

1 RAREsim: A simulation method for very rare genetic variants

2 Megan Null^{1,2*}, Josée Dupuis³, Christopher R. Gignoux^{4,5}, Audrey E. Hendricks^{1,4,5}

3 ¹*Mathematical and Statistical Sciences, University of Colorado Denver*; ²*Mathematics and Physical Sciences, The*
4 *College of Idaho* ³*Department of Biostatistics, Boston University School of Public Health*; ⁴*Human Medical Genetics*
5 *and Genomics Program, University of Colorado Anschutz Medical Campus*; ⁵*Colorado Center for Personalized*
6 *Medicine, University of Colorado Anschutz Medical Campus*

7
8 *Corresponding author email: mnull@collegeofidaho.edu

1 **Abstract**

2 Identification of rare variant associations is crucial to fully characterize the genetic architecture
3 of complex traits and diseases. Essential in this process is the evaluation of novel methods in
4 simulated data that mirrors the distribution of rare variants and haplotype structure in real data.
5 Additionally, importing real variant annotation enables in silico comparison of methods that
6 focus on putative causal variants, such as rare variant association tests, and polygenic scoring
7 methods. Existing simulation methods are either unable to employ real variant annotation or
8 severely under- or over-estimate the number of singletons and doubletons reducing the ability to
9 generalize simulation results to real studies. We present RAREsim, a flexible and accurate rare
10 variant simulation algorithm. Using parameters and haplotypes derived from real sequencing
11 data, RAREsim efficiently simulates the expected variant distribution and enables real variant
12 annotations. We highlight RAREsim's utility across various genetic regions, sample sizes,
13 ancestries, and variant classes.

1 **Introduction**

2 Studies of rare variants are important to gain a full understanding of the genetics of health
3 and disease, informing targeted drug development and precision medicine. Rare variants (minor
4 allele frequency (MAF) <1%) have been associated with traits across many diseases including
5 cancer, kidney, neurodevelopmental, cardiovascular, and infectious disease. With decreasing
6 sequencing costs, rare variant data are increasingly accessible¹ resulting in large sequencing
7 studies (e.g. >35,000; 45,000; and 70,000 subjects), databases including the UKBiobank,
8 GenomeAsia, and NIH programs such as the Genome Sequencing Program (GSP) and Trans-
9 Omics for Precision Medicine (TOPMed). Rare variant methods continue to be developed (e.g.
10 SKAT-O, iECAT, ProxECAT, and ACAT) to take advantage of the ever increasing sequencing
11 data.

12 Simulation studies enable evaluation of methods and study design (e.g. power and sample
13 size estimates) in known and controlled settings. Simulations that do not adequately mirror
14 essential properties of real data may have issues generalizing to real data, potentially resulting in
15 incorrect conclusions of method efficacy or power. In general, four qualities are necessary to
16 emulate in rare variant simulations of a genetic region: (1) allele frequency spectrum (AFS), (2)
17 total number of variants, (3) haplotype structure, and (4) variant annotation. To our knowledge,
18 no rare variant simulation method currently exists that incorporates all four qualities.

19 (1) **AFS** is the distribution of variant allele frequencies within a genetic region. Numerous
20 studies have shown that the AFS is skewed towards very rare variants with the vast
21 majority of variants being singletons and doubletons²⁻⁴.

1 (2) The **total number of variants**, especially very rare variants, differs by ancestry and
2 sample size. The total number of known variants is expected to increase as more
3 ancestrally diverse and larger samples are sequenced^{5,6}.

4 (3) **Haplotype structure**, the linkage disequilibrium (LD) and probability that rare single
5 nucleotide variants (SNVs) appear on the same haplotype background, varies across the
6 genome and by ancestry.

7 (4) **Variant annotation** is often used in rare variant methods and is thus essential for
8 accurate evaluation of those methods. For instance, weighting functional variants has
9 been shown to increase power to detect rare variant association in a gene region^{1,7,8}.

10 Other variant annotation, such as association with disease, is used to evaluate pleiotropy
11 and in genetic correlation and polygenic risk scores. A great variety of variant annotation
12 exists such as functional consequences, conservation score, chromatin state, eQTLs,
13 epigenetic information, and prior disease associations, among others⁹. While *in silico*
14 simulation of variant annotation can capture and emulate some annotation patterns,
15 simulations derived from real data can easily incorporate precise empirical patterns from
16 multiple annotation types, even those unique to a specific genetic region of interest,
17 providing a more direct link between simulations and real data.

18

19 Population genetics simulation methods, such as Wright-Fisher^{10,11} and coalescent¹²,
20 require only demographic and recombination information as input and often achieve an AFS,
21 total number of variants, and LD structure similar to real data. However, these methods can be
22 extremely computationally expensive or are not designed to emulate existing genetic regions
23 resulting in an inability to use real variant annotations. Alternatively, resampling methods create

1 haplotype mosaics from real genetic data using techniques that mimic recombination and
2 mutations, maintaining the ability to use existing annotations. These methods, such as
3 HAPGEN2¹³ derived from the original work of Li and Stephens¹⁴, are relatively computationally
4 efficient and maintain the appropriate AFS, expected number of variants, and haplotype structure
5 when simulating common variants^{13,15}. However, as we and others¹⁶ show, HAPGEN2 does not
6 simulate the correct total number or AFS for rare variants simulating too few rare and very rare
7 variants (e.g. singletons and doubletons). There is currently no available software to simulate
8 rare variant genetic data with a realistic AFS while retaining variant annotation.

9 To address this gap, we present RAREsim, a flexible and scalable genetic simulation
10 method designed for accurate simulation of rare variants. We assess and show the utility of
11 RAREsim across a variety of genetic regions, datasets, ancestries, and sample sizes. We provide
12 RAREsim as an R package to enable easy implementation and appropriate simulation of rare
13 variant data.

14

15 **Results**

16 *Algorithm*

17 RAREsim uses two primary datasets: input simulation data and target data. The *input*
18 *simulation dataset* is a sample of haplotypes (e.g. 1000 Genomes haplotypes³) with minor alleles
19 coded as 1 and all reference alleles, including monomorphic bases, coded as 0. The *target*
20 *dataset* is summary level data used to estimate RAREsim parameters. The target data has two
21 components: the allele count at each variant and the total number of variants in a genetic region
22 of interest at various sample sizes (e.g. downsamplings from gnomAD²). While the *input*

1 *simulation dataset* is required, the *target dataset* is not necessary if default or user defined
2 parameters are used.

3 RAREsim has three main steps: (1) simulating haplotypes, (2) estimating expected
4 number of variants, and (3) pruning rare variants to match expected. A flowchart summarizing
5 the RAREsim algorithm is in **Figure 1**.

6 *(1) Simulating an abundance of rare variants*

7 RAREsim uses HAPGEN2¹³ to simulate haplotypes for N_{sim} individuals. HAPGEN2 simulates
8 haplotypes by creating mosaics of input and already simulated haplotypes, using recombination
9 information so that regional LD is retained^{13,15}. When all sequencing bases, including
10 monomorphic, are included in the input haplotypes more de novo variants are simulated than
11 expected. By sampling from previously simulated haplotypes to create a new haplotype,
12 HAPGEN2 resamples de novo variants resulting in inflated numbers of rare variants.

13 *(2) Estimating expected number of variants per MAC bin*

14 The number of variants per MAC bin is estimated using two functions. The parameters for
15 these functions are estimated using target data (described below). Alternatively, user-defined or
16 default parameters can be used, eliminating the need for the user to provide and fit target data.
17 The default parameters were derived using the default target data (**Methods**).

18 *(2a) Number of Variants Function*

19 The total number of variants in a region depends on the sample size. RAREsim estimates the
20 expected number of variants per kilobase (Kb) for a sample size n using the *Number of Variants*
21 function. Specifically,

$$22 \quad f_{N_{variant}}(n) = \phi n^{\omega},$$

1 where $f_{Nvariant}(n)$ is the number of variants per Kb for n individuals. The parameters ϕ and ω
2 are estimated to modify the scale and shape of the function, respectively. When simulating N_{sim}
3 individuals, RAREsim calculates the total number of variants in the region by multiplying the
4 size of the region in Kb, S_{Kb} , by the expected number of variants per Kb, $f_{Nvariant}(n = N_{sim})$.

5 The target data used for the *Number of Variants* function provides the observed number
6 of variants per Kb, T_n , in the simulation region observed at sample size n . The parameters are
7 optimized by minimizing a least squares loss function summed over all observed sample sizes in
8 the target data:

$$9 \quad \min_{\phi, \omega} \left(\sum_n (T_n - \phi n^\omega)^2 \right).$$

10 Sequential quadratic programming (SQP) via the *slsqp* function in the *nloptr* R package¹⁷ is used
11 with constraints $0 < \omega < 1$ and $\phi > 0$ to minimize the loss function. The initial starting values
12 for the algorithm are $\omega = 0.45$ and ϕ such that the largest observed sample size, n_{max} , fits the
13 observed number of variants per Kb at n_{max} , $0.45n_{max}^\omega = T_{n_{max}}$. If the initial starting values do
14 not result in a sufficient fit (loss > 1,000), a range of starting values are evaluated: $\omega \in$
15 $\{0.15, 0.25, 0.35, 0.45, 0.55, 0.65\}$.

16 **(2b) Allele Frequency Spectrum Function (AFS Function)**

17 The AFS function, $f_{AFS}(z)$, estimates the proportion of variants with MAC = z . The
18 largest MAC, z_{max} , has MAF $\approx 1\%$ in the target dataset. For a target dataset with N_{target}
19 individuals, $z_{max} = floor(N_{target} * 2 * 0.01)$. Specifically,

$$20 \quad f_{AFS}(z) = b \times \frac{1}{(z + \beta)^\alpha},$$

21 with b such that $\sum_{z=1}^{z_{max}} f_{AFS}(z) = p_{RV}$

1 $for\ z \in \{1,2,3,\dots,z_{max}\}.$

2 The scale parameter b ensures the sum of all individual rare MAC proportions equals the total
 3 proportion of rare variants observed in the target data, p_{RV} . Parameters α and β determine the
 4 shape of the distribution.

5 Because individual $MAC = z$ may have no observed minor alleles in the target data,
 6 particularly for higher rare MACs (e.g. $MAC=10, 11, 12$), MAC bins are used to optimize
 7 parameters. MAC bins are mutually exclusive, exhaustive groups of rare MACs. Seven MAC
 8 bins are used here: singletons, doubletons, $MAC = 3-5$, $MAC = 6-10$, $MAC = 11-20$, $MAC = 21-$
 9 $MAF = 0.5\%$, $MAF = 0.5\% - 1\%$, denoted as $Bin_1, Bin_2, \dots, Bin_7$, respectively. The total
 10 number and thresholds to define the bins can be modified by the user. Within each bin j , the
 11 estimated proportion of variants, $\sum_{z \in Bin_j} \left(\frac{b}{(z+\beta)^\alpha}\right)$, is compared to the observed proportion of
 12 variants in the target data, $\sum_{z \in Bin_j} A_z$. A_z is the observed proportion of variants with $MAC = z$.
 13 Parameter estimates for α and β are found by minimizing the least squares loss over all bins
 14 using SQP¹⁷ constraining $\alpha > 0$,

15
$$\min_{\alpha, \beta} \left(\sum_{j=1}^B \left(\sum_{z \in Bin_j} \left(\frac{b}{(z + \beta)^\alpha} \right) - \sum_{z \in Bin_j} A_z \right)^2 \right).$$

16

17 ***(2c) Expected Number of Variants per Minor Allele Count Bin***

18 RAREsim uses the total number of variants within a genetic region for N_{sim} individuals,
 19 $f_{Nvariant}(N_{sim}) * S_{Kb}$, and the proportion of variants in Bin_j to obtain the expected number
 20 variants in each MAC bin, $E_{Bin_j}[v]$,

$$\mathbb{E}_{Bin_j}[v] = f_{Nvariant}(N_{Sim}) * S_{Kb} * \sum_{z \in Bin_j} f_{AFS}(z)$$

The expected total number of rare variants is calculated by summing across all rare MAC bins.

$$\mathbb{E}[v] = \sum_j \mathbb{E}_{Bin_j}[v]$$

The total number of simulated rare variants (M_{Sim}) is calculated by summing over all MAC bins.

$$M_{Sim} = \sum_j M_{Sim, Bin_j}$$

(3) *Pruning*

As described in (1), simulations using HAPGEN2 usually result in a larger total number of simulated rare variants than expected from step (2). Simulated variants are pruned by retuning all or a subset of alternate alleles to reference alleles. Within HAPGEN2 and similar to real haplotypes, rare alleles have a high probability of being on the same haplotype background. Pruning alternate alleles preserves the high likelihood that rare alleles are on the same haplotype. Variants are probabilistically pruned, creating variability over the simulation replicates in the number of variants per MAC bin.

RAREsim sequentially prunes variants from high to low MAC bins starting with the highest rare MAC bin that has at least 10% more simulated variants than expected (i.e. $M_{Sim, Bin_j} > 1.1 \mathbb{E}_{Bin_j}[v]$).

(1) **For MAC bins with more simulated variants than expected**, a simulated variant is pruned with probability $P(rem)_j$, where

$$P(rem)_j = 1 - \frac{\mathbb{E}_{Bin_j}[v]}{M_{Sim, Bin_j}}$$

1 For each variant within the bin, RAREsim randomly draws from a Uniform(0,1)
2 distribution. If the draw is within $[0, P(rem)_j]$, the variant is pruned. The location of
3 a pruned variant is stored to allow variants to be added back at lower MAC bins that
4 have fewer simulated variants than expected, as described in the next section.

5 (2) **For MAC bins with fewer simulated variants than expected** (i.e. $M_{Sim, Bin_j} <$
6 $\mathbb{E}_{Bin_j}[v]$), each of the K previously pruned variants from higher MAC bins are
7 added to Bin_j with probability,

$$8 \quad P(add)_j = \frac{\mathbb{E}_{Bin_j}[v] - M_{Sim, Bin_j}}{K}.$$

9 Random draws from Uniform(0,1) are used to determine which variants to add. The
10 variant is added if the draw is within $[0, P(add)_j]$. The MAC for each added variant
11 is determined with a random sample from all possible MACs within Bin_j . RAREsim
12 then randomly samples without replacement the necessary number of haplotypes
13 containing the alternate allele for the given variant. The allele for all other
14 haplotypes is returned to reference.

15

16 **Results**

17 *Evaluation of Number of Variants and Allele Frequency Spectrum Functions*

18 *AFS* and *Number of Variants* functions were fit using target data from gnomAD v2.1 for
19 four ancestry/sample-size groups (African, N = 8,128; East Asian, N = 9,197; Non-Finnish
20 European, N=56,885; South Asian, N = 15,308) (**Supplemental Table 1**). Data from
21 chromosome 19 was divided into 1 cM blocks; six blocks were merged with the preceding
22 adjacent block (**Methods**), resulting in 101 blocks for simulation (**Supplemental Table 2**).

1 The ancestry/sample-size specific fitted *Number of Variants* function closely matches the
2 observed values for all four ancestries (**Figure 2, Supplemental Figure 1**). Ninety percent of cM
3 blocks had a relative difference within 2.42% for the estimated vs. observed number of variants
4 per Kb. The average relative difference for all ancestries was -1.10% (90% CI =
5 $(-1.17\%, -1.02\%)$), with ancestry specific averages of -0.65% (African, 90% CI =
6 $(-0.83\%, -0.47\%)$), -1.09% (East Asian, 90% CI = $(-1.24\%, -0.94\%)$), -1.64% (Non-
7 Finnish European, 90% CI = $(-1.74\%, -1.55\%)$), and -1.00% (South Asian, 90% CI =
8 $(-1.13\%, -0.87\%)$) (**Supplemental Figure 2**). The negative mean relative differences indicate
9 a slight but systematic overestimation of the number of variants per Kb for most blocks. Within a
10 given target dataset, the *Number of Variants* function appears to slightly overestimate the
11 observations of larger sample sizes and underestimate the observations of smaller sample sizes
12 (**Supplemental Figure 1**).

13 Observed and estimated variation in the number of variants per Kb across cM blocks
14 increases with sample size (**Figure 2A**). However, even for the largest available target data
15 sample size (Non-Finnish European, $N=56,885$), the variability of the number of variants per Kb
16 remains low with 90% of the block-specific estimates within 35 variants of the median estimate.

17 The *AFS* function matched the observed data well with no apparent systematic bias. The
18 average absolute difference between the observed and estimated proportion of variants in each
19 MAC bin over all ancestries, blocks, and MAC bins was 0.53% (90% CI = $(0.51\%, 0.55\%)$)
20 (**Figure 2B**). Ninety percent of the estimated proportions were within 1.3% of that observed. The
21 maximum absolute difference in MAC bin proportion was 4.50% (observed in East Asian, MAC
22 3-5). Singleton counts matched particularly well, with a maximum absolute difference of 0.73%.

23 Despite different ancestries and widely different sample sizes (from $N = 8,128$ to $N =$

1 56,885 for African and Non-Finnish European respectively) the proportion of variants per MAC
2 bin were similar (**Supplemental Figure 3**). There is more variation between ancestry/sample-
3 size groups for the total proportion of rare variants (i.e. proportion of all variants with MAF
4 <1%) (**Supplemental Figure 4**). Regardless, within each ancestry, variation of the AFS between
5 cM blocks remains small. For instance, the MAC bin with the most variation (East Asian
6 singleton bin) has 90% of the blocks having estimated proportions within 6.3% of the median.

8 *Evaluation of Simulation Results*

9 One hundred replicates of each block were simulated using RAREsim and HAPGEN2
10 matching the gnomAD sample size for each ancestry group ($N_{African} = 8,128$, $N_{East\ Asian} =$
11 $9,197$, $N_{Non-Finnish\ European} = 56,885$, $N_{South\ Asian} = 15,308$) (**Figure 3**). RAREsim
12 produced a similar number of variants relative to gnomAD across all ancestry groups and MAC
13 bins indicating that the total number of variants and AFS are representative of real sequencing
14 data. Conversely, HAPGEN2¹³ with only polymorphic SNVs greatly underestimated the total
15 number of rare variants, especially very rare variants. HAPGEN2 simulations including all
16 sequencing bases produced many more rare variants than observed. These results are consistent
17 across all cM blocks and for the cumulative chromosome 19 coding region (**Supplemental**
18 **Figures 5-8, Supplemental Table 3-6**).

20 *Generalizability of Default Parameters*

21 Ancestry/sample-size specific default parameters for the *Number of Variants* and *AFS*
22 functions (**Table 1**) were estimated using the median observation over cM blocks for each
23 ancestry/sample-size group. RAREsim default parameters performed well, matching the

1 observed number of variants and AFS in a wide variety of situations including three GENCODE
2 regions on chromosomes 1, 6 and 9¹⁸, in non-coding regions within blocks, and in other
3 datasets/sample sizes - gnomAD v3 and UKBioBank¹⁹ (**Methods, Figure 4, Supplemental**
4 **Figures 9 – 14**). The simulated sample sizes evaluated were up to ~3x larger (gnomAD
5 v3African) and ~2x smaller (gnomAD v3 Non-Finnish European) than the sample sizes used to
6 derive the default parameters. The default parameters often performed similarly to the cM block
7 specific simulation parameters and always outperformed HAPGEN2.

8

9 *Stratified Simulation of Functional and Synonymous Variants*

10 As expected, we observed more functional SNVs than synonymous², with the largest
11 differences observed at $MAC \leq 5$ (**Supplemental Figure 15**). This resulted in substantially
12 different fitted *Number of Variants* functions for the two types of variants (**Supplemental**
13 **Figure 16**). Stratified simulation of functional and synonymous variants closely approximated
14 the number of variants observed in each MAC bin and suggests utility in separately simulating
15 different groups of variants (**Supplemental Figure 17**).

16

17 *Simulation of Large Sample Sizes*

18 As discussed previously (*Generalizability of Default Parameters*), RAREsim accurately
19 simulated 21,042 African samples to match gnomAD v3 using ancestry specific default
20 parameters derived from African gnomAD v2.1 (N=8,128). We currently assume that the AFS
21 does not change with sample size. Consistent shape of AFS was observed over the gnomAD v2.1
22 ancestry/sample size groups ($N_{African} = 8,128$, $N_{East\ Asian} = 9,197$, $N_{Non-Finnish\ European} =$
23 $56,885$, $N_{South\ Asian} = 15,308$). Further, ancestry specific fitted *Number of Variants* functions

1 extrapolated to larger sample sizes than observed were similar to the shape of the fitted *Number*
2 *of Variants* function for the total gnomAD v2.1 sample (N=125,748) (**Supplemental Figure 18**).
3 Therefore, we believe simulating sample sizes up to ~125,000 is likely reasonable.

4 5 *Computation Time*

6 The time to simulate one replicate using a desktop with 32GB RAM for a cM block on
7 chromosome 19 varied between 15 seconds for the smallest block (3,183 bp) and sample size
8 (N_{African}=8,128) and 12 hours 32 minutes 20 seconds for the largest block (81,235 bp) and sample
9 size (N_{Non-Finnish European}=56,5885). The median run time was 4 minutes 16 seconds
10 (**Supplemental Table 7**).

11 The amount of time to simulate haplotypes with RAREsim is dependent on the number of
12 samples being simulated and the size of the region. Simulating a region with ~19 Kb varied
13 between 1 minute 34 seconds for N=8,128 and 37 minutes 34 seconds for N=56,885. When
14 simulating N=15,308 individuals, RAREsim simulations took between 16 seconds for a region of
15 ~3 Kb and 11 minutes 31 seconds for a region of ~81 Kb. The rate limiting step in large
16 simulations was HAPGEN2. For the largest region and sample size, HAPGEN2 took over 11
17 hours to simulate when using a machine with 32 GB RAM. The same region was simulated in ~1
18 hour and 19 minutes using 192 GB RAM indicating that memory capacity was reached using the
19 original computing specs (**Methods**).

20 21 **Discussion**

22 Here we present RAREsim, a rare variant simulation algorithm. Unlike HAPGEN2,
23 which either severely under or over simulates the proportion of very rare variants, RAREsim

1 simulates the expected proportion of rare and very rare variants across a variety of genetic
2 regions, ancestries, and sample sizes. RAREsim produces simulations that match the expected
3 AFS, total number of variants, and haplotype structure while enabling variant annotation. To our
4 knowledge, no other existing simulation software is able to emulate real data in all of these areas.
5 We show that RAREsim's ancestry specific default parameters derived from the coding regions
6 of chromosome 19 generalize to other chromosomes, datasets, sample sizes¹⁹, and non-coding
7 regions, approximating the number of variants per MAC bin with remarkable accuracy. We offer
8 user flexibility by enabling use of RAREsim with default parameters, user defined parameters, or
9 parameters estimated to match user provided target data.

10 For typical uses of simulated genetic data (i.e. evaluating or comparing methods and
11 general power analysis) we recommend simulating with the default parameters. Default
12 parameters were shown to be robust across sample sizes, chromosomes, coding and intergenic
13 regions, and datasets. It is possible, although we believe unlikely, that the default parameters will
14 perform poorly when used in scenarios not evaluated here. If precise matching of a particular
15 empirical data characteristic such as functional variant type, genetic region, sample size, or
16 ancestry is important, we recommend re-estimating the simulation parameters using RAREsim
17 functions. Additionally, parameters are able to be specified without fitting target data. For
18 example, to simulate a specific ancestry, a user could make an educated decision on the total
19 number of variants based on the relationship to the ancestries evaluated here (e.g. total number of
20 variants between that of African and European).

21 RAREsim simulates haplotypes in the same form as HAPGEN2: hap/leg/sample files²⁰.
22 Haplotype files can be converted to vcf files using the SHAPEIT *convert* command²¹ or bcftools
23 *haplegendsample2vcf* command²². Genetic association with disease can be simulated from a

1 sample of generated haplotypes using an existing software such as PhenotypeSimulator²³.
2 Simulation of families or large pedigrees can be performed with a pedigree simulation software
3 such as ped-sim²⁴.

4 It has been shown that sample sizes in the tens to hundreds of thousands are needed to
5 have sufficient power to detect associations with rare variants²⁵. Due to the lack of very large
6 (>100,000), ancestry specific, publicly available target data at the time of publication, we could
7 not assess the accuracy of the RAREsim simulations for very large sample sizes. As genetic
8 sequencing resources continue to increase in size, RAREsim will be ideally suited to simulate
9 large sample sizes with estimation of new simulation parameters. For the *Number of Variants*
10 function, the extrapolated, ancestry specific *Number of Variants* functions were compared with
11 that of the full sample available in gnomAD v2.1. We believe that the *Number of Variants*
12 function is able to accurately simulate sample sizes up to what is observed in gnomAD v2.1
13 (~125,000). Alternatively, users can use population-genetics theory or other resources to make
14 an informed decision for the total number of variants expected for very large samples. One such
15 resource is the Capture-Recapture²⁶ algorithm, which can be used to estimate the number of
16 segregating sites given allele count data. While Capture-Recapture can extrapolate to larger
17 sample sizes, the software cannot be easily used for sample sizes that are smaller than the
18 observed target data. RAREsim does not currently modify the *AFS* function as sample size
19 increases. Consistent *AFS* were observed over the gnomAD sample sizes (N = 8,128 - 56,88).
20 However, we expect the *AFS* to deviate with very large sample sizes. A user can update the *AFS*
21 function parameters if desired, and research is ongoing to estimate the expected *AFS* for very
22 large sample sizes. We believe that RAREsim can currently accurately simulate sample sizes up
23 to ~125,000.

1 There are several limitations to RAREsim. First, RAREsim is only as good as the data on
2 which the simulations are based. Errors or inconsistencies in the target data or input simulation
3 haplotypes will be propagated through the simulations. Secondly, the default parameters were
4 developed and evaluated on autosomes. A user may fit sex chromosome target data or assume
5 parameters extend to sex chromosomes. Finally, RAREsim requires additional memory (e.g. >32
6 GB RAM) when simulating large regions and sample sizes. For efficient simulation of large
7 regions and sample sizes, highmem computing or breaking up the simulation region into smaller
8 portions and combining after simulation is needed. We are actively working on extensions for
9 these limitations.

10 One of the primary benefits of RAREsim is its ability to match real target data, either
11 provided by the user or as done here with gnomAD v2.1. Matching observed data allows
12 RAREsim to adapt as sequencing data evolves due to technological advances or improved
13 genetic resources from additional ancestral populations and increased sample sizes. For example,
14 RAREsim will be able to approximate TopMED²⁷ and ALFA, aggregated allele frequencies from
15 dbGaP²⁸ once these resources are released. RAREsim can also simulate unique characteristics of
16 a specific genetic region such as haploinsufficiency, contribution to a polygenic risk score, or
17 selection. The flexibility of RAREsim to emulate real data allows users to assess methods and
18 complete power analyses for in relevant and realistic genetic regions and samples.

19

20 **Acknowledgements**

21 The UK Biobank data was gathered using the UK Biobank Resource under Application
22 Number 42614. We would like to thank Achilleas Pitsillides for obtaining the UK Biobank allele
23 counts. We would also like to thank Robert Goedman for providing programming insight. We

1 thank Dr. Ferdinand Baer for support of this project. This work was supported by the National
2 Human Genome Research Institute (R35HG011293 and U01HG009080 to AEH and CGR;
3 U01HG009080-05S1 to CGR).

4

5 **Methods**

6 *Input Simulation Datasets*

7 For the input simulation dataset, we used haplotypes, legend files (an accompanying
8 variant list), and a recombination map from 1000 Genomes Phase 3 (hg19)³. The files were
9 modified to include information at each sequencing base. The recombination map was derived
10 from the combined sample of all ancestries (OMNI)³. Links to these resources are in

11 **Supplemental Table 8.**

12 For the simulation haplotypes, African, East Asian, Non-Finnish European, and South
13 Asian global ancestries were used with admixed African samples (African Caribbeans in
14 Barbados (ACB) and Americans of African ancestry in Southwest Utah (ASW)) excluded.
15 HAPGEN2 simulates biallelic SNVs. Hence, we omitted indels present in the 1000G haplotypes.
16 For multiallelic SNVs, the first alternate allele in the legend file with at least one observed
17 alternate allele was kept.

18

19 *Target Datasets*

20 Exome sequencing data from gnomAD v2.1² on chromosome 19 were used as target data to
21 estimate the simulation parameters. For the *Number of Variants* function, the observed number
22 of loss-of-function, synonymous, and missense variants by gene, ancestry group, and sample size
23 (e.g. 500, 1000, 2000, 5000, etc.) from Karczewski et al.² was used

1 (<https://storage.googleapis.com/gnomad-public/papers/2019-flagship->
2 [lof/v1.0/gnomad.v2.1.1.lof_metrics_downsamplings.txt.bgz](https://storage.googleapis.com/gnomad-public/papers/2019-flagship-lof/v1.0/gnomad.v2.1.1.lof_metrics_downsamplings.txt.bgz), **Supplemental Table 8**). The
3 number of variants for all three function classification groups were summed to obtain the total
4 number of variants observed per gene. The GENCODE v19 file¹⁸ (link in **Supplemental Table**
5 **8**) contains the genomic positions of the canonical transcript coding regions used by Karczewski
6 et al.² to define genes. The total number of variants within the simulation region of interest was
7 found by summing over all genes in the region. When a region contained overlapping genes, the
8 proportion of overlap was calculated and removed from one gene, as to not count variants twice.

9 For the AFS target data, allele counts for biallelic SNVs in the coding region of canonical
10 transcripts for four ancestry groups from gnomAD v2.1 were used (African, N = 8,128; East
11 Asian, N = 9,197; Non-Finnish European, N=56,885; South Asian, N = 15,308)²

12 ([https://storage.googleapis.com/gnomad-](https://storage.googleapis.com/gnomad-public/release/2.1.1/vcf/exomes/gnomad.exomes.r2.1.1.sites.vcf.bgz)
13 [public/release/2.1.1/vcf/exomes/gnomad.exomes.r2.1.1.sites.vcf.bgz](https://storage.googleapis.com/gnomad-public/release/2.1.1/vcf/exomes/gnomad.exomes.r2.1.1.sites.vcf.bgz), **Supplemental Table 8**).

14 We observed slight discrepancies for some regions in the total number of variants between the
15 gnomAD v2.1 data used for the AFS target data and the gnomAD downsamplings data used for
16 the Number of Variants target data (**Supplemental Figure 19**). Differences in the number of
17 variants per gene likely arise due to inconsistencies between the two datasets with respect to
18 classification of variant function, removal of variants in overlapping genes, as well as other
19 differences. The discrepancies did not substantially affect the simulation results, as shown when
20 the simulated haplotypes are compared to the gnomAD data (**Results**).

21

22 *Centimorgan Blocks*

1 Chromosome 19 was divided into 1 cM blocks for simulation. The cM blocks were
2 defined using the 1000 Genomes Project recombination map estimated from the combined set of
3 all ancestries³ (*Input Simulation Datasets*). Blocks were restricted to the coding region of the
4 canonical transcript for each gene. Genes that overlapped multiple blocks or were between cM
5 blocks (the recombination map did not contain information at all bp) were included with the
6 previous cM block. Of the 107 cM blocks, two blocks did not meet the requirement of containing
7 at least two genes (blocks 17 and 23). Additionally, there were four blocks (blocks 8, 50, 57, and
8 92) with fewer than 100 SNVs in at least one ancestry in the gnomAD target data. These six
9 blocks were merged with the preceding adjacent block, resulting in 101 blocks for simulation
10 (**Supplemental Table 2**). Blocks ranged from 3,183 bp to 81,253 bp (median = 19,029; Q1 =
11 11,037; Q3 = 27,204).

13 *Implementation of RAREsim*

14 RAREsim was implemented for a genetic region of interest using the computing
15 flowchart shown in **Supplemental Figure 20**. First, the input haplotype and legend files were
16 modified to include all bp, including monomorphic bases. Then, haplotypes were simulated for
17 each cM block using HAPGEN2¹³ with default parameters. The relative risk was set to 1.0, and
18 hence no disease loci were simulated. The random seed in HAPGEN is set by time. Therefore,
19 simulation replicates cannot be run in parallel across multiple cores starting at the same time. To
20 avoid this, we simulated replicates for the same simulation scenario on the same computing core
21 in series. Alternatively, the *pause* Bash command could be used when simulating haplotypes in
22 parallel. Then, the expected number of variants per MAC bin was calculated using the
23 *expected_variants* RAREsim function with either default simulation parameters or region-

1 specific simulation parameters estimated using *fit_afs* and *fit_nvariants* functions. Finally, the
 2 simulated haplotypes were pruned using the *prune_variants* function to identify pruning
 3 locations and a supplemental Bash script to efficiently prune the identified locations.

4

5 *Evaluation of Allele Frequency Spectrum and Number of Variants Functions*

6 In our application, parameters for the *Number of Variants* and *AFS* functions were
 7 estimated for each of the 101 blocks and four ancestry groups from gnomAD. To evaluate how
 8 well the *Number of Variants* function fit the target data, we calculated the relative difference for
 9 the observed target data at the sample size available in gnomAD, $T_{N_{gnomAD}}$, to the *Number of*
 10 *Variants* function estimate, $f_{N_{variant}}(n = N_{gnomAD})$. The relative difference was calculated as

$$11 \quad \frac{T_{N_{gnomAD}} - f_{N_{variant}}(n = N_{gnomAD})}{T_{N_{gnomAD}}} = \frac{T_{N_{gnomAD}} - \hat{\phi}(N_{gnomAD})^{\hat{\omega}}}{T_{N_{gnomAD}}}.$$

12 We evaluated fit of the *AFS* function with the difference between the estimated proportion of
 13 variants, $\hat{f}_{AFS}(z)$, and the observed proportion in gnomAD, A_z , for each MAC Bin_j ,

$$14 \quad \sum_{z \in Bin_j} \hat{f}_{AFS}(z) - \sum_{z \in Bin_j} A_z = \sum_{z \in Bin_j} \left(\frac{\hat{b}}{(z + \hat{\beta})^{\hat{\alpha}}} \right) - \sum_{z \in Bin_j} A_z.$$

15

16 *Default Function Parameters*

17 Ancestry specific default parameters were calculated using the median target data over all
 18 blocks (i.e. median number of variants per Kb at each sample size (*Number of Variants* function)
 19 and median proportion of variants in each MAC bin (*AFS* function)). For the *Number of Variants*
 20 function, the 5th and 95th percentile observations over the 101 blocks were used to estimate 5th
 21 and 95th percentile functions.

1

2 *Evaluation of Simulation Results*

3 One hundred replicates of each block were simulated for the gnomAD sample size of
4 each ancestry group ($N_{African} = 8,128$, $N_{East\ Asian} = 9,197$, $N_{Non-Finnish\ European} = 56,885$,
5 $N_{South\ Asian} = 15,308$). The matched sample size enabled a direct comparison between
6 gnomAD and simulated data as sample size greatly influences the number of variants expected in
7 MAC bins. We compared RAREsim to the default implementation of HAPGEN2 with only
8 polymorphic SNVs in the input simulation data and to HAPGEN2 using all bp, including
9 monomorphic bp. Each block was simulated and pruned independently. To evaluate simulations
10 from chromosome 19 as a whole, the variant counts for each MAC bin were summed over all cM
11 blocks.

12

13 *Generalizability of Default Parameters*

14 Generalizability of default parameters was assessed on different chromosomes, other
15 sample sizes, in an intergenic region, and simulating data to match another dataset. To evaluate
16 the performance of the default parameters for other chromosomes, we simulated GENCODE
17 regions on chromosomes 1, 6, and 9 that were chosen to be representative of the genome¹⁸
18 (**Supplemental Table 9**). As with the blocks on chromosome 19, the regions were restricted to
19 canonical coding exons. These blocks were each 500 Kb, but when restricted to the coding
20 region were 24,918; 12,519; and 17,051 bp on chromosomes 1, 6, and 9 respectively.

21 Whole genome sequencing data from gnomAD v3 was used to evaluate the performance
22 of default parameters for different sample sizes, and for intergenic regions. To evaluate default
23 parameters for different sample sizes, we simulated three blocks (5th, 50th, and 95th percentile

1 block for number of variants) for the African ancestry group (N=21,042 for v3 compared to
2 N=8,128 for v2.1) and Non-Finnish European ancestry group (N=32,299 for v3 compared to
3 N=56,885 for v2.1). To evaluate the utility of default parameters for intergenic regions, we
4 simulated intergenic regions within the three blocks limiting the original coding region size.

5 Finally, to evaluate performance of the default parameters in another dataset and sample
6 size, we simulated a Non-Finnish European sample to match the UK Biobank¹⁹. Due to an error
7 in the UK Biobank 50K release, the 95th percentile block contained missing data; thus, we used
8 the 5th, 50th, and 94th percentile blocks instead. We simulated 41,246 Non-Finnish European
9 individuals, which was the number of individuals in exome sequencing British sample after
10 removing ethnic outliers.

11

12 *Stratified Simulation of Functional and Synonymous Variants*

13 To demonstrate RAREsim's ability to simulate different types of variants, such as
14 variants in different functional classes, variants were stratified and simulated by functional and
15 synonymous status. The reference and alternate allele are required for variant annotation. For
16 polymorphic variants within gnomAD, the observed reference and alternate allele were used.
17 Monomorphic bp in gnomAD were annotated using all possible alternate alleles with the
18 *convert2annovar* function in ANNOVAR²⁹. To restrict to one alternate allele, each allele was
19 first annotated as a transition or transversion. Within the exome, Wang et al.³⁰ observed
20 transition to transversion ratios (Ti/Tv) between 2.79 and 2.84 across ancestries. Here, the
21 average Ti/Tv of 2.815 was used to calculate the probability (0.7379) of a transition for each
22 variant. For each monomorphic bp, we performed a random draw from a Uniform(0,1)
23 distribution. If the random draw was within [0,0.7379], the transition alternate allele was used.

1 Otherwise, the variant was annotated as a transversion and the alternate allele was assigned
2 randomly from the two possible alternate alleles.

3 Variants were annotated using Ensembl Variant Effect Predictor (VEP)³¹ release 100. For
4 variants with multiple annotations, the most severe consequence was chosen. Synonymous
5 variants were those annotated as synonymous. Matching gnomAD's annotation², functional
6 variants were those annotated as missense, frameshift, splice site disrupting, and stop gained.
7 Stratified simulation with RAREsim by variant class, including refitting target data and pruning,
8 was performed for the block with the median number of variants.

9

10 *Simulation of Large Sample Sizes*

11 We fit the *Number of Variants* function to the median cM block for the total gnomAD
12 v2.1 sample (N = 125,748) and compared the fitted, ancestry specific *Number of Variants*
13 functions extrapolated to large sample sizes.

14

15 *Computing Time*

16 Computing time was evaluated on an Ubuntu 18.04.2 LTS desktop with Intel® Core™
17 i7-6700 CPG at CPU 8 x 3.40Ghz. The desktop is 64-bit with 1.1 TB (disk) GNOME 3.28.2 and
18 32GB RAM. For each ancestry, the simulation time for the cM block with the minimum (3,183),
19 median (19,029), and maximum (81,235) bp was recorded. To re-evaluate the simulation of
20 haplotypes with HAPGEN2 using more memory, a Dual Intel Xeon E5-2670v2 (2.5 Ghz x 10
21 cores, each), 192GB PC3-12800R RAM (12x16GB sticks) was used.

22

23 **Data Availability**

1 All data used are publicly available with links found in **Supplemental Table 8**. The reference
2 haplotype and legend files with all monomorphic sequencing bases within the coding regions
3 included are available at https://github.com/meganmichelle/RAREsim_Example.

4

5 **Code Availability**

6 RAREsim is an open-source R package, and all code can be found at
7 <https://github.com/meganmichelle/RAREsim>. A small example simulation with the necessary
8 script is available at https://github.com/meganmichelle/RAREsim_Example. Code to complete
9 the majority of the analyses included here can also found at
10 https://github.com/meganmichelle/RAREsim_Example.

11

12 **Author Contributions**

13 M.N. and A.E.H. developed the algorithm and evaluation of the algorithm with insight from J.D.
14 and C.R.G. All analyses were performed by M.N., and the R package was developed by M.N.
15 M.N. and A.E.H. drafted the paper. A.E.H. supervised the project. All authors reviewed the final
16 manuscript.

17

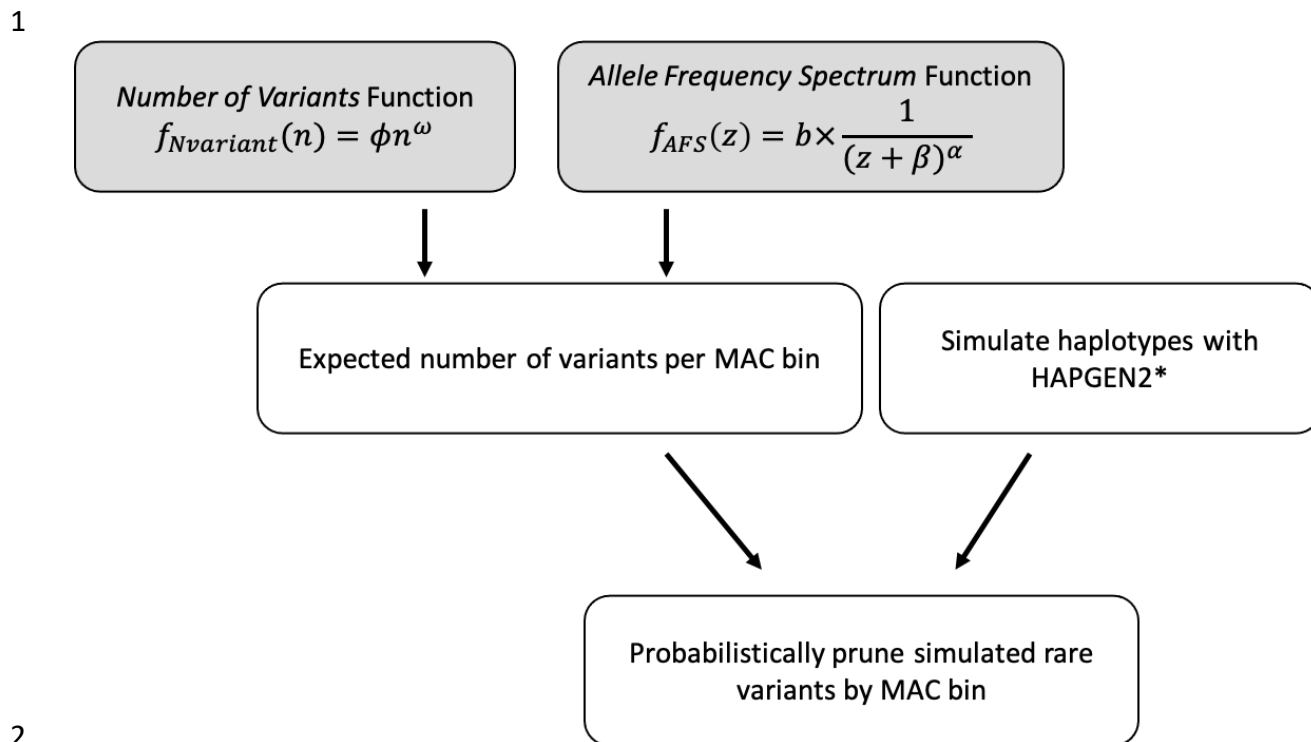
18 **Competing Interests**

19 The authors declare no competing interests.

References

- 1
- 2
- 3 1. Povysil, G. *et al.* Rare-variant collapsing analyses for complex traits: guidelines and applications. *Nat Rev Genet* **20**, 747-759 (2019).
- 4
- 5 2. Karczewski, K.J. *et al.* The mutational constraint spectrum quantified from variation in
- 6 141,456 humans. *Nature* **581**, 434-443 (2020).
- 7 3. Genomes Project, C. *et al.* A global reference for human genetic variation. *Nature* **526**,
- 8 68-74 (2015).
- 9 4. Consortium, U.K. *et al.* The UK10K project identifies rare variants in health and disease.
- 10 *Nature* **526**, 82-90 (2015).
- 11 5. Barbitoff, Y.A. *et al.* Whole-exome sequencing provides insights into monogenic disease
- 12 prevalence in Northwest Russia. *Mol Genet Genomic Med* **7**, e964 (2019).
- 13 6. Coventry, A. *et al.* Deep resequencing reveals excess rare recent variants consistent with
- 14 explosive population growth. *Nat Commun* **1**, 131 (2010).
- 15 7. Madsen, B.E. & Browning, S.R. A groupwise association test for rare mutations using a
- 16 weighted sum statistic. *PLoS Genet* **5**, e1000384 (2009).
- 17 8. Hendricks, A.E. *et al.* Rare Variant Analysis of Human and Rodent Obesity Genes in
- 18 Individuals with Severe Childhood Obesity. *Sci Rep* **7**, 4394 (2017).
- 19 9. The NHGRI Genome Sequencing Program, G. Functional Annotation of Variants - Online
- 20 Resource (FAVOR) Server. (2020).
- 21 10. Fisher, R.A. On the dominance ratio. *Proceedings of the Royal Society of Edinburgh* **42**,
- 22 321-331 (1923).
- 23 11. Wright, S. Evolution in Mendelian Populations. *Genetics* **16**, 97-159 (1931).
- 24 12. Kingman, J. On the genealogy of large populations. *Journal of Applied Probability* **19(A)**,
- 25 27-43 (1982).
- 26 13. Su, Z., Marchini, J. & Donnelly, P. HAPGEN2: simulation of multiple disease SNPs.
- 27 *Bioinformatics* **27**, 2304-5 (2011).
- 28 14. Li, N. & Stephens, M. Modeling linkage disequilibrium and identifying recombination
- 29 hotspots using single-nucleotide polymorphism data. *Genetics* **165**, 2213-33 (2003).
- 30 15. Hendricks, A.E., Dupuis, J., Gupta, M., Logue, M.W. & Lunetta, K.L. A comparison of gene
- 31 region simulation methods. *PLoS One* **7**, e40925 (2012).
- 32 16. Moutsianas, L. *et al.* The power of gene-based rare variant methods to detect disease-
- 33 associated variation and test hypotheses about complex disease. *PLoS Genet* **11**,
- 34 e1005165 (2015).
- 35 17. Johnson, S.G. The NLOpt nonlinear-optimization package, <http://ab-initio.mit.edu/nlopt>.
- 36 18. Frankish, A. *et al.* GENCODE reference annotation for the human and mouse genomes.
- 37 *Nucleic Acids Res* **47**, D766-D773 (2019).
- 38 19. Van Hout, C.V. *et al.* Whole exome sequencing and characterization of coding variation
- 39 in 49,960 individuals in the UK Biobank. *bioRxiv*, 572347 (2019).
- 40 20. Su, Z. HAPGEN version 2. Vol. 2020 (2011).
- 41 21. O'Connell, J. *et al.* Haplotype estimation for biobank-scale data sets. *Nat Genet* **48**, 817-
- 42 20 (2016).
- 43 22. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**,
- 44 2078-9 (2009).

- 1 23. Meyer, H.V. & Birney, E. PhenotypeSimulator: A comprehensive framework for
2 simulating multi-trait, multi-locus genotype to phenotype relationships. *Bioinformatics*
3 **34**, 2951-2956 (2018).
- 4 24. Caballero, M. *et al.* Crossover interference and sex-specific genetic maps shape identical
5 by descent sharing in close relatives. *PLoS Genet* **15**, e1007979 (2019).
- 6 25. Zuk, O. *et al.* Searching for missing heritability: designing rare variant association
7 studies. *Proc Natl Acad Sci U S A* **111**, E455-64 (2014).
- 8 26. Gravel, S., National Heart, L. & Blood Institute, G.O.E.S.P. Predicting discovery rates of
9 genomic features. *Genetics* **197**, 601-10 (2014).
- 10 27. Taliun, D. & al., e. Sequencing of 53,831 diverse genomes from the NHLBI TOPMed
11 Program. (2019).
- 12 28. Phan, L. *et al.* ALFA: Allele Frequency Aggregator. *National Center for Biotechnology*
13 *Information, U.S. National Library of Medicine* (2020).
- 14 29. Wang, K., Li, M. & Hakonarson, H. ANNOVAR: functional annotation of genetic variants
15 from high-throughput sequencing data. **38**, e164 (2010).
- 16 30. Wang, J., Raskin, L., Samuels, D.C., Shyr, Y. & Guo, Y. Genome measures used for quality
17 control are dependent on gene function and ancestry. *Bioinformatics* **31**, 318-23 (2015).
- 18 31. McLaren, W. *et al.* The Ensembl Variant Effect Predictor. *Genome Biol* **17**, 122 (2016).
- 19



*Haplotypes must include information at all sequencing bases to allow for an abundance of variants for subsequent pruning.

Figure 1. Flowchart of RAREsim Flowchart describing RAREsim simulation process. Simulation parameters can be estimated using target data, default parameters, or user defined parameters (gray).

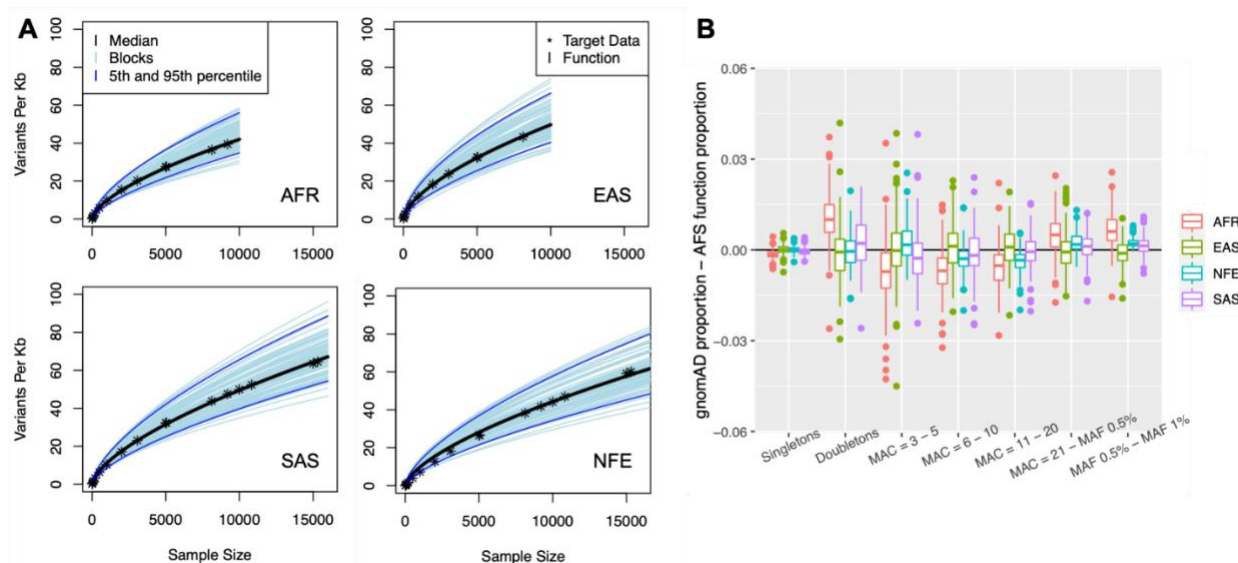
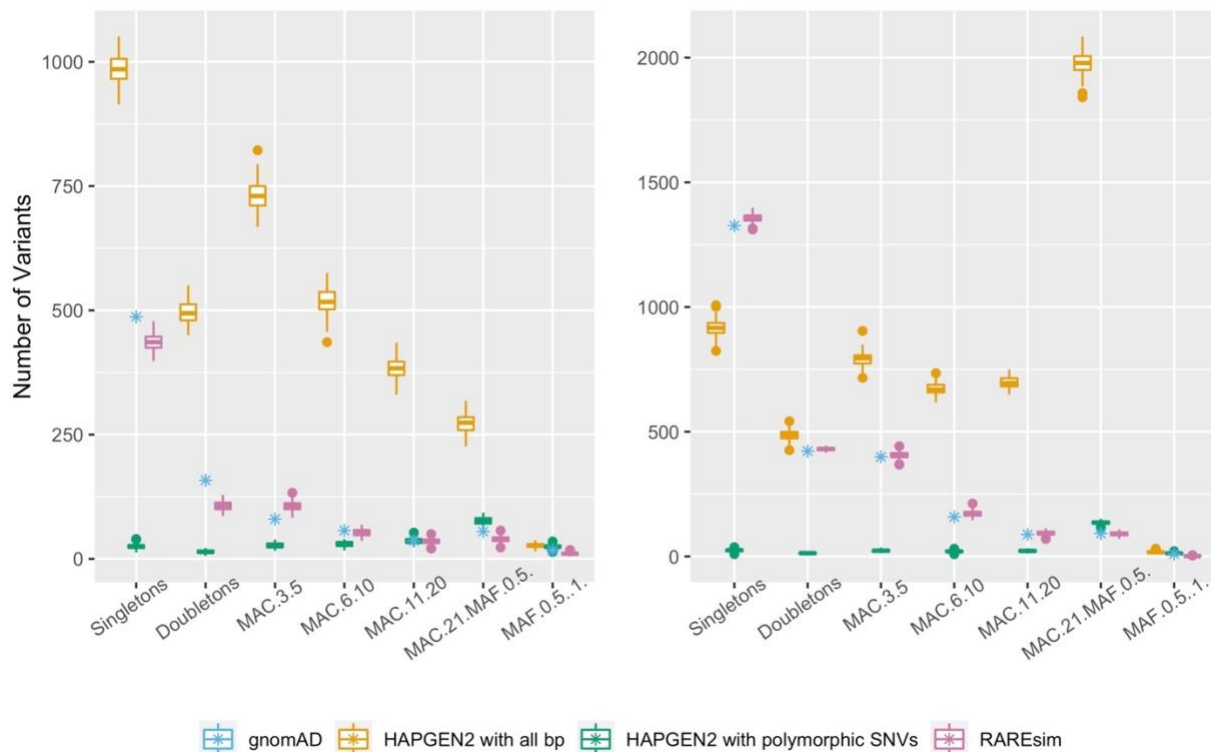
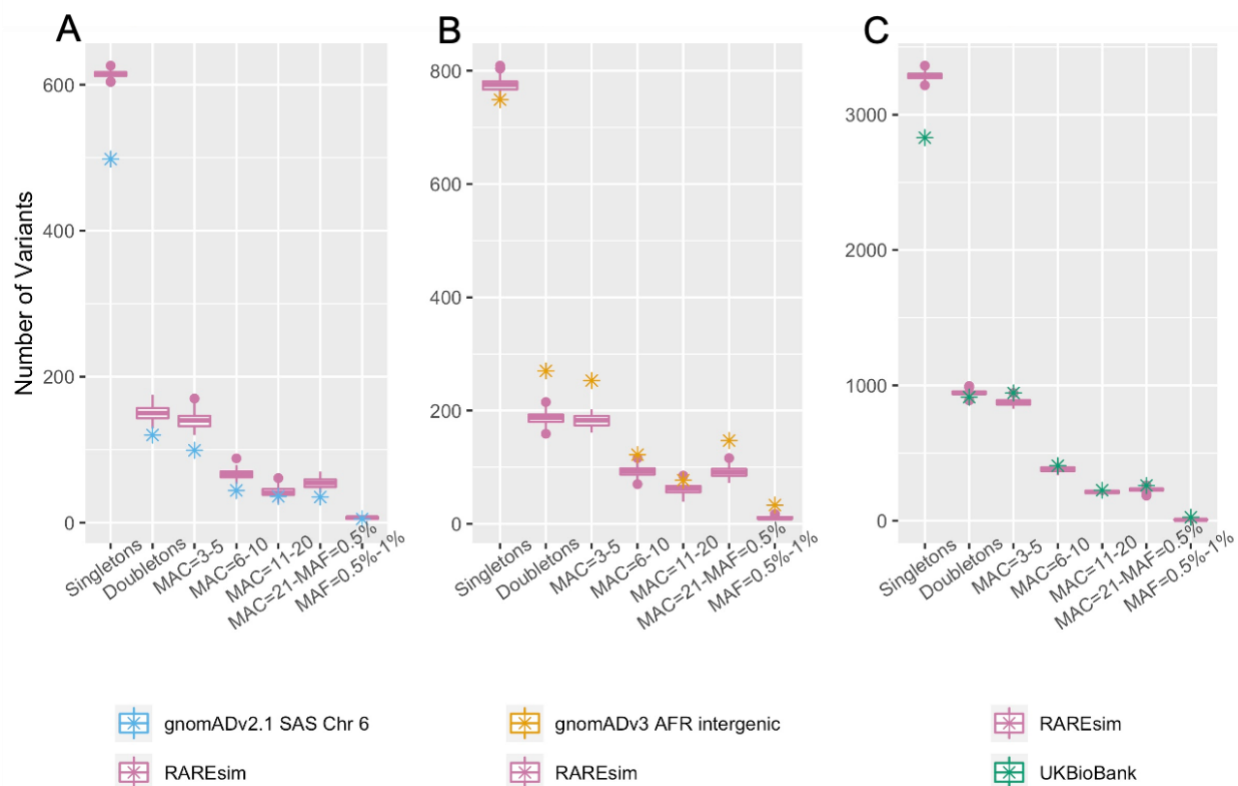


Figure 2. Evaluation of function fit **A) Number of Variants function:** Fitted *Number of Variants* functions for all cM blocks. The median block is shown in black and the 5th and 95th blocks are shown in dark blue. The observed target data (*) for the median block are close to the fitted function for all four ancestries. Sample sizes up to 15,000 are shown here. The full sample size for the Non-Finnish European sample is in Supplemental Figure 1. **B) AFS function:** Difference between the gnomAD target data and estimates from the *AFS* function for the proportion of variants in each MAC bin for all chromosome 19 blocks by ancestry and MAC bin. All absolute differences are within 0.05, indicating the *AFS* function fits the target data well.



1

Figure 3. Evaluation of RAREsim. The distribution and number of variants simulated using RAREsim (green), HAPGEN2 with only polymorphic SNVs (default, pink), and HAPGEN2 with all sequencing bases (blue) is compared to gnomAD (yellow) for the cM block with the median number of bp. Ancestry specific simulations are shown for African (N = 8,128; left) and Non-Finnish European (N = 56,885; right) matching the sample size observed in gnomAD v2.1. RAREsim emulates the expected number of variants within each MAC bin, while the other simulation methods either grossly underestimate (HAPGEN2 with polymorphic SNVs) or overestimate (HAPGEN2 with all sequencing bp) the number of variants.



1
2

Figure 2. Generalizability of Default Parameters Utility of RAREsim’s ancestry specific default parameters for different chromosomes (A), sample sizes (B & C), intergenic regions (B), and other target datasets (B & C). RAREsim simulations closely approximate the observed number of variants (y-axis) in each MAC bin (x-axis) in all scenarios. **A)** Simulations using South Asian default parameters for Chromosome 6 GENCODE region. **B)** Simulating a sample size of 21,042 using African ancestry default parameters (derived from N=8,128) for in an intergenic region from gnomAD v3. **C)** Simulating a sample size of 41,246 to match a British sample from the UK Biobank using Non-Finnish European default parameters (derived from N=56,885).

3

1 **Table 1: Default function parameters estimates.**

	Number of Variants $f_{Nvariant}(x) = \phi x^\omega$		Allele Frequency Spectrum $f_{AFS}(z) = b \times \frac{1}{(z + \beta)^\alpha}$		
	$\hat{\phi}$	$\hat{\omega}$	$\hat{\alpha}$	$\hat{\beta}$	\hat{b}
African	0.1576	0.6247	1.5883	-0.3083	0.2872
East Asian	0.1191	0.6369	1.6656	-0.2951	0.3137
Non-Finnish European	0.1073	0.6539	1.9470	0.1180	0.6676
South Asian	0.1249	0.6495	1.6977	-0.2273	0.3564

2