputation are increasingly valuable, because most studies are likely to be collecting increasing numbers of samples and using this inferential process rather than sequencing full genomes.

Overall, population-specific panels might have been thought to be ‘private’, with potential discoveries limited to that population. Instead, the effectiveness of population-specific reference panels can be appreciable for other populations, but will vary depending on the size of the panels and the demographic history of the isolate. Intuitively in Europe, their value may be greater for populations like Basques and Greeks, who are relatively genetically distant from the European samples selected for the 1000 Genomes Project. Here, we show that sequencing efforts from the Sardinian founder population can, when coupled with available panels, improve rare variant imputation accuracy in other population backgrounds as well. This reinforces the value of isolated populations for discovery of variants that are locally enriched but rarer and thus harder to detect in international surveys.²⁵

CONFLICT OF INTEREST

The authors declare no conflict of interest.

ACKNOWLEDGEMENTS

This work was supported by the Intramural Research Program of the National Institute of Health, National Institute on Aging (N01-AG-1-2109 and HHSN271201100005C), National Human Genome Research Institute grants (HG005581, HG005552, HG006513 and HG007022 to GRA), the National Institute on Drug Abuse (DA 024417 and DA 034606) and the Italian FISM (2011/R/13 to FC). We thank Frederic Reinier, Riccardo Berutti, Rossano Atzeni, Goo Jun, Alan Kwong, Maria Valentini, Roberto Cusano, Manuela Oppo, Rosella Pilu and Brendan Tarrrier for additional help on generating and managing sequencing data; Mariano Dei, Monia Lobina and Francesca Deidda

for sample preparation; and Lidia Leoni, Carlo Podda and Antonio Concas for their technical support on the high-performance computing cluster at CRS4.

- 1 Li Y, Sidore C, Kang HM, Boehnke M, Abecasis GR: Low-coverage sequencing: implications for design of complex trait association studies. *Genome Res* 2011; **21**: 940–951.
- 2 Porcu E, Sanna S, Fuchsberger C, Fritsche LG: Genotype imputation in genome-wide association studies. *Curr Protoc Hum Genet* 2013; **Chapter 1**: Unit1.25.
- 3 Marchini J, Howie B: Genotype imputation for genome-wide association studies. *Nat Rev Genet* 2010; **11**: 499–511.
- 4 Li Y, Willer CJ, Ding J, Scheet P, Abecasis GR: MaCH: using sequence and genotype data to estimate haplotypes and unobserved genotypes. *Genet Epidemiol* 2010; **34**: 816–834.
- 5 Marchini J, Howie B, Myers S, McVean G, Donnelly P: A new multipoint method for genome-wide association studies by imputation of genotypes. *Nat Genet* 2007; **39**: 906–913.
- 6 Pilia G, Chen WM, Scuteri A *et al*: Heritability of cardiovascular and personality traits in 6,148 Sardinians. *PLoS Genet* 2006; **2**: e132.
- 7 Iacono WG, McGue M, Krueger RF: Minnesota Center for Twin and Family Research. *Twin Res Hum Genet* 2006; **9**: 978–984.
- 8 Miller MB, Basu S, Cunningham J *et al*: The Minnesota Center for Twin and Family Research genome-wide association study. *Twin Res Hum Genet* 2012; **15**: 767–774.
- 9 Voight BF, Kang HM, Ding J *et al*: The metabochip, a custom genotyping array for genetic studies of metabolic, cardiovascular, and anthropometric traits. *PLoS Genet* 2012; **8**: e1002793.
- 10 Cortes A, Brown MA: Promise and pitfalls of the ImmunoChip. *Arthritis Res Ther* 2011; **13**: 101.
- 11 Goldstein JI, Crenshaw A, Carey J *et al*: zCall: a rare variant caller for array-based genotyping: genetics and population analysis. *Bioinformatics* 2012; **28**: 2543–2545.
- 12 Scuteri A, Sanna S, Chen WM *et al*: Genome-wide association scan shows genetic variants in the FTO gene are associated with obesity-related traits. *PLoS Genet* 2007; **3**: e115.
- 13 Naitza S, Porcu E, Steri M *et al*: A genome-wide association scan on the levels of markers of inflammation in Sardinians reveals associations that underpin its complex regulation. *PLoS Genet* 2012; **8**: e1002480.
- 14 Vrieze SI, Feng S, Miller MB *et al*: Rare nonsynonymous exonic variants in addiction and behavioral disinhibition. *Biol Psychiatry* 2013; **75**: 783–789.
- 15 Sanna S, Pitzalis M, Zoledziwska M *et al*: Variants within the immunoregulatory CBLB gene are associated with multiple sclerosis. *Nat Genet* 2010; **42**: 495–497.
- 16 Pitzalis M, Zavattari P, Murru R *et al*: Genetic loci linked to type 1 diabetes and multiple sclerosis families in Sardinia. *BMC Med Genet* 2008; **9**: 3.
- 17 Orru V, Steri M, Sole G *et al*: Genetic variants regulating immune cell levels in health and disease. *Cell* 2013; **155**: 242–256.
- 18 Jun G, Wing MK, Abecasis GR, Kang HM: An efficient and scalable analysis framework for variant extraction and refinement from population scale DNA sequence data. *Genome Res* 2014 (submitted).
- 19 Delaneau O, Zagury JF, Marchini J: Improved whole-chromosome phasing for disease and population genetic studies. *Nat Methods* 2013; **10**: 5–6.
- 20 Su Z, Marchini J, Donnelly P: HAPGEN2: simulation of multiple disease SNPs. *Bioinformatics* 2011; **27**: 2304–2305.
- 21 Danecek P, Auton A, Abecasis G *et al*: The variant call format and VCFtools. *Bioinformatics* 2011; **27**: 2156–2158.
- 22 Nelson SC, Doheny KF, Pugh EW *et al*: Imputation-based genomic coverage assessments of current human genotyping arrays. *G3 (Bethesda)* 2013; **3**: 1795–1807.
- 23 Francioli LC, Menelaou A, Pulit SL *et al*: Whole-genome sequence variation, population structure and demographic history of the Dutch population. *Nat Genet* 2014; **46**: 818–825.
- 24 Lindquist KJ, Jorgenson E, Hoffmann TJ, Witte JS: The impact of improved microarray coverage and larger sample sizes on future genome-wide association studies. *Genet Epidemiol* 2013; **37**: 383–392.
- 25 Holm H, Gudbjartsson DF, Sulem P *et al*: A rare variant in MYH6 is associated with high risk of sick sinus syndrome. *Nat Genet* 2011; **43**: 316–320.

Supplementary Information accompanies this paper on European Journal of Human Genetics website (<http://www.nature.com/ejhg>)