1

## Sequencing and Imputation in GWAS: Cost-Effective Strategies to Increase Power and Genomic Coverage Across Diverse Populations

4

5  Corbin Quick,[1] Pramod Anugu,[2] Solomon Musani,[2] Scott T. Weiss,[3,4,5] Esteban G. Burchard,[6,7]
6  Marquitta J. White,[6] Kevin L. Keys,[6] and NHLBI Trans-Omics for Precision Medicine (TOPMed)
7  Consortium[†], Francesco Cucca,[8,9] Carlo Sidore,[8] Michael Boehnke,[1,‡] and Christian
8  Fuchsberger[1,10,11,‡,*]

9

10  [1] Department of Biostatistics and Center for Statistical Genetics, University of Michigan School
11  of Public Health, Ann Arbor, MI

12  [2] University of Mississippi Medical Center, Jackson, MS

13  [3] Harvard Medical School, Boston, MA

14  [4] Channing Department of Network Medicine, Brigham and Women's Hospital, Boston, MA

15  [5] Partners HealthCare Personalized Medicine, Boston, MA

16  [6] Department of Medicine, University of California San Francisco, San Francisco, California

17  [7] Department of Bioengineering and Therapeutic Sciences, University of California San
18  Francisco, San Francisco, California

19  [8] Istituto di Ricerca Genetica e Biomedica (IRGB), CNR, Monserrato, Italy
20  [9] Dipartimento di Scienze Biomediche, Università di Sassari, Sassari, Italy
21  [10] Division of Genetic Epidemiology, Department of Medical Genetics, Molecular and Clinical
22  Pharmacology, Medical University of Innsbruck, Innsbruck, Austria

23  [11] Institute for Biomedicine, Eurac Research, Affiliated Institute of the University of Lübeck,
24  Bolzano, Italy
25  [†] A complete list of NHLBI Trans-Omics for Precision Medicine (TOPMed) Consortium co-
26  investigators is provided in Supplemental Materials
27  [‡] These authors jointly supervised this work
28  * Correspondence: Christian.Fuchsberger@eurac.edu
29

**ABSTRACT**

A key aim for current genome-wide association studies (GWAS) is to interrogate the full spectrum of genetic variation underlying human traits, including rare variants, across populations. Deep whole-genome sequencing is the gold standard to capture the full spectrum of genetic variation, but remains prohibitively expensive for large samples. Array genotyping interrogates a sparser set of variants, which can be used as a scaffold for genotype imputation to capture variation across a wider set of variants. However, imputation coverage and accuracy depend crucially on the reference panel size and genetic distance from the target population.

Here, we consider a strategy in which a subset of study participants is sequenced and the rest array-genotyped and imputed using a reference panel that comprises the sequenced study participants and individuals from an external reference panel. We systematically assess how imputation quality and statistical power for association depend on the number of individuals sequenced and included in the reference panel for two admixed populations (African and Latino Americans) and two European population isolates (Sardinians and Finns). We develop a framework to identify powerful and cost-effective GWAS designs in these populations given current sequencing and array genotyping costs. For populations that are well-represented in current reference panels, we find that array genotyping alone is cost-effective and well-powered to detect both common- and rare-variant associations. For poorly represented populations, we find that sequencing a subset of study participants to improve imputation is often more cost-effective than array genotyping alone, and can substantially increase genomic coverage and power.

1    **INTRODUCTION**

2    Genome-wide association studies (GWAS) have detected thousands of common genetic

3    variants associated with hundreds of complex diseases and traits[1]. A key aim for the next wave of

4    GWAS is to interrogate the full spectrum of genetic variation underlying human genetic traits,

5    including rare (minor allele frequency [MAF] < 0.5%) variants, across a wide range of human

6    populations. Detecting association at rare variants requires both more comprehensive genomic

7    coverage and sufficient sample size. Deep whole genome sequencing (WGS) is the gold standard

8    method for capturing rare variation; however, even in the era of the $1,000 genome, large WGS

9    association studies remain prohibitively expensive. Genotype imputation has been a mainstay of

10   GWAS, providing increased genomic coverage from inexpensive array-based genotype call sets.

11   While initial imputation studies only surveyed common variants (e.g.[2]), larger and more diverse

12   reference panels now enable more accurate and comprehensive imputation of rare and low-

13   frequency (0.5% < MAF < 5%) variants across a wide range of populations (e.g.[3]).

14   Imputation algorithms model haplotypes in the study sample as mosaics of haplotypes in a

15   reference panel (e.g. from the International HapMap Project[4] or 1000 Genomes Project[5]) to predict

16   genotypes at untyped variants[6]. By increasing genomic coverage and accuracy, imputation

17   increases statistical power to detect association, enables more complete meta-analysis of results

18   from multiple studies, and facilitates the identification of causal variants through genetic fine-

19   mapping[6; 7]. Imputation coverage and accuracy depend crucially on the size of the reference panel

20   and the genetic distance between reference and target populations[6; 8]. The largest current broadly

21   available reference panels, e.g. from the Haplotype Reference Consortium[9] (HRC) and UK10K

22   project[10], include tens of thousands of predominantly European individuals. These panels provide

23   near complete imputation of genetic variation down to MAF~0.1% for many European

1  populations, but lower imputation quality for non-European and admixed populations and

2  population isolates, particularly for rare and low-frequency variants[11-13]. The 1000 Genomes

3  Project and HapMap panels include individuals from diverse worldwide populations, but provide

4  more limited imputation coverage and accuracy due to their smaller sample sizes.

5  Capturing rare variation across diverse populations is crucial to detect population

6  differences in genetic risk factors, accurately predict genetic risk, and identify causal variants and

7  biological mechanisms through trans-ethnic fine-mapping[14; 15]. Population-matched or multi-

8  ethnic reference panels can improve imputation quality and coverage for rare variants in GWAS

9  of diverse populations[11-13; 16-18]; this approach has enabled discovery of novel loci and refinement

10  of association signals for multiple populations and complex traits[12; 19; 20].

11  Here, we consider an approach in which a subset of study participants is whole genome

12  sequenced and the rest are array-genotyped and imputed using an augmented reference panel that

13  comprises the sequenced participants and individuals from an external reference panel[21; 22]. This

14  hybrid sequencing-and-imputation strategy provides more comprehensive coverage than array

15  genotyping alone, and is less costly than whole genome sequencing the entire sample. We and

16  others have used this strategy[18; 23-25], but no analysis of coverage, power, and cost-effectiveness

17  has been carried out to date. Here, we assess how imputation coverage and power to detect

18  association vary across genotyping arrays and as a functions of the number of population-matched

19  individuals sequenced and included in the reference panel for two admixed populations (African

20  Americans and Latino Americans) and two European population isolates (Sardinians and Finns)

21  to identify powerful and cost-effective GWAS strategies in these populations. We also describe an

22  interactive web-based tool to assist researchers in the design and planning of their own GWAS.

1    **MATERIALS AND METHODS**

2    We first describe WGS data sources used in our analysis. Next, we describe imputation

3    strategies, and outline procedures and imputation quality metrics to compare these strategies.

4    Finally, we present a novel method to estimate power for the sequencing-only, imputation-only,

5    and sequencing-and-imputation strategies. For ease of presentation, we assume a dichotomous trait

6    and a multiplicative disease model, although our findings generalize easily to continuous traits and

7    other genetic models.

8    **Data Resources**

9    We used WGS data on 11,920 individuals to assess imputation quality across reference

10    panel configurations and genotyping arrays for admixed populations and population isolates. For

11    our analysis of admixed populations, we used WGS data on 3,412 African Americans (participants

12    from the Jackson Heart Study) and 2,068 Latino Americans (participants of Puerto Rican and

13    Mexican descent from the GALA II study and Costa Rican descent from the Genetic Epidemiology

14    of Asthma in Costa Rica and CAMP studies) in the National Heart, Lung, and Blood Institute

15    (NHLBI) Trans-Omics for Precision Medicine (TOPMed) WGS program. For our analysis of

16    isolated populations, we used WGS data on 2,995 Finns (participants of the GoT2D, 1KGP, SISu,

17    and Kuusamo studies) and 3,445 Sardinians (participants of the SardiNIA study) in the HRC.

18    **Procedures to Evaluate Imputation Coverage and Accuracy**

19    We considered three imputation strategies: (1) using sequenced study participants as a

20    study-specific reference panel, (2) using an external reference panel alone (for this comparison,

21    the HRC or HRC subset excluding individuals from the target population), and (3) using an

22    augmented panel that comprises sequenced study participants and an external panel.

1   For African Americans, who are underrepresented in the current version 1.1 of the HRC,

2  we constructed population-specific and HRC-augmented reference panels with 0 to 2,000 African

3  Americans. For Latino Americans, we used the same approach but restricted the study-specific

4  panel size to $\leq$1,500 due to the more limited available sample of sequenced Latino American

5  individuals. For Finns and Sardinians, which are present in the HRC, we constructed augmented

6  reference panels that comprised the 29,470 non-Finnish or 29,020 non-Sardinian individuals in the

7  HRC together with 0 to 2,000 Finns or Sardinians from the HRC.

8

9  **Table 1. Genotyping Arrays used for Comparisons**

| Array | No. Marker Variants | List Cost per Sample[26] |
|---|---|---|
| Illumina Infinium Core | 307K | $49 |
| Illumina Infinium OmniExpress | 710K | $94 |
| Illumina Infinium Omni2.5 | 2.5M | $172 |

10

11   For each population, each imputation strategy, and each of three commonly-used

12  genotyping arrays (Table 1), we used sequence-based genotype calls at marker variants present on

13  the array as a scaffold for imputation using Minimac3, masking the remaining sequence-based

14  genotype calls[7]. We then compared the imputed genotype dosages to the true (masked) genotypes

15  to estimate (a) imputation $r^2$, the squared Pearson correlation between true genotype and imputed

16  dosage, and (b) imputation coverage, the proportion of variants with imputation $r^2 > 0.3$ and minor

17  allele count (MAC) $\geq$ 5 (the MAC threshold used by the HRC panel[9]) in the reference panel.

1     **Estimating Power to Detect Association using Empirical Imputation Quality Data**

2         When sequenced individuals are included in the reference panel, power calculations should

3     account for the interdependence between imputation $r^2$ and the number of participants sequenced

4     $n$, and for the possibility that the variant is not imputable (absent in the reference panel or not

5     imputed due to insufficient MAC, or filtered prior to association analysis due to imputation $r^2$

6     falling below a given threshold). While common variant associations are likely to be captured by

7     LD proxy SNPs even when the causal variant is not directly genotyped or imputed, rare variant

8     associations are much less likely to be captured by proxy SNPs[27]. Here, we assume that power to

9     detect association for variants that are not imputable is zero. This assumption affects power

10     calculations almost exclusively for rare variants, since common variants are almost uniformly

11     imputable with large reference panels[7; 9].

12         We assume that the $n$ participants who are sequenced are randomly subsampled from the

13     overall sample of $n + m$ study participants, and that test statistics are calculated separately for the

14     sequenced and imputed subsamples and combined using the effective sample size weighted meta-

15     analysis test statistic $Z_{nm} = c_{nm}^{1/2} Z_n^{seq} + (1 - c_{nm})^{1/2} Z_m^{imp}$, where $c_{nm} = n/(n + r^2 m)$. The

16     asymptotic distribution of $Z_{nm} - \eta \sqrt{n + r^2 m}$ is normal with mean 0 and variance 1, where $r^2$ is

17     the squared correlation between imputed dosages and true genotypes, and $\eta$ is an effect size

18     parameter which is equal to 0 under the null hypothesis of no association. The form of $\eta$ depends

19     on the association model (e.g. additive, dominant, multiplicative), relative risk or odds ratio, MAF,

20     and population prevalence and, for binary traits, the case-control ratio. Under an arbitrary

21     association model for binary traits, we can write

22
$$\eta = \frac{2(p_{case} - p_{control})}{\sqrt{(1 + s)\left(v_{case} + \frac{1}{s} v_{control}\right) + 4(p_{case} - p_{control})^2}}$$

1    where $p_{case}$ and $p_{control}$ are the alternate allele frequencies in the disease-positive and disease-

2    negative populations, $v_{case}$ and $v_{control}$ are the variances of genotypes in the disease-positive

3    and disease-negative populations, and $s$ is the GWAS case-control ratio.

4        To estimate power while accounting for variability in imputation $r^2$ and the possibility that

5    a variant is not imputable, we average empirical imputation $r^2$ values and MACs across variants

6    from experiments with real data described in the previous section. Specifically, we estimate power

7    to detect association when $n$ individuals are sequenced and $m$ are genotyped and imputed as

$$\widehat{Power}(m,n) = \frac{1}{\sum_j w_j^{MAF} w_j^{PS}} \sum_j w_j^{MAF} w_j^{PS} C_{nj} \int_{-z_{1-\alpha/2}}^{z_{1-\alpha/2}} \phi\left(u - \eta\sqrt{n + r_{nj}^2 m}\right) du$$

9    where $\phi(u) = e^{-\frac{u^2}{2}}/\sqrt{2\pi}$ is the standard normal density function, $z_{1-\alpha/2}$ is the $\alpha$-level

10    significance threshold, $r_{nj}^2$ is the imputation $r^2$ value for the $j^{th}$ variant, $C_{nj} = I(MAC_{nj}^{panel} \geq$

11    $5, r_{nj}^2 \geq 0.3)$ is an indicator equal to 1 if the $j^{th}$ variant was imputable and 0 otherwise, and

12    $MAC_{nj}^{panel}$ is the reference panel MAC for the $j^{th}$ variant when the $n$ sequenced individuals from

13    the target population were included in the reference panel.

14        We define the first weight term $w_j^{MAF} = P_N^{GWAS}(\hat{p}_j)/\hat{P}_N(\hat{p}_j)$, where $N$ is the total number

15    of samples used in our analysis for the given population (e.g. $N =$3,412 for African Americans),

16    $\hat{p}_j$ is the sample MAF for the $j^{th}$ variant in the total sample, $\hat{P}_N(x)$ is the proportion of variants

17    with MAF $= x$, and $P_N^{GWAS}(x)$ is the probability of observing sample MAF $= x$ in a sample of size

18    $N$ given the specified association model. For example, in a GWAS with sample size $N$ and case-

19    control ratio $s$, the sample MAC (which is equal to $2N\hat{p}$, where $\hat{p}$ is the sample MAF) is

20    approximately Poisson distributed with mean $2N(sp_{case} + p_{control})/(s + 1)$, where $p_{case} =$

21    $p\gamma/[1 + p(\gamma - 1)]$ and $p_{control} = (p - Kp_{case})/(1 - K)$ for a variant with population MAF $p$

8

1    and relative risk $\gamma$ for a disease with prevalence $K$. This weighting approach adjusts for differences

2    between the empirical distribution of MACs across variants in real data, and the theoretical MAC

3    distribution for a variant with the specified MAF, effect size, prevalence in a GWAS with sample

4    size $N$ and case-control ratio $s$.

5        The second weighting term $w_j^{\text{PS}}$ accounts for the probability that a variant with the

6    specified population MAF $p$ is population-specific (monomorphic outside the target population),

7    and is defined

8
$$w_j^{PS} = \begin{cases} \hat{P}_{PS}(p)/\hat{P}_{PS}(\hat{p}_j), & \text{variant } j \text{ is population-specific,} \\ [1 - \hat{P}_{PS}(p)]/[1 - \hat{P}_{PS}(\hat{p}_j)], & \text{otherwise,} \end{cases}$$

9    where $\hat{P}_{PS}(x)$ is the fraction of variants that are population-specific among variants with MAF=$x$

10   in the target population. This adjustment factor ensures that the weight assigned to population-

11   specific variants in power calculations reflects the probability that a variant with the specified

12   population MAF $p$ is population-specific.

13

**RESULTS**

First, we compare strategies to improve imputation using study-specific WGS data for African Americans, Latino Americans, Sardinians, and Finns. Next, we assess the effects of genotyping array on imputation quality and coverage for each population and reference panel. We then use these results to estimate statistical power to detect association as a function of study-specific panel size, number of participants imputed, external reference panel, and genotyping array. Finally, we identify cost-effective study designs by comparing statistical power and total experimental (sequencing and genotyping) costs for sequencing-only, imputation-only, and sequencing-and-imputation GWAS designs for each population and genotyping array.

**Strategies to Improve Imputation using Study-Specific WGS Data**

We compared imputation $r^2$ and coverage (proportion of variants with imputation $r^2 > 0.3$ and reference MAC $\geq 5$) for three imputation strategies: (1) using an external reference panel (the HRC or HRC subset) alone, (2) using an augmented reference panel that combines the study-specific and external panels, and (3) using a study-specific reference panel alone.

The external panel alone (HRC for Latino Americans and African Americans, and HRC subset that excludes individuals from the target population for Finns and Sardinians) provided 96% imputation coverage for MAF $\geq 0.25\%$ variants (where MAF is calculated separately within each population) for Finns, 84% coverage for Sardinians, 86% coverage for Latino Americans, and 77% coverage for African Americans (Figure 1, top row). The relatively lower coverage for African Americans is expected since the HRC consists primarily of Central and Northern Europeans, who are genetically closer to Finns and Sardinians, and includes relatively few Africans or African Americans. Despite the small number of Latino or Native Americans included in the HRC, imputation coverage was slightly higher for Latino Americans than for Sardinians. This may

1    reflect the high degree of European admixture in many Latino American populations[28], and the

2    abundance of population-specific rare and low-frequency variants in the Sardinian population[24].

3         Augmenting an external reference panel with even a relatively small number of sequenced

4    individuals substantially increased coverage, particularly for African Americans and Sardinians,

5    and for variants with lower MAF. For example, augmenting the external panel with 500 sequenced

6    individuals from the study population improved overall imputation coverage for MAF=0.25-0.5%

7    variants by 4% for Finns, 9% for Latino Americans, 16% for African Americans, and 23% for

8    Sardinians genotyped using the OmniExpress relative to the external panel alone (Figure 1).

9    Similarly, augmenting the external reference panel with even 200 individuals increased imputation

10   coverage for MAF=0.1-0.25% variants by 3%, 4%, 6%, 10% relative to the external panel alone

11   for Finns, Latino Americans, African Americans, and Sardinians, respectively.

12        With 2,000 individuals from the target population (or 1,500 for Latino Americans),

13   population-specific panels provided roughly equivalent imputation $r^2$ compared to augmented

14   panels (Supplemental Figure 1A); however, augmented panels provided higher imputation

15   coverage overall for low MAF variants (Supplemental Figure 1B). For example, augmented panels

16   with 2,000 individuals from the target population (or 1,500 for Latino Americans) provided 86%,

17   80%, 79%, and 86% coverage for 0.1-0.25% MAF variants for Finns, Latino Americans, African

18   Americans, and Sardinians respectively, whereas population-specific panels alone provided 72%,

19   51%, 78%, and 72% coverage using the Omni Express array. However, imputation coverage for

20   variants with MAF>0.25% differed by <1% between augmented and population-specific panels

21   with 2,000 individuals from the target population (or 1,500 for Latino Americans) for all

22   populations and genotyping arrays. When a smaller number (less than 500) of individuals from the

23   target population are sequenced, augmented reference panels provided substantially higher

1    imputation coverage and $r^2$ than population-specific panels alone. For example, augmented panels

2    with 500 individuals from the target population provided 90%, 85%, 65%, and 85% coverage for

3    0.25-0.5% MAF variants for Finns, Latino Americans, African Americans, and Sardinians

4    respectively, whereas population-specific panels of 500 individuals provided ≤30% coverage

5    using the Omni Express array.

6        Even very rare variants (MAF=0.1-0.25%) attained high coverage across all populations

7    given a sufficient number of population-matched individuals in the reference panel. For example,

8    attaining ≥70% imputation coverage for MAF=0.1-0.25% variants required a study-specific panel

9    of ≥1,800 individuals for African Americans, 1,000 for Latino Americans, 700 for Sardinians, and

10   0 for Finns using the OmniExpress. These increases in imputation coverage primarily reflect

11   increasing numbers of population-specific variants captured in the reference panel, which are

12   absent from or present in low copy number in the external panel.

13   **Imputation Coverage and Quality across Genotyping Arrays**

14       Imputation coverage was generally similar for the OmniExpress and Omni2.5 arrays, but

15   consistently lower for the less dense Core array. Coverage differed by <7% between the

16   OmniExpress and Omni2.5 across all MAF bins, populations, and reference panels, whereas the

17   Core provided up to 24% lower coverage than the Omni2.5 (Figure 1, upper panels). Imputation

18   coverage was more heterogeneous across arrays for populations with greater genetic distance from

19   the external reference panel (e.g., African Americans and the HRC panel), particularly with smaller

20   (or absent) study-specific panels. Because we used the same reference panels for each genotyping

21   array, differences in imputation coverage between arrays are solely due to differences in the

22   proportion of variants that attained imputation $r^2 \geq 0.3$. Imputation $r^2$ varied more across

23   genotyping arrays than did imputation coverage (Figure 1, lower versus upper panels); however,

12

1  **Figure 1. Imputation Quality by Population and Genotyping Array.**



2

3  *Imputation coverage (upper panels) and mean imputation $r^2$ (lower panels) as functions of the number of population-matched*
4  *individuals included in augmented reference panels (Number Sequenced, x-axis). Here and elsewhere, MAF is calculated separately*
5  *within each population.*
6

1     the magnitude of differences in imputation $r^2$ between arrays was still generally modest,

2     particularly for the Finns and Sardinians.

3     **Powerful and Cost-Effective Strategies for GWAS across Populations**

4     We compared the cost-effectiveness of sequencing-only, imputation-only, and sequencing-

5     and-imputation strategies by analyzing statistical power to detect association as a function of

6     numbers of study participants sequenced and imputed, genotyping array, and reference panel

7     across a range of genetic models. Here, we define the most *cost-effective* strategy as either (1)

8     minimizing total experimental (sequencing and genotyping) cost while attaining power at or above

9     a given threshold, or equivalently (2) maximizing power while maintaining cost no greater than a

10    specified constraint.

11    The cost-effectiveness of sequencing a subset of study participants varied greatly across

12    populations. For Finns, imputation-only designs were most powerful to detect association and

13    adding sequenced individuals increased power only minimally, even for low-frequency and rare

14    variants. For Sardinians, Latino Americans, and African Americans, sequencing a subset of study

15    participants was optimal, and often achieved substantially greater power than imputation-only or

16    sequencing-only studies. For example, a GWAS of African Americans with equal numbers of

17    cases and controls in which 400 participants are sequenced and 11,100 are imputed using the

18    Illumina Infinium Core array has 90% power to detect a risk variant with MAF = 0.5% and RR =

19    4 for a disease with prevalence 1%, whereas an imputation-only GWAS with the same total cost

20    (19,250 participants) has only 68% power (Figure 2). Even for populations in which optimal

21    sequencing-and-imputation designs had substantially greater power than imputation-only, the

22    optimal number to sequence was often modest. For example, only 210 participants are sequenced

23    under the optimal design using the Illumina OmniExpress to attain 80% power in the previous

1     example (Figure 3). This is expected because even a relatively small study-specific panel can

2     substantially increase imputation coverage (Figure 1, upper panels).

3     **Denser Genotyping Arrays vs. Sequencing: Which is More Cost-Effective to Increase Power?**

4     Imputation coverage and power to detect association can be increased by using denser

5     genotyping arrays, which provide a more informative framework for imputation, or by sequencing

6     population-matched individuals and augmenting the reference panel. We assessed the cost-

7     effectiveness of these two strategies by comparing power to detect association across genotyping

8     arrays for study designs that have the same total cost assuming $1000 for WGS and current list

9     prices for genotyping arrays (Table 1). As expected, the optimal number of participants sequenced

10    to maximize power given fixed total cost generally decreased with increasing array density. For

11    example, the optimal number sequenced to maximize power to detect association was 500, 300,

12    and 90 for the Infinium Core, OmniExpress, and Omni2.5 respectively for Sardinians given total

13    sequencing and genotyping budget of $2M for a risk variant with $RR = 2$, $MAF = 1\%$, and disease

14    prevalence 1%. Power to detect association under the optimal design given a fixed total cost was

15    generally greater for sparser arrays; in the previous example, power under the optimal design was

16    98%, 91%, and 55% for the Infinium Core, OmniExpress, Omni2.5.

17    We also compared optimal designs to attain power above a given threshold at minimum

18    total cost across genotyping arrays based on the per-sample array genotyping costs reported in

19    Table 1. Generally, sparser arrays were more cost-effective (reached the power threshold with

20    lower total cost) than dense arrays. In fact, the sparsest genotyping array in our analysis, the

21    Infinium Core, was most cost-effective across all disease models and populations apart from

22    African Americans, for whom the Infinium OmniExpress was most cost-effective for some rare-

23    variant disease models. This last result is unsurprising given the substantial difference in

15

## Figure 2. Power and Optimal Design by Population and Genotyping Array.



*Power to detect association for case-control studies with equal numbers of cases and controls as a function of sequenced subsample size (x-axis) and imputed subsample size (y-axis) for a variant with MAF 0.5% and relative risk 4 for a disease with prevalence 1%. Axes are scaled to reflect costs of genotyping arrays (Table 1) and sequencing ($1K per sample). Dashed diagonal lines indicate study designs with the same total cost, given by $y = a - bx$ where $a = (Total\ Cost)\ /(Array\ Cost)$ and $b = (Sequencing\ Cost)/(Array\ Cost)$. Circled points indicate optimal study designs, which attain the indicated power level at minimum total experimental cost (or, maximize power at the indicated total experimental cost), shown only for optimal designs with total genotyping cost ≤ $2M ($1.5M for Latino Americans).*

**Figure 3. Power as a Function of Minor Allele Frequency and Effect Size.**



*Statistical power (y-axis) to detect a rare large-effect variant (MAF=0.25%, RR=3; top row) and common modest-effect variant (MAF=5%, RR=1.3; bottom row) for a disease with prevalence 1% as a function of the number of participants array-genotyped and imputed (x-axis) when 0, 500, or 2,000 participants are sequenced and included in an augmented reference panel. The number of participants sequenced has a far greater impact on statistical power for the rare variant association. Importantly, statistical power is bounded above by the probability that the variant is imputable ($r^2 > 0.3$ and reference MAC $\geq 5$), causing power to asymptote below 1 as a function of the number of imputed participants (e.g., upper-left panel).*

1    imputation coverage between the Infinium Core and Omni arrays for African Americans (Figure

2    1). Importantly, our analysis assumes 1) a direct trade-off between the GWAS sample size and

3    sequencing/array genotyping costs, and 2) no additional costs per GWAS sample other than

4    sequencing/genotyping. Under these assumptions, we found that denser arrays are generally less

5    cost-effective than sparser arrays; of course, denser arrays provide higher imputation coverage

6    given a fixed GWAS sample size.

**Optimal Study Design as a Function of Minor Allele Frequency and Effect Size**

8        Power to detect association under a given study design depends on MAF, effect size

9    (relative risk or odds ratio), and population prevalence[29]. These parameters also influence the

10    relative cost-effectiveness of sequencing and imputation. While common variants can be

11    accurately imputed with small reference panels, large population-matched reference panels are

12    needed to capture rare (population-specific) variants. In Figure 3, we illustrate the impact of

13    sequencing on statistical power for two combinations of MAF and effect size in each of the four

14    study populations.

15        The optimal percentage of study participants sequenced to attain ≥80% power to detect

16    association at minimum total cost increases with decreasing MAF (Figure 4). This is expected,

17    since larger reference panels are needed to capture variants with lower frequency. Finally, the

18    optimal percentage of study participants sequenced to attain ≥80% power decreases with

19    increasing effect size magnitude. This is expected, since the expected number of risk alleles

20    captured in the reference panel increases with effect size magnitude.

1 **Figure 4. Optimal Design as a Function of Minor Allele Frequency and Effect Size.**



2

3 *Percentage of participants sequenced (x-axis) and total sample size (y-axis) under optimal designs to attain statistical power ≥80% for*
4 *rare and common variants across two effect size values for each of the four study populations using the Infinium Core array. Here,*
5 *effect size refers to the χ-squared non-centrality parameter (NCP) for single-variant association tests given perfect genotype accuracy,*
6 *which is defined as $\eta^2$ in Methods. Relative risk (RR) values corresponding to each combination of MAF and NCP are indicated in the*
7 *far-right panel (for Sardinians). With NCP held constant, differences in optimal design for different MAF values are solely due to*
8 *differences in imputation coverage and quality across the MAF spectrum.*

## DISCUSSION

While the cost of genome sequencing has fallen dramatically[29], large genome sequencing studies remain prohibitively expensive. Large imputation reference panels are now enabling accurate imputation of even very rare variants (MAF>0.001)[9; 13; 30], making imputation-based GWAS viable and cost-effective for detecting associations across much of the allele frequency spectrum. For populations with limited reference panel data, we have shown that sequencing a subset of study participants can substantially increase imputation coverage and accuracy, particularly for rare and population-specific variants, at a fraction of the cost of sequencing the entire study cohort. Our results also suggest that it is almost always advantageous to augment existing reference panels, except when the study-specific sequenced panel is large or the target population has high genetic distance from the external panel.

Complementary sequencing-and-imputation GWAS strategies have been applied to refine association signals and discover novel associations for several populations and complex traits[12; 19; 20]. While most sequencing-and-imputation studies to date have been carried out in European isolated populations, our results suggest that this strategy can also be powerful and cost-effective for admixed and non-European populations. In addition to increasing genomic coverage and power to detect association for the study itself, sequencing a subset of study participants provides a data resource that can be used to enhance imputation in future studies of the same or related populations so long as the sequence data can be shared.

Directly augmenting an existing reference panel with study-specific sequence data is not always feasible due to technical, logistical, and privacy constraints. However, we and others have found that the distributed reference panel approach (separately imputing with two or more reference panels and combining the results) provides nearly equivalent imputation quality

1    (Supplemental Figure 2). Thus, study-specific WGS data can be used to improve imputation even

2    when directly augmenting an external panel is not feasible.

3        While large reference panels enable accurate imputation across a wide range of the allele

4    frequency spectrum[9; 13], the extent of genetic variation that can be captured through imputation is

5    limited relative to WGS. For example, *de novo* mutations cannot be imputed regardless of

6    reference panel size. This is particularly salient for monogenic disorders; for example, over 80%

7    of achondroplasia cases occur from recurrent *de novo* mutations in *FGFR3*[31]. Thus, imputation

8    may be unable to detect causative alleles for traits with extreme genetic architectures, even with

9    very large reference panels.

10        As increasingly large and diverse sequencing projects are conducted, larger and more diverse

11    reference panels will become available. In the design and planning of GWAS, it may be prudent

12    to consider resources under development and pending release in addition to resources that are

13    currently available. More broadly, our analysis highlights the utility of collaboration and

14    coordination across institutions for effective study design and resource allocation. For example,

15    the optimal design to maximize power in an individual study does not necessarily maximize meta-

16    analysis power across multiple studies of the same trait and population.

17        Our analysis of cost-effectiveness and optimal design depends crucially on the relative per-

18    sample costs of sequencing and array genotyping. Both sequencing and array genotyping costs

19    have fallen markedly in recent years, and are likely to continue to do so. Depending on the relative

20    rates of change, cost-effectiveness and optimal design also may change. In addition, the cost of

21    participant recruitment and DNA sample collection may alter the relative cost-effectiveness of

22    sequencing and genotyping. Finally, our cost-effectiveness analysis assumes that sample size is

23    unconstrained; this may not apply for small populations or rare diseases.

21

1      While our results are illustrative, investigators may wish to explore questions of the relative

2      cost-effectiveness of sequencing and array genotyping strategies in the context of their own study

3      and relevant assumptions about population, reference panels, and sequencing and array genotyping

4      costs. To enable this exploration, we have developed a flexible, easy-to-use tool, APSIS (Analysis

5      of Power for Sequencing and Imputation Studies), which is open source and freely available (see

6      Web Resources).

7 **Conclusions**

8      Here, we assessed the genomic coverage, statistical power, and cost-effectiveness of

9      sequencing and imputation-based designs for GWAS in four populations across a range of genetic

10     models. We developed a novel method to account for available reference haplotype data in power

11     calculations using empirical data, which can be applied to inform GWAS planning and design. For

12     European populations that are well-represented in current reference panels, our results suggest that

13     imputation-based GWAS is cost-effective and well-powered to detect both common- and rare-

14     variant associations. For populations with limited representation in current reference panels, we

15     found that sequencing a subset of study participants can substantially increase genomic coverage

16     and power to detect association, particularly for rare and population-specific variants. Our results

17     also suggest that larger and more diverse reference panels will be important to facilitate array-

18     based GWAS in global populations.

19 **WEB RESOURCES**

20 APSIS (Analysis of Power for Sequencing-and-Imputation Studies):

21 http://github.com/corbinq/APSIS

22

24

1    MD, Adam Davis, MA, MPH, Michael A. LeNoir, MD, Kelley Meade, MD, Saunak Sen, PhD and

2    Fred Lurmann, MS. The authors also wish to thank the staffs and participants contributed to the

3    GALA II study.

25

6    **Disclaimer**

7        The views expressed in this manuscript are those of the authors and do not necessarily

8    represent the views of the National Heart, Lung, and Blood Institute; the National Institutes of

9    Health; or the U.S. Department of Health and Human Services.

10

1 **REFERENCES**

2 1. MacArthur, J., Bowler, E., Cerezo, M., Gil, L., Hall, P., Hastings, E., Junkins, H., McMahon,
3     A., Milano, A., and Morales, J. (2016). The new NHGRI-EBI Catalog of published
4     genome-wide association studies (GWAS Catalog). Nucleic acids research 45, D896-
5     D901.

6 2. Scott, L.J., Mohlke, K.L., Bonnycastle, L.L., Willer, C.J., Li, Y., Duren, W.L., Erdos, M.R.,
7     Stringham, H.M., Chines, P.S., Jackson, A.U., et al. (2007). A genome-wide association
8     study of type 2 diabetes in Finns detects multiple susceptibility variants. Science 316,
9     1341-1345.

10 3. Mahajan, A., Wessel, J., Willems, S.M., Zhao, W., Robertson, N.R., Chu, A.Y., Gan, W.,
11     Kitajima, H., Taliun, D., and Rayner, N.W. (2018). Refining the accuracy of validated
12     target identification through coding variant fine-mapping in type 2 diabetes. Nature
13     genetics 50, 559.

14 4. Consortium, I.H. (2010). Integrating common and rare genetic variation in diverse human
15     populations. Nature 467, 52.

16 5. Consortium, G.P. (2015). A global reference for human genetic variation. Nature 526, 68.

17 6. Li, Y., Willer, C., Sanna, S., and Abecasis, G. (2009). Genotype imputation. Annual review of
18     genomics and human genetics 10, 387-406.

19 7. Das, S., Forer, L., Schonherr, S., Sidore, C., Locke, A.E., Kwong, A., Vrieze, S.I., Chew,
20     E.Y., Levy, S., McGue, M., et al. (2016). Next-generation genotype imputation service
21     and methods. Nat Genet 48, 1284-1287.

22 8. Roshyara, N.R., and Scholz, M. (2015). Impact of genetic similarity on imputation accuracy.
23     BMC Genet 16, 90.

24 9. McCarthy, S., Das, S., Kretzschmar, W., Delaneau, O., Wood, A.R., Teumer, A., Kang, H.M.,
25     Fuchsberger, C., Danecek, P., Sharp, K., et al. (2016). A reference panel of 64,976
26     haplotypes for genotype imputation. Nat Genet 48, 1279-1283.

27 10. consortium, U.K. (2015). The UK10K project identifies rare variants in health and disease.
28     Nature 526, 82.

29 11. Deelen, P., Menelaou, A., Van Leeuwen, E.M., Kanterakis, A., Van Dijk, F., Medina-
30     Gomez, C., Francioli, L.C., Hottenga, J.J., Karssen, L.C., and Estrada, K. (2014).

Improved imputation quality of low-frequency and rare variants in European samples using the 'Genome of The Netherlands'. European Journal of Human Genetics 22, 1321.

12. Pistis, G., Porcu, E., Vrieze, S.I., Sidore, C., Steri, M., Danjou, F., Busonero, F., Mulas, A., Zoledziewska, M., Maschio, A., et al. (2015). Rare variant genotype imputation with thousands of study-specific whole-genome sequences: implications for cost-effective study designs. Eur J Hum Genet 23, 975-983.

13. Zhou, W., Fritsche, L.G., Das, S., Zhang, H., Nielsen, J.B., Holmen, O.L., Chen, J., Lin, M., Elvestad, M.B., and Hveem, K. (2017). Improving power of association tests using multiple sets of imputed genotypes from distributed reference panels. Genetic epidemiology 41, 744-755.

14. Kichaev, G., and Pasaniuc, B. (2015). Leveraging Functional-Annotation Data in Trans-ethnic Fine-Mapping Studies. Am J Hum Genet 97, 260-271.

15. Popejoy, A.B., and Fullerton, S.M. (2016). Genomics is failing on diversity. Nature News 538, 161.

16. Ahmad, M., Sinha, A., Ghosh, S., Kumar, V., Davila, S., Yajnik, C.S., and Chandak, G.R. (2017). Inclusion of Population-specific Reference Panel from India to the 1000 Genomes Phase 3 Panel Improves Imputation Accuracy. Sci Rep 7, 6733.

17. Lencz, T., Yu, J., Palmer, C., Carmi, S., Ben-Avraham, D., Barzilai, N., Bressman, S., Darvasi, A., Cho, J., and Clark, L. (2017). High-depth whole genome sequencing of a large population-specific reference panel: Enhancing sensitivity, accuracy, and imputation. bioRxiv, 167924.

18. Van Leeuwen, E.M., Kanterakis, A., Deelen, P., Kattenberg, M.V., Abdellaoui, A., Hofman, A., Schönhuth, A., Menelaou, A., de Craen, A.J., and van Schaik, B.D. (2015). Population-specific genotype imputations using minimac or IMPUTE2. Nature protocols 10, 1285.

19. Auer, P.L., and Lettre, G. (2015). Rare variant association studies: considerations, challenges and opportunities. Genome Med 7, 16.

20. Holm, H., Gudbjartsson, D.F., Sulem, P., Masson, G., Helgadottir, H.T., Zanon, C., Magnusson, O.T., Helgason, A., Saemundsdottir, J., and Gylfason, A. (2011). A rare variant in MYH6 is associated with high risk of sick sinus syndrome. Nature genetics 43, 316.

21. Hu, Y.J., Li, Y., Auer, P.L., and Lin, D.Y. (2015). Integrative analysis of sequencing and array genotype data for discovering disease associations with rare mutations. Proc Natl Acad Sci U S A 112, 1019-1024.

22. Zeggini, E. (2011). Next-generation association studies for complex traits. Nat Genet 43, 287-288.

23. Fuchsberger, C., Flannick, J., Teslovich, T.M., Mahajan, A., Agarwala, V., Gaulton, K.J., Ma, C., Fontanillas, P., Moutsianas, L., McCarthy, D.J., et al. (2016). The genetic architecture of type 2 diabetes. Nature 536, 41-47.

24. Sidore, C., Busonero, F., Maschio, A., Porcu, E., Naitza, S., Zoledziewska, M., Mulas, A., Pistis, G., Steri, M., Danjou, F., et al. (2015). Genome sequencing elucidates Sardinian genetic architecture and augments association analyses for lipid and blood inflammatory markers. Nat Genet 47, 1272-1281.

25. Steinthorsdottir, V., Thorleifsson, G., Sulem, P., Helgason, H., Grarup, N., Sigurdsson, A., Helgadottir, H.T., Johannsdottir, H., Magnusson, O.T., Gudjonsson, S.A., et al. (2014). Identification of low-frequency and rare sequence variants associated with elevated or reduced risk of type 2 diabetes. Nat Genet 46, 294-298.

26. Illumina, I. (2018). Microarray kits for genotyping and epigenetic analysis. *https://wwwilluminacom/products/by-type/microarray-kitshtml*

27. Montpetit, A., Nelis, M., Laflamme, P., Magi, R., Ke, X., Remm, M., Cardon, L., Hudson, T.J., and Metspalu, A. (2006). An evaluation of the performance of tag SNPs derived from HapMap in a Caucasian population. PLoS genetics 2, e27.

28. Bryc, K., Velez, C., Karafet, T., Moreno-Estrada, A., Reynolds, A., Auton, A., Hammer, M., Bustamante, C.D., and Ostrer, H. (2010). Genome-wide patterns of population structure and admixture among Hispanic/Latino populations. Proceedings of the National Academy of Sciences, 200914618.

29. Sham, P.C., and Purcell, S.M. (2014). Statistical power and significance testing in large-scale genetic studies. Nat Rev Genet 15, 335-346.

30. Mahajan, A., Taliun, D., Thurner, M., Robertson, N.R., Torres, J.M., Rayner, N.W., Steinthorsdottir, V., Scott, R.A., Grarup, N., and Cook, J.P. (2018). Fine-mapping of an expanded set of type 2 diabetes loci to single-variant resolution using high-density imputation and islet-specific epigenome maps. bioRxiv, 245506.

1    31. Bellus, G.A., Hefferon, T.W., de Luna, R.O., Hecht, J.T., Horton, W.A., Machado, M.,

2         Kaitila, I., McIntosh, I., and Francomano, C.A. (1995). Achondroplasia is defined by

3         recurrent G380R mutations of FGFR3. American journal of human genetics 56, 368.

4