1    **Profiling and leveraging relatedness in a precision medicine cohort**

2    **of 92,455 exomes**

3

4    Jeffrey Staples,[1] Evan K. Maxwell,[1] Nehal Gosalia,[1] Claudia Gonzaga-Jauregui,[1] Christopher

5    Snyder,[2] Alicia Hawes,[1] John Penn,[1] Ricardo Ulloa,[1] Xiaodong Bai,[1] Alexander E. Lopez,[1]

6    Cristopher V. Van Hout,[1] Colm O'Dushlaine,[1] Tanya M. Teslovich,[1] Shane E. McCarthy,[1]

7    Suganthi Balasubramanian,[1] H. Lester Kirchner,[3] Joseph B. Leader,[3] Michael F. Murray,[3] David

8    H. Ledbetter,[3] Alan R. Shuldiner,[1] George Yancoupolos,[1] Frederick E. Dewey,[1] David J. Carey,[3]

9    John D. Overton,[1] Aris Baras,[1] Lukas Habegger,[1] and Jeffrey G. Reid[1,*]

10   **Affiliations:**

11   [1]Regeneron Genetics Center, Regeneron Pharmaceuticals, Tarrytown, NY, 10591, USA

12   [2]Rochester Institute of Technology, Rochester, NY 14623, USA

13   [3]Geisinger Health System, Danville, PA 17822, USA

14   [*]Corresponding author: email: jeffrey.reid@regeneron.com; twitter: @JGReid

15

16   **Abstract**

17       Large-scale human genetics studies are ascertaining increasing proportions of

18   populations as they continue growing in both number and scale. As a result, the amount

19   of cryptic relatedness within these study cohorts is growing rapidly and has significant

20   implications on downstream analyses. We demonstrate this growth empirically among

21   the first 92,455 exomes from the DiscovEHR cohort and, via a custom simulation

22   framework we developed called SimProgeny, show that these measures are in-line with

23   expectations given the underlying population and ascertainment approach. For example,

1    we identified ~66,000 close (first- and second-degree) relationships within DiscovEHR

2    involving 55.6% of study participants. Our simulation results project that >70% of the

3    cohort will be involved in these close relationships as DiscovEHR scales to 250,000

4    recruited individuals. We reconstructed 12,574 pedigrees using these relationships

5    (including 2,192 nuclear families) and leveraged them for multiple applications. The

6    pedigrees substantially improved the phasing accuracy of 20,947 rare, deleterious

7    compound heterozygous mutations. Reconstructed nuclear families were critical for

8    identifying 3,415 *de novo* mutations in ~1,783 genes. Finally, we demonstrate the

9    segregation of known and suspected disease-causing mutations through reconstructed

10   pedigrees, including a tandem duplication in *LDLR* causing familial hypercholesterolemia.

11   In summary, this work highlights the prevalence of cryptic relatedness expected among

12   large healthcare population genomic studies and demonstrates several analyses that are

13   uniquely enabled by large amounts of cryptic relatedness.

14

19

20

21

22

23

# 1   **Introduction/Background**

2      The number and scale of large human sequencing projects is rapidly growing,

3   including DiscovEHR[1], UK Biobank[2], the US government's All of Us (part of the

4   Precision Medicine Initiative)[3], TOPMed (Web Resources), ExAC/gnomAD[4], and many

5   others. Many of these studies are collecting samples from integrated healthcare

6   populations that have accompanying phenotype-rich electronic health records (EHRs)

7   with a goal of combining the EHRs and genomic sequence data to catalyze translational

8   discoveries and precision medicine[1]. These large-scale healthcare population-based

9   genomic (HPG) studies are recruiting participants through healthcare systems where

10   volunteers donate DNA and provide medically relevant metrics recorded in their EHRs.

11   A major difference between the new HPG study design and traditional population-based

12   studies is ascertainment, both in how participants are recruited and in the proportion of

13   the population within in a geographical area that participates (Figure 1A). Traditionally,

14   the high-expense of large-scale genetic studies and the limited resources of individual

15   investigators has generated study populations exhibiting shallow ascertainment of

16   individuals from a variety of geographical areas. To improve statistical power, samples

17   from many different collection centers are combined into a larger cohort, and these

18   cohorts are often merged into a much larger consortium consisting of tens to hundreds of

19   thousands of individuals. While the total number of individuals sampled is often high,

20   these studies typically only sample a relatively small portion of individuals in any given

21   geographic area. In contrast, planned and on-going HPG studies are sampling tens to

22   hundreds of thousands of participants from individual healthcare systems[1].

1    The difference in these two ascertainment approaches results in different patterns of

2    genetic relatedness among individuals in these cohorts. Relatedness is a continuum that

3    manifests itself within a cohort in a variety of ways depending on the population and how

4    individuals are sampled from it. Because traditional population-based studies have

5    generally collected samples from multiple geographical areas, they most commonly

6    exhibit the broadest "class" of relatedness: *population structure*. Population structure

7    (often referred to as "substructure" or "stratification") within a genetic study results from

8    it containing different ancestral groups or "genetic demes" in which allele frequencies are

9    more similar within genetic demes than between demes. Genetic demes arise due to

10   more recent genetic isolation, drift, and migration patterns. Classic examples of

11   homogenous genetic demes include founder populations such as the Ashkenazi Jewish[5,6]

12   and Old Order Amish[7], but they also occur at many levels of genetic isolation such as

13   continental, sub-continental[6,8,9], and even within the same urban community[5,10]. When

14   genetic demes begin to mix, they create genetically admixed populations. Studies that

15   contain an admixed population or more than one genetic deme are likely to contain

16   population structure. Ascertainment of individuals within genetic demes can generate

17   *distant cryptic relatedness*[5,6], the second "class" of relatedness[11], defined here as third- to

18   ninth-degree relatives. These distant relatives are unlikely to be identifiable from the

19   EHR but are important because usually one or more large segments of their genomes are

20   identical-by-descent (IBD), depending on their degree of relatedness and the

21   recombination and segregation of alleles[12]. Distant cryptic relatedness is usually limited

22   in study cohorts built from small samplings of large populations, but the level of cryptic

23   relatedness increases substantially as the effective population size decreases and the

1    sample size increases. Finally, unless designed to collect families, traditional population-

2    based studies typically have very little *family structure:* the third "class" of relatedness,

3    consisting of first- and second-degree relationships[2,4,13-15] (Figure 1B, C).

4        In contrast, the HPG study design enriches for family structure in several ways. First,

5    HPG studies heavily sample from specific healthcare system regions, and the number of

6    pairs of related individuals ascertained increases combinatorially as more individuals are

7    sampled from a single region (Figure 1A). Second, families who live in the same

8    geographic area likely receive medical care from the same doctors at the same healthcare

9    system due to referrals, shared insurance coverage, and convenience. Third, families who

10   have visited a healthcare system for many years with multiple encounters will have

11   extensive medical records, making them more likely to be included in a study compared

12   to transient residents with brief medical records and fewer encounters. Both family

13   structure and distant cryptic relatedness are more pronounced in populations with low

14   migration rates[5]. Conversely, confounding population substructure may be less of a factor

15   in HPG studies if the sampled healthcare system's population is a single homogenous

16   genetic deme[1]. As a result, we expect to see an enrichment of family structure in HPG

17   studies compared to random ascertainment of a population. In this article, we focus on

18   family structure and its prevalence in a HPG study using both simulated and real data.

19       The increase in family structure within HPG studies has significant implications when

20   choosing and executing downstream analyses and must be considered thoughtfully[16-22].

21   Some tools assume all individuals are unrelated (e.g. PCA), some effectively handle

22   estimates of pairwise-wise relationships (e.g. linear mixed models), and others can

23   directly leverage pedigree structures (e.g. linkage and TDT analyses). The following is a

1    description of some common analysis tools and their use cases based on the varying

2    levels of family structure within large population-based datasets (Figure 1D).

3        Removal of family structure (i.e. selectively excluding samples to eliminate

4    relationships) is a viable option if a dataset has few closely related samples[4,13-15] and if

5    the size of the unrelated subset is acceptable for the statistical analysis being performed.

6    For example, principal components (PCs) can be computed on an unrelated subset of the

7    data, and then all samples can be projected onto these PCs[1]. A number of methods exist

8    to compute the maximally sized unrelated set of individuals[23,24]. However, this strategy

9    reduces the sample size and power while discarding potentially valuable relationship

10    information. In practice, the degree of information loss is unacceptable for many analyses

11    if the dataset has even a moderate level of family structure.

12        Several statistical methods have been developed that explicitly model estimates of

13    pairwise relationships. For example, mixed models provide better power for genome-

14    wide association studies and reduce type-one error compared to methods that do not

15    model the confounding relatedness[19,25-27], but mixed models become computationally

16    expensive when applied at scale. Pairwise relationships can also be used in a pedigree-

17    free QTL linkage analysis[28]. Additional software packages exist that model population

18    structure and family structure for pairwise relationship estimation (PCrelate)[29] and

19    principal component analysis (PC-AiR)[30].
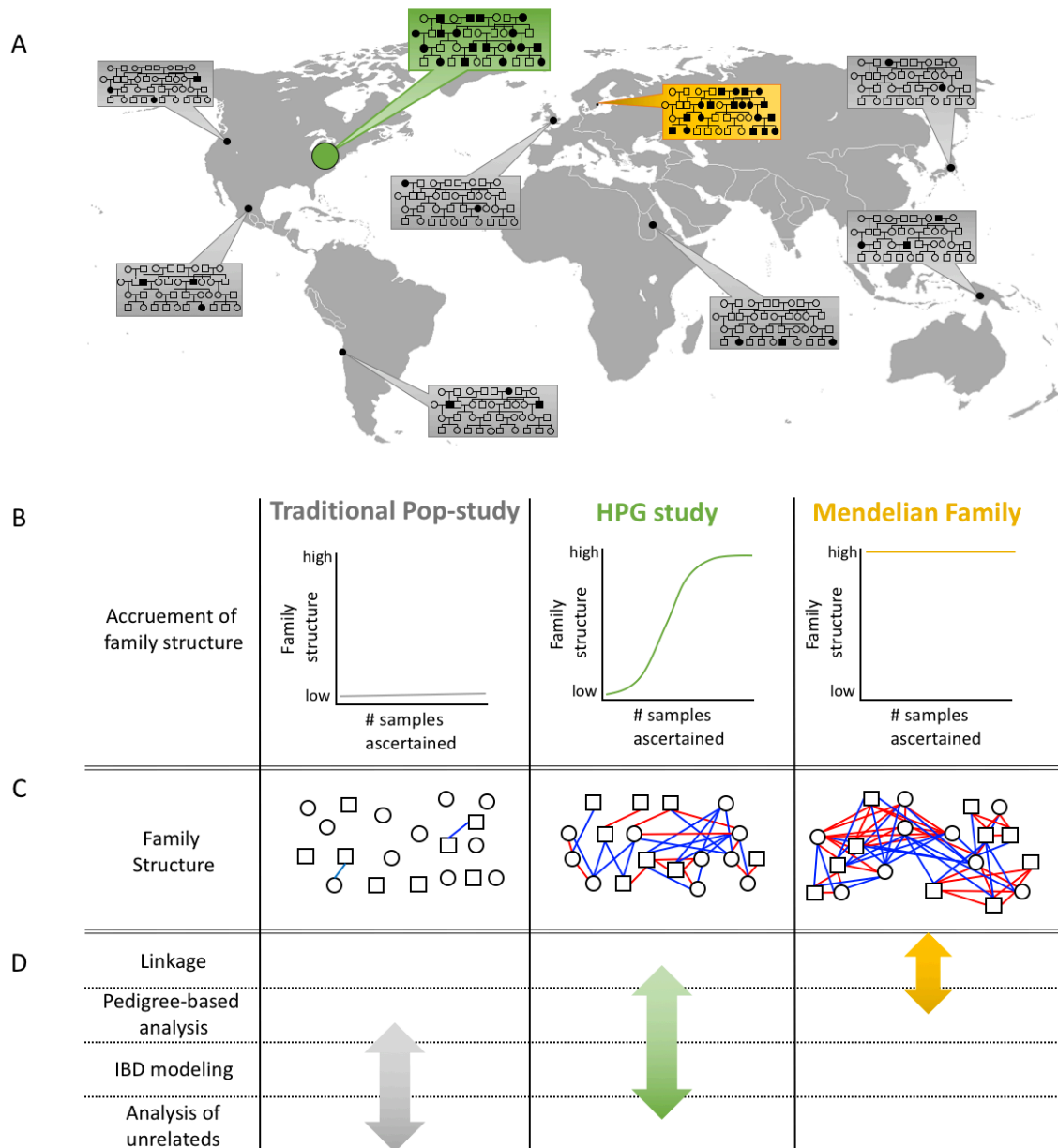
20        Extensive family structure can be used to reconstruct informative pedigree structures

21    directly from the genetic data with tools like PRIMUS[31] and CLAPPER[32], opening the

22    door to a number of pedigree-based methods and analyses that are particularly useful for

23    investigating the correlation between rare variation and disease[33]. Methods that use

1    pedigree structures to aid in identifying the genetic cause of a given phenotype typically

2    involve innovative variations on association mapping, linkage analysis, or both, including

3    MORGAN[34], pVAAST[17], FBAT (Web Resources), QTDT (Web Resources),

4    ROADTRIPS[35], rareIBD[36], and RV-GDT[37]. The appropriate method to use depends on

5    the phenotype, mode of inheritance, ancestral background, pedigree structure/size,

6    number of pedigrees, and size of the unrelated dataset[38]. Genetically reconstructed

7    pedigrees and estimated relationships can be used in a number of other ways beyond

8    association analyses, including pedigree-aware imputation, pedigree-aware phasing,

9    Mendelian error checking, compound heterozygous knockout detection, *de novo* mutation

10   calling, empirical validation of variant calling methods, and rare variant family-based

11   segregation analysis.

12      We demonstrate the value of identifying family structure in a large clinical cohort as

13   part of the DiscovEHR study. This cohort of 92,455 exomes originated from a

14   collaborative, ongoing study by the Regeneron Genetics Center (RGC) and the Geisinger

15   Health System (GHS) initiated in 2014[1]. DiscovEHR is a dense sample of patient-

16   participants from a single healthcare system that serves a largely rural population in

17   central Pennsylvania with low migration rates. We identify a tremendous amount of

18   family structure within the DiscovEHR cohort, and our simulations project that 70%-80%

19   of the individuals in our sequenced cohort will have a first- or second-degree relative as

20   we continue sequencing up to 250K individuals. This has significant implications on

21   downstream analyses, but also affords us the opportunity to leverage the rich family

22   structure through pedigree reconstruction, phasing compound heterozygous mutation

23   (CHM), and detecting *de novo* mutations (DNM).

1

2



3

*Figure 1. Ascertaining a high proportion of the population in a geographical area increases*

*family structure and impacts statistical analysis approaches that should be used. (A) Traditional*

*population-based studies (gray boxes) typically sample a small portion of individuals from*

*several populations. HPG studies (green box) more densely sample individuals from one or more*

1    *populations. Family-based studies (yellow box) heavily sample within extended families, but do*

2    *not sample nearly as many individuals as the other two study designs. (B) The three study designs*

3    *result in very different levels of individuals with one or more close relatives in the dataset. (C)*

4    *These three ascertainment approaches result in very different amounts of family structure. Red*

5    *and blue lines indicate first- and second-degree pairwise relationships, respectively. HPG studies*

6    *are expected to contain a level of family structure between the other two designs. (D) For this*

7    *study, statistical analysis approaches were binned into four categories based on the level of*

8    *family structure required to effectively use the approach (first column): "Linkage" refers to*

9    *traditional linkage analyses using one or more informative pedigrees; "Pedigree-based analysis"*

10    *refers to statistical methods beyond linkage that use pedigree structures within a larger cohort*

11    *that includes unrelated individuals; "IBD modeling" refers to analysis that model the pairwise*

12    *relationships between individuals without using the entire pedigree structure; "Analysis of*

13    *Unrelateds" refers to analyses that assume all individuals in the cohort are unrelated. The*

14    *amount of family structure impacts the approaches that can be used, and the arrows indicate the*

15    *analysis ranges for which the three study designs are best suited.*

16

17    **Subjects and Methods**

18    **Patients and samples**

19    We sequenced the exomes of 93,368 de-identified patient-participants from the

20    Geisinger Health System (GHS). Participants were consented into the MyCode®

21    Community Health Initiative[38] and contributed DNA samples for genomic analysis as

22    part of the Regeneron-GHS DiscovEHR collaboration[1]. Each patient has their exome

23    linked to a corresponding de-identified electronic health record (EHR). A more detailed

24    description of the first 50,726 sequenced individuals has been previously published[1,39].

1    The DiscovEHR Study did not specifically target families to participate in the

2    study, but implicitly enriched for adults with chronic health problems who interact

3    frequently with the healthcare system (and may be related to each other), as well as

4    participants from the Coronary Catheterization Laboratory and the Bariatric Service from

5    GHS.

6    **Sample preparation, sequencing, variant calling, and sample QC**

7    Sample prep and sequencing for the first ~61K samples have been previously

8    described[1] and this set of samples are referred to in this manuscript as the "VCRome set".

9    The remaining set of ~31K samples were prepared in the same process, except in place of

10   the NimbleGen probed capture we used a slightly modified version of IDT's xGen probes;

11   supplemental probes were added to capture regions of the genome well-covered by the

12   NimbleGen VCRome capture reagent but poorly covered by the standard xGen

13   probes. Captured fragments were bound to streptavidin-conjugated beads and non-

14   specific DNA fragments were removed by a series of stringent washes according to the

15   manufacturer's recommended protocol (IDT). We refer to this second set of samples as

16   the "xGen set". Variant calls were produced using the Genome Analysis Toolkit (GATK;

17   Web Resources). GATK was used to conduct local realignment of the aligned, duplicate-

18   marked reads of each sample around putative indels. GATK's HaplotypeCaller was then

19   used to process the INDEL-realigned, duplicate-marked reads to identify all exonic

20   positions at which a sample varied from the genome reference in the genomic VCF

21   format (GVCF). Genotyping was accomplished using GATK's GenotypeGVCFs on each

22   sample and a training set of 50 randomly selected samples outputting a single- sample

23   VCF file identifying both SNVs and indels as compared to the reference. The single

1    sample VCF files were used to create a pseudo-sample that contained all variable sites

2    from the single sample VCF files in both sets. Independent pVCF files were created for

3    the VCRome set by joint calling 200 single-sample gVCF files with the pseudo-sample to

4    force a call or no-call for each sample at all variable sites across the two capture sets. All

5    200-sample pVCF files were combined to create the VCRome pVCF file. This process

6    was repeated to create the xGen pVCF file. The VCRome and xGen pVCF files were

7    then combined to create the union pVCF. We aligned sequence reads to GRCh38 and

8    annotated variants using Ensembl 85 gene definitions. We restricted the gene definitions

9    to 54,214 transcripts that are protein-coding with an annotated start and stop,

10   corresponding to 19,467 genes. After the previously described sample QC process,

11   92,455 exomes remained for analysis.

12

13   **Principal components and ancestry estimation**

14           We used PLINKv1.9[24] to merge the union datasets with HapMap3[40] and kept only

15   SNPs that were in both datasets based on rsID. We also applied the following PLINK

16   filters: --maf 0.1 --geno 0.05 --snps-only --hwe 0.00001 to obtain a set of high-quality

17   common variants. We calculated principal component (PC) analysis for the HapMap3

18   samples and then projected each sample in our dataset onto those PCs using PLINK. We

19   used the PCs for the HapMap3 samples to train a kernel density estimator (KDE) for each

20   of the five ancestral super classes: African (AFR), admixed American (AMR), east Asian

21   (EAS), European (EUR), and south Asian (SAS). We used the KDEs to calculate the

22   likelihood that each sample belongs to each of the super classes. For each sample, we

23   assigned the ancestral superclass based on the likelihoods. If a sample has two ancestral

1    groups with a likelihood > 0.3, then we assigned AFR over EUR, AMR over EUR, AMR

2    over EAS, SAS over EUR, AMR over AFR; otherwise "UNKNOWN" (this was done to

3    provide stringent estimates of the EUR and EAS populations and inclusive estimates for

4    the more admixed populations in our dataset). If zero or more than two ancestral groups

5    had a high enough likelihood, then the sample was assigned "UNKNOWN" for ancestry.

6    Samples with unknown ancestry were excluded from the ancestry-based identity-by-

7    descent (IBD) calculations.

8

9    **IBD estimation**

10   Genome-wide identity-by-descent (IBD) estimates are a metric to quantify the

11   level of relatedness between pairs of individuals[28]. We applied the same Hardy-Weinberg

12   equilibrium, minor allele frequency, and variant level missingness that we applied during

13   the PCA analysis. Next, we used a two-pronged approach to obtain accurate IBD

14   estimates from the DiscovEHR cohort exomes. First, we calculated IBD estimates among

15   individuals within the same ancestral superclass (e.g. AMR, AFR, EAS, EUR, and SAS)

16   as determined from our ancestry analysis. We used the following PLINK flags to obtain

17   IBD estimates out to second-degree relationships: --genome --min 0.1875. This allows

18   for more accurate relationship estimates because all samples share similar ancestral

19   alleles; however, this approach is unable to predict relationships between individuals with

20   different ancestral backgrounds, e.g. a child of a European father and Asian mother.

21   Second, in order to catch the first-degree relationships between individuals with

22   different ancestries, we calculated IBD estimates among all individuals using the --min

23   0.3 PLINK option. We then grouped individuals into first-degree family networks where

1   network nodes are individuals and edges are first-degree relationships. We ran each first-

2   degree family network through the prePRIMUS pipeline[31], which matches the ancestries

3   of the samples to appropriate ancestral minor allele frequencies to improve IBD

4   estimation. This process accurately estimates first- and second-degree relationships

5   among individuals within each family network (minimum PI_HAT of 0.15).

6          Finally, we combined the IBD estimates from the two previously described

7   approaches by adding in any missing relationships from family network derived IBD

8   estimates to the ancestry-based IBD estimates. This approach resulted in accurate IBD

9   estimates out to second-degree relationships among all samples of similar ancestry and

10  first-degree relationships among all samples.

11         IBD proportions for third-degree relatives are challenging to accurately estimate

12  from large exome sequencing dataset with diverse ancestral backgrounds because the

13  analysis often results in an excess number of predicted 3$^{rd}$ degree relationships due to

14  artificially inflated IBD estimates. We used a --min 0.09875 cutoff during the ancestry

15  specific IBD analysis to get a sense of how many third-degree relationships we may have

16  in the DiscovEHR cohort, but these were not used in any of the phasing or pedigree-

17  based analyses. Rather for the relationships-based analyses reported in this paper, we

18  only used high-confidence third-degree relationships we identified within first- and

19  second-degree family networks.

20         During QC prior to creating the final set of 92,455 individuals, we removed all

21  identical pairs of samples (PI_HAT > 0.9) unless GHS was able to find evidence through

22  a chart review that the two corresponding individuals appear to be different people, share

1    the same birthdate, and had one or more additional pieces of information that they were

2    related (e.g. same last name, shared parent, same address, or listed the other as a relative).

3

4    **Pedigree reconstruction**

5        We reconstructed all first-degree family networks identified within the

6    DiscovEHR cohort with PRIMUSv1.9.0[31]. The combined IBD estimates were provided

7    to PRIMUS along with the genetically derived sex and EHR reported age. We specified a

8    relatedness cutoff of PI_HAT > 0.375 to limit the reconstruction to first-degree family

9    networks, and a minimum cutoff of 0.1875 to define second-degree networks.

10

11    **Allele-frequency-based phasing**

12        We phased all bi-allelic variants from the VCRome and xGen exome datasets

13    separately using EAGLEv2.3[41]. In order to parallelize our analysis, we divided the

14    genome into overlapping segments of ~40K variants with a minimum overlap of 500

15    variants and 250K base-pairs. Since our goal was to phase putative compound

16    heterozygous mutations within genes, we took care to have the segment break points

17    occur in intergenic regions.

18        We used the UCSC LiftOver program to lift-over EAGLE's provided

19    genetic_map_hg19.txt.gz file from hg19 to GRCh38 and removed all variants that

20    switched chromosomes or changed relative order within a chromosome resulting in the

21    cM position to not be increasing when sorting on increasing chromosome position. In

22    most cases, this QC step removed inversions around centromeres. We also removed all

23    SNPs that mapped to an alternate chromosome. In total, only 2,783 of the 3.3 million

1    SNPs were removed from the genetic map file. We provided the data for each segment to

2    EAGLE as PLINK formatted files and ran it on DNAnexus with the following EAGLE

3    command line parameters:

4    --geneticMapFile=genetic_map_hg19_withX.txt.GRCh38_liftover.txt.gz

5    --maxMissingPerIndiv 1

6    --genoErrProb 0.01

7    --numThreads=16

8

9    **Compound heterozygous calling**

10    Our goal was to obtain high confidence compound heterozygous mutation (CHM)

11    calls of putative loss-of-function (pLoF) variants to identify humans with both copies of

12    genes potentially knocked out or disrupted. We classify variants as pLoFs if they result in

13    a frameshift, stop codon gain, stop codon loss, start codon gain, start codon loss, or

14    splicing acceptor or donor altering variant. We created a second, expanded set of

15    potentially harmful variants that included the pLOFs as well as likely disruptive missense

16    variants, which are variants predicted to be deleterious by all five of the following

17    methods: SIFT[41] (damaging), PolyPhen2 HDIV[42] (damaging and possibly damaging),

18    PolyPhen2 HVAR (damaging and possibly damaging), LRT[43] (deleterious), and

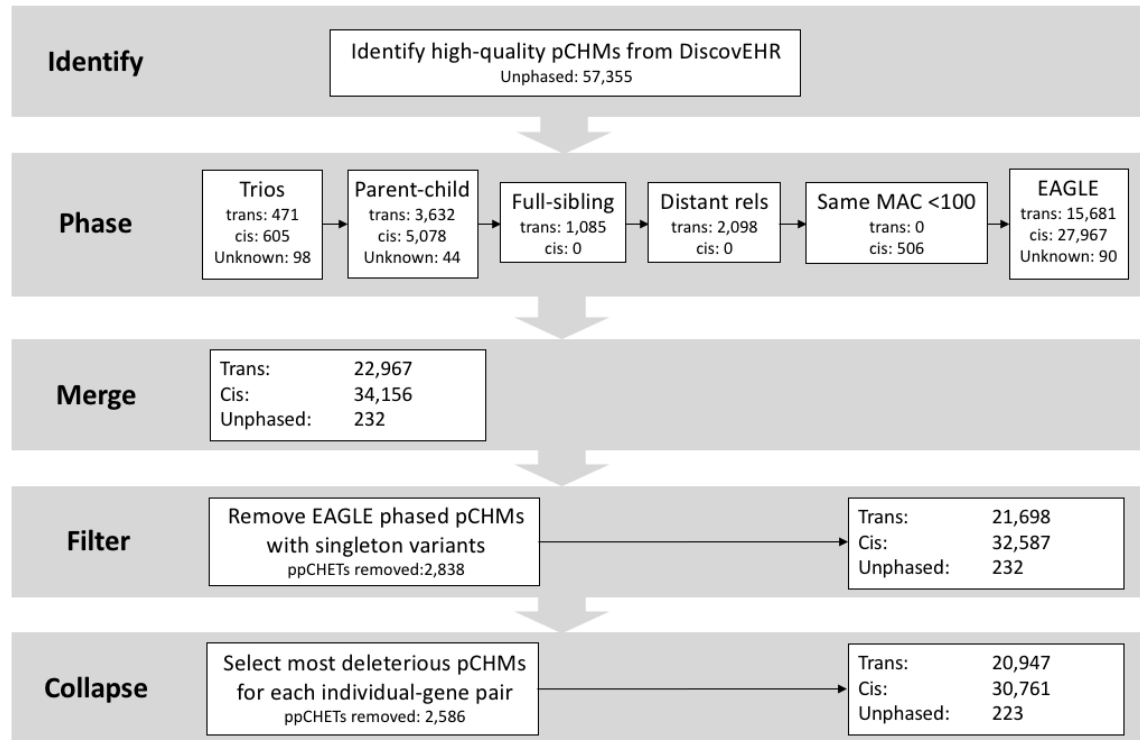19    MutationTaster[44] (disease causing automatic and disease causing).

20    We identified rare (alternate allele frequency < 1%) potential compound

21    heterozygous mutations (pCHMs) by testing all possible combinations of heterozygous

22    pLoFs and/or deleterious missense variants within a gene of the same person. We

23    excluded all variants that were out of Hardy-Weinberg equilibrium (p-value < $10^{-15}$

1   calculated with PLINKv1.9[24]), that exceeded 10% missingness within the individuals

2   capture specific dataset (i.e. VCRome or xGen sets), or that had another variant within 10

3   base-pairs in the same individual. We also excluded SNPs with quality by depth (QD) < 3,

4   alternate allele balance (AB) < 15%, and read depth < 7, and we excluded indels with QD

5   < 5, AB < 20%, and read depth < 10. After filtering, we had 57,355 high-quality pCHMs

6   distributed among 36,739 individuals that could knockout or disrupt the function of both

7   copies of a person's gene if the pCHM variants are phased in *trans*.

8       The next step was to phase the pCHMs. We used a combination of population

9   allele-frequency-based phasing with EAGLE and pedigree/relationship-based phasing to

10  determine if the pCHMs were in *cis* or *trans*. Figure 2 diagrams the pCHM phasing

11  workflow we employed to obtain the most accurate phasing for each pCHM. Trios and

12  relationships with individuals in both the VCRome and xGen datasets were used only if

13  both variants in the pCHM were on both the VCRome and our modified xGen capture

14  designs. Trio and relationship phasing proved to be more accurate than EAGLE phasing

15  (Table S1), so we preferentially used the pedigree and relationship data for phasing.

16  Table S2 describes our logic used to determine phase of the pCHMs for the different

17  types of familial relationships. For all remaining pCHMs, we used the EAGLE phased

18  data described above. We excluded any EAGLE phased pCHM where one or both of the

19  variants was a singleton because EAGLE phasing accuracy with singletons was not

20  significantly different than random guessing (Table S3). We found that if the two variants

21  in the pCHM have the same minor allele count (MAC) less than 100, then they are in *cis*

22  (22 out of 22 occurrences in child of trios) in our dataset.

1   We used the trio-phased pCHMs as the truth set to evaluate the overall phasing

2   accuracy of EAGLE. However, if we included the parents of the children of trios in the

3   EAGLE phasing dataset, then EAGLE will use their haplotypes to more accurately phase

4   the children's variants. Having the parental haplotypes in the dataset improves the

5   phasing accuracy of the children's pCHMs but does not provide an accurate estimate of

6   phasing accuracy across the entire dataset. To obtain a good measure of accuracy for the

7   EAGLE pCHM phasing across the entire cohort, we reran EAGLE on the entire dataset

8   as before but excluded all first-degree relatives of one child in each nuclear family before

9   phasing. We then compared the EAGLE phased pCHMs to the trio-phased pCHMs to

10  estimate the overall EAGLE phasing accuracy.

11  Finally, if there were more than one pCHM within the same gene of an individual,

12  then only the pCHM with the most deleterious profile was retained (Table S4). Using the

13  approach outlined above, we were able to phase >99% of all pCHMs, and identify 20,947

14  rare compound heterozygous mutations (CHMs) that are predicted to be function altering.

1

*Figure 2. Decision cascade for determining the phase of potential compound heterozygous*

*mutations (pCHMs) among the 92K DiscovEHR participants. 25.1% of pCHMs and 33.8% of the*

*CHMs (trans) were phased with trios or relationships data.*

5

**Compound heterozygous mutation validation**

We evaluated phasing accuracy by comparing phasing predictions to phasing

done with trios and with Illumina reads. We performed Sanger validation on a subset of

the incorrectly phased pCHMs to see if the variants were false positive calls.

First, phasing accuracy of the pCHMs was evaluated by using the trio phased

pCHMs as truth. Since the phasing approach of each familial relationship is performed

independently from the trio phasing, we can get a good measure of phasing accuracy of

each of the relationship classes as long as the pCHM carrier is a child in a trio. Table S1

shows that the accuracy of family-based phasing was 99.6% (1060/1064 pCHMs) for rare

1 pCHMs. EAGLE phasing was less accurate at 89.1% (766/860 pCHMs; Table S1). We

2 evaluated the accuracy of EAGLE at phasing pCHMs in different minor allele frequency

3 ranges, and found that it consistently attains an accuracy greater than 90% with a MAC

4 greater than 9 and ~77% for a MAC between 2-9 (Table S3). EAGLE phasing performed

5 poorly with singletons.

6      Second, we attempted to validate 200 pCHMs with short Illumina reads (~75 bp)

7 by looking at the read stacks in the Integrative Genomics Viewer (IGV)[45] to see if the two

8 variants occur on the same read or independently. We were able to decisively phase 190

9 (115 *cis* and 79 *trans*; 126 EAGLE phased and 74 pedigree/relationship phased) selected

10 pCHMs using short reads. The remaining ten showed read evidence of both *cis* and *trans*

11 phasing, most likely due to one or both of the variants being a false positive call. Visual

12 validation showed an overall accuracy of 95.8% and 89.9% for pedigree/relationships and

13 EAGLE phasing, respectively (Table S5). While the Illumina read-based validation

14 results are in line with the trio validation results, we do note that the Illumina read-based

15 validation accuracy results are lower than the phasing accuracy determined by phasing

16 with trios. The difference is likely due to the enrichment for false positive pCHMs in

17 small problematic regions of exons prone to sequencing and variant calling errors.

18

19

1    *De novo* **mutation (DNM) detection**

2        We merged the results from two different approaches for detecting DNMs. The first

3    method is TrioDeNovo[46], which reads in the child's and parents' genotype likelihoods at

4    each of the child's variable sites. These likelihoods are input into a Bayesian framework

5    to calculate a posterior likelihood that a child's variant is a DNM. The second program is

6    DeNovoCheck (Web Resources), which is described in the supplemental methods of de

7    Ligt, et al.[47]. DeNovoCheck takes in a set of candidate DNMs identified as being called

8    in the child and not in either parent. It then verifies the presence of the variant in the child

9    and absence in both of the parents by examining the BAM files. We filter these potential

10   DNMs and evaluate a confidence level for each DNM in the union set using a variety of

11   QC metrics. Figure S1 illustrates this DNM calling process, shows the variant filters we

12   applied, and provides the criteria we used to classify each DNM as either low-confidence,

13   moderate-confidence, or high-confidence. We excluded all low-confidence and non-

14   exonic DNMs from the summary results of this paper, but we considered them when

15   doing visual validation to estimate the false negative rate of excluding them. We also

16   excluded the DNM calls for one extreme outlying participant who had an order of

17   magnitude more DNMs called than any other sample.

18

19   **LDLR tandem duplication distant pedigree estimation**

20        Although we cannot know the true family history of the de-identified individuals

21   in our cohort, we have used PRIMUS[31] reconstructed pedigrees, ERSA[12] distant

22   relationship estimate, and PADRE[48] to connect the pedigrees to identify the best pedigree

23   representation of the mutation carriers of a novel tandem duplication in LDLR[49]. We used

1  HumanOmniExpress array data (available for 25 out of the 37 carriers) to estimate the

2  more distant relationships and used the method as described in the PADRE to connected

3  the PRIMUS reconstructed pedigrees.

4

5  **SimProgeny**

6      We developed a forward simulation framework (SimProgeny) to simulate a wide

7  variety of populations, including a population served by a healthcare system like GHS.

8  SimProgeny also simulates sample ascertainment used by HPG studies (Figure S2).

9  SimProgeny can simulate populations of millions of people dispersed across one or more

10  sub-populations based on user specified population parameters (Table S6). The

11  simulation progresses year-to-year simulating couplings, births, separations, migrations,

12  deaths, and movement between sub-populations based on specified parameters. This

13  process generates realistic pedigree structures and populations that represent a wide

14  variety of HPG studies. The default values have been tuned to model the DiscovEHR

15  cohort, but these parameters can be easily customized to model different populations by

16  modifying the configuration file included with the SimProgeny code available at (Web

17  Resources). See Supplemental Methods for a detailed description of SimProgeny.

18    In addition to modeling populations, SimProgeny simulates two ascertainment

19  approaches to model selecting individuals from a population for a genetic study: random

20  ascertainment and clustered sampling. Random ascertainment gives each individual in the

21  population an equal chance of being ascertained without replacement. Clustered sampling

22  is an approach to enrich for close relatives, and it is done by selecting an individual at

23  random along with a number of their first- and second-degree relatives. The number of

1    first-degree relatives is determined by sampling a value from a Poisson distribution with

2    a user specified first-degree ascertainment lambda (default is 0.2). The number of second-

3    degree relatives is determined in the same way and the default second-degree

4    ascertainment lambda is 0.03. See Supplemental Methods for additional information on

5    SimProgeny's ascertainment options.

6

7    **Simulation of the underlying DiscovEHR population and its ascertainment**

8        Our DiscovEHR simulations contained individual populations with starting sizes of

9    200K, 300K, 350K, 400K,475K, 500K, and 550K. We tuned the SimProgeny parameters

10    (Table S6) with publically available country, state, and county level data as well as our

11    own understanding of how individuals were ascertained through GHS consenting and

12    sample collection. Sources for the selected parameters are available in supplemental file

13    Simulation_parameters.xls. We reduced the immigration and emigration rates from the

14    state-wide Pennsylvania (PA) average given that GHS primarily serves rural areas that

15    tend to have lower migration rates than more urban areas. Simulations were run with a

16    burn-in period of 120 years and then progressed for 101 years. Simulated populations

17    grew by ~15%, which is similar to the growth of PA since the mid-20[th] century.

18        We performed both random and clustered ascertainment. For both ascertainment

19    approaches, we shuffled the ascertainment order of the first 5% of the population

20    (specified with the ordered_sampling_proportion parameter) to model the random

21    sequencing order of the individuals in GHS biobank at the beginning of our collaboration.

22    While the selection of this parameter has no effect on random ascertainment and a

23    negligible effect on the accumulation of pairwise relationships in clustered ascertainment,

1  it does affect the proportion of individuals with one or more relatives in the clustered

2  sampling dataset by creating an inflection point at 5% population ascertainment in the

3  simulation results plots (Figure S3B, D). This inflection point would be less pronounced

4  if we were to model the freeze process of the real data or model a smoother transition

5  between sequencing samples from the biobank and newly ascertained individuals.

6  Notably, the inflection point is more pronounced with higher values of lambda from the

7  Poisson distribution.

8

9

10  **Results**

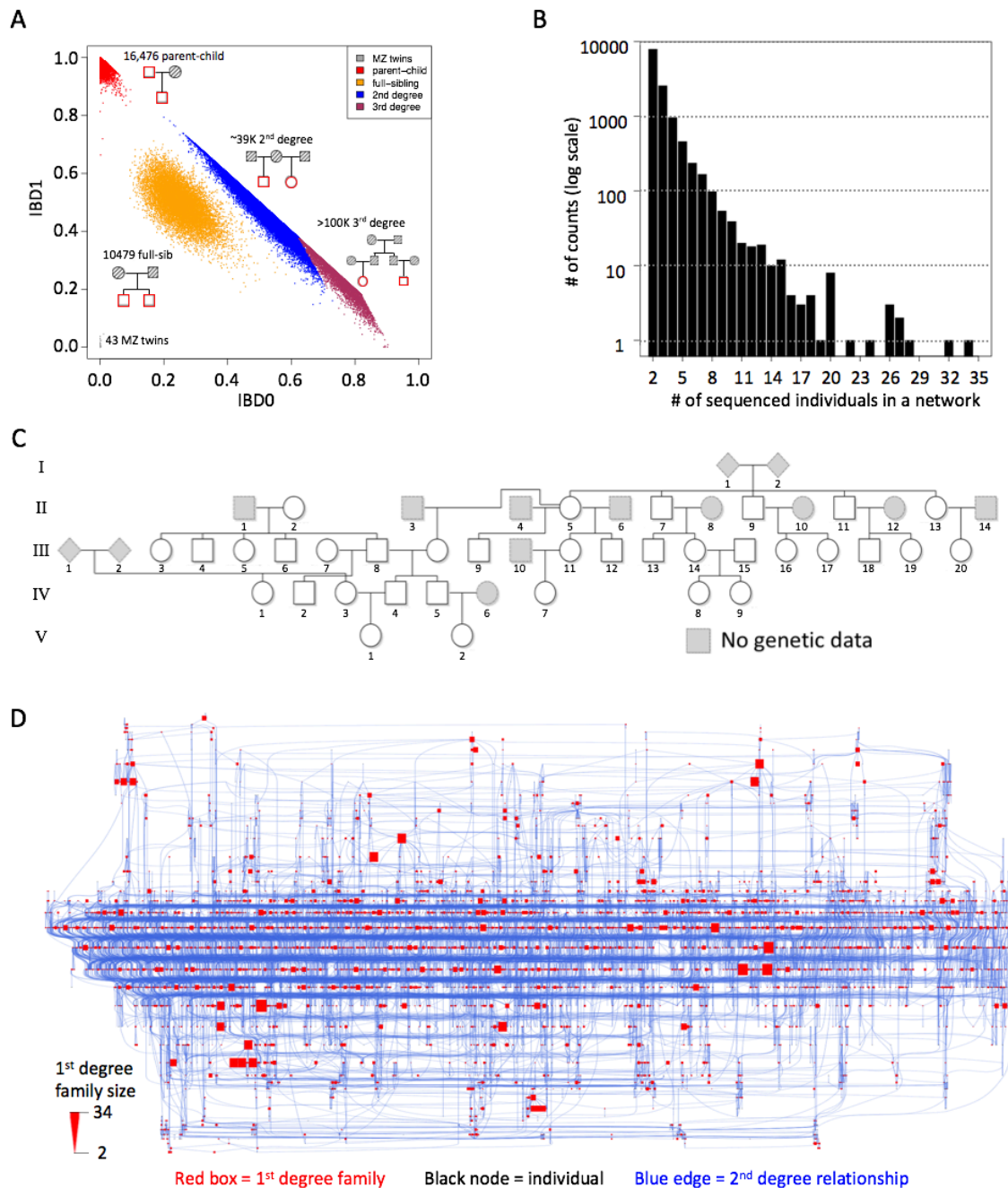11  **Relationship estimation and relatedness in DiscovEHR**

12  In the current dataset of 92,455 individuals, we identified 43 monozygotic twins,

13  16,476 parent-child relationships, 10,479 full-sibling relationships, and ~39,000 second-

14  degree relationships (Figure 3A). Next, we treated individuals as nodes and relationships

15  as edges to generate undirected graphs. Using only first-degree relationships, we

16  identified 7,684 connected components, which we refer to as first-degree family networks.

17  Figure 3B shows the distribution in size of the first-degree family networks, which range

18  from 2 to 25 sequenced individuals. Similarly, we found 10,173 second-degree family

19  networks; the largest containing 19,968 individuals (~22% of the overall dataset; Figure

20  3C). We were able to identify ~5,300 third-degree relationships within the second-degree

21  family networks. Using a lower IBD cutoff (PI_HAT > 0.09875) for the IBD estimations

22  within ancestral groups without consideration of second-degree family networks, we

23  found well over 100,000 third-degree relationships within the DiscovEHR cohort. Given

1    that 95.9% of DiscovEHR individuals are of European ancestry (Table S7), it is not

2    surprising that the vast majority (98.6%) of the pairwise relationships found are between

3    two individuals of European ancestry (Table S8). Nonetheless, we identified many

4    relationships between people of the same, non-European ancestry and between

5    individuals with different ancestries; for example, there are several trios having one

6    European parent, one East Asian parent, and a child whose ancestry is unassigned to a

7    super-population given the ad-mixed nature of his/her genome.

8          Importantly, we show both empirically (Figure 4A) and through simulation

9    (Figure 5) that the rate of accumulating relatives far exceeds the rate of ascertaining

10   samples. This is expected since there are combinatorially increasing numbers of possible

11   pairwise relationships within the dataset as the size increases, and the likelihood that a

12   previously unrelated individual in the dataset becomes involved in a newly identified

13   relationship also increases. Currently, 39% of individuals in the DiscovEHR cohort have

14   at least one first-degree relative in the dataset, and 56% of the participants have one or

15   more first- or second-degree relatives in the dataset (Figure 4B).
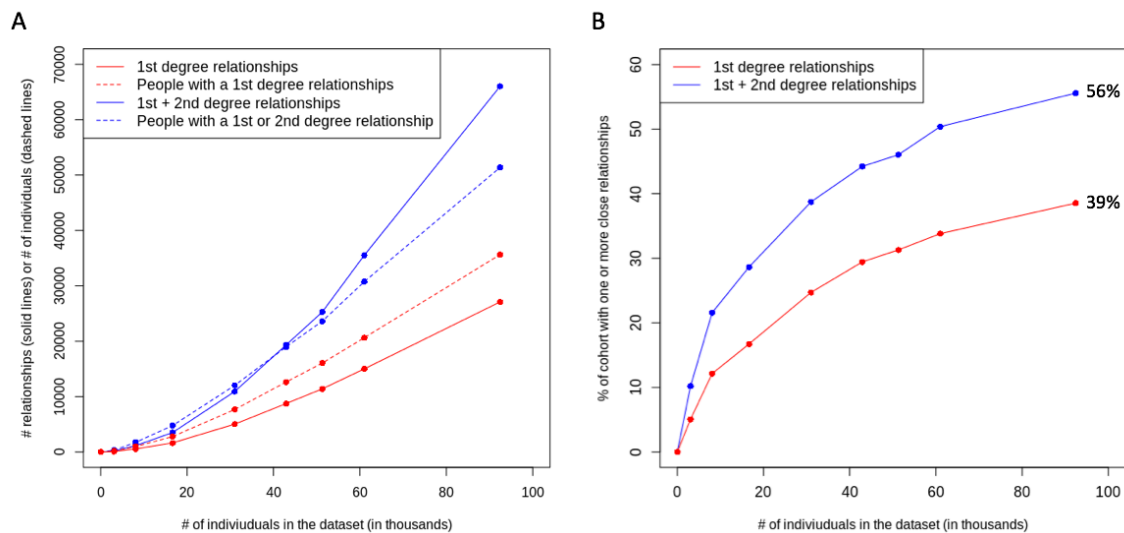
16

17

1

2    *Figure 3. First 92K sequenced individuals from the DiscovEHR cohort contain an extensive*

3    *amount of relatedness. (A) IBD0 vs IBD1 plot shows pairwise relationships segregating into*

4    *different familial relationship classes. The IBD sharing distributions of second- and third-degree*

5    *relationships overlap with each other, so a hard cutoff halfway between the two expected means*

6    *was selected. (B) The distribution of size of first-degree family networks ranges between 2 and 34*

1   *sequenced individuals, with a vast majority of smaller family networks. (C) More than 99.98% of*

2   *the first-degree family networks' pedigree structures were reconstructed using the pairwise-IBD*

3   *estimates, including this pedigree of 34 sequenced individuals. (D) The largest second-degree*

4   *family network of 19,968 (~22% of the dataset) shows 4,062 first-degree family networks (red*

5   *boxes proportionally sized to the number of individuals in the network; including the pedigree*

6   *shown in C) and 5,584 additional individuals (black nodes) connected by 11,430 second-degree*

7   *relationships (blue edges). *Third-degree relationships are challenging to accurately estimate*

8   *due to technical limitations of exome data as well as the widening and overlapping variation*

9   *around the expected mean IBD proportions of more distant relationship classes (e.g. fourth-*

10  *degree and fifth degree). We provided a lower bound estimate of the number of third-degree*

11  *relationships.*

12

13



15  *Figure 4. Accumulation of relatedness within the DiscovEHR cohort at consecutive data freezes.*

16  *A) The number of pairwise relationships has grown rapidly. B) The proportion of individuals in*

17  *the cohort having a first- or second-degree relative identified in the cohort.*

18

1  **Simulations with SimProgeny and relatedness projections**

2  Prior to the launch of the DiscovEHR collaboration, it was unclear how much

3  relatedness we should expect to see and whether it would follow the levels of relatedness

4  seen in previous population-based genomic studies. However, it became clear early on

5  that the cohort contained far more family structure than typically seen in population-

6  based studies, and projections estimated that the proportion of the cohort involved in

7  close relationships would eventually involve the majority of our dataset. Given the

8  impact of this relatedness on downstream analyses, we set out to determine whether this

9  amount of relatedness is expected, whether it is unique to our dataset, and how much it

10  would grow as the sequenced cohort expands.

11  To answer these questions, we developed a flexible simulation framework

12  (SimProgeny) to model a wide variety of study populations and sampling approaches to

13  estimate the amount of relatedness researchers should expect to find for a given set of

14  populations and sampling parameters. While we apply this framework to the DiscovEHR

15  cohort, it is flexible enough that it also can be applied to modeling shallower

16  ascertainment of more transient populations.

17  We used SimProgeny to simulate the DiscovEHR population and the ascertainment of

18  the first 92,455 participants. As expected, the simulations show that DiscovEHR

19  participants were not randomly sampled from the population, but rather the dataset is

20  enriched for close relatives (Figure S4). Therefore, we used a clustered ascertainment

21  approach (see Methods) that more accurately models ascertainment from a healthcare

22  system study population and the subsequent enrichment of close relatives observed in the

23  real data (Figure 5). These simulation results suggest that the effective population size for

1    the first 60K participants was ~475K individuals, and a Poisson distribution having

2    lambda of 0.2 most closely matches the enrichment of first-degree relatives. However,

3    the departure of the real data line (Figure 4, faint red line) from the ~475K simulation

4    line (solid green line) at 90K ascertained samples suggests that the DiscovEHR cohort's

5    effective population size may have increased after ascertaining the first 60K samples.

6    These estimates are consistent with our knowledge that the majority of the first 30K-60K

7    DiscovEHR participants reside in the counties surrounding the GHS headquarters in

8    Danville, and the participant base subsequently expanded to more heavily include pockets

9    of individuals from north-central and northeast rural Pennsylvania (Figure S5). Most

10   notably, ascertainment was not evenly distributed across the entire GHS catchment area

11   (>2.5 million individuals).

12       After identifying simulation parameters that reasonably fit the real data, we used

13   SimProgeny to obtain a projection of the amount of first degree relationships we should

14   expect as DiscovEHR expands to our goal of 250K participants. If we continue to

15   ascertain participants in the same way, we expect to obtain ~150K first-degree

16   relationships (Figure 5C) involving ~60% of DiscovEHR participants (Figure 5D). We

17   then expanded our simulation analysis to include second-degree relationships, and the

18   simulation results suggest that with 250K participants we should expect well over 200K

19   combined first- and second-degree relationships involving over 70% of the individuals in
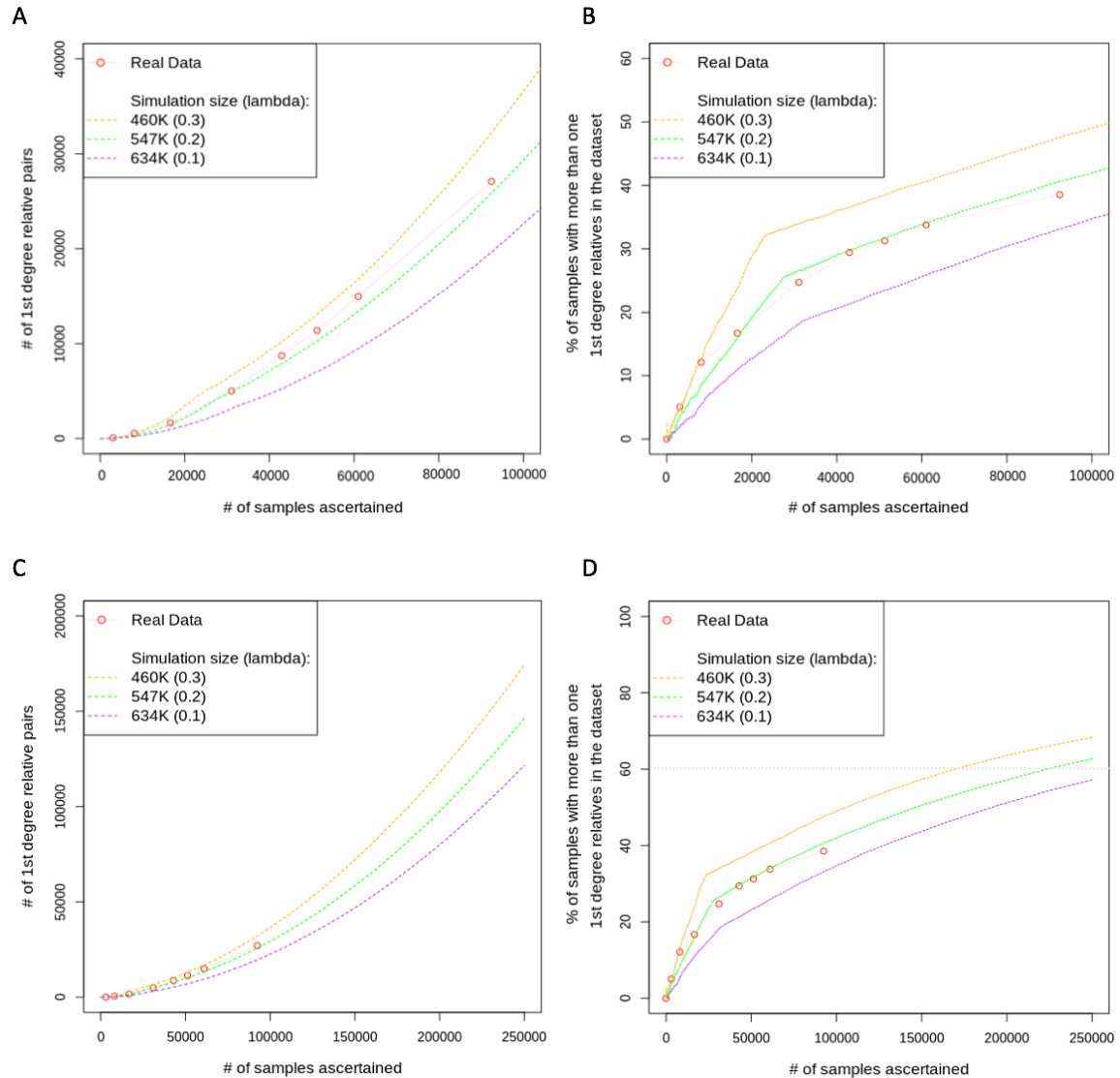
20   DiscovEHR (Figure S3).

21       These projections of relatedness in DiscovEHR assume that we continue ascertaining

22   participants in the same way we did for the first 60K-90K participants. However, if

23   DiscovEHR expands ascertainment of participants to additional GHS clinics and

1  hospitals in other regions, then these relatedness estimates are likely to drop, because

2  expanding the participant base increases the size of the effective sampling population and

3  taps into new genetic demes or distant branches of the same demes. The level of the

4  relatedness will depend on the proportion of the total population we ascertain and the

5  underlying population demographics of the regions, both of which can be simulated with

6  SimProgeny.

7      While SimProgeny is designed to reasonably model a real population, it has its

8  limitations. For example, the Poisson distribution accurately models the clustered

9  sampling of first-degree relationships we observe in our real data, but it underestimates

10  the clustering of second-degree relationships compared to what is observed in

11  DiscovEHR. Thus, there are bound to be nuances that cannot be easily modeled with a

12  fixed distribution, and there are likely to be other confounding aspects to how participants

13  were ascertained in the real dataset.

14      Regardless, our simulation results demonstrate a clear enrichment of relatedness in

15  the DiscovEHR HPG study as well as provide key insights into the tremendous amount of

16  relatedness we expect to see as we continue to ascertain additional participants, assuming

17  future ascertainment is reasonably well modeled by SimProgeny. These observations can

18  also be extrapolated to other large HPG studies, and the flexibility built into the model

19  provides the ability to tune the model to a wide variety of different populations and

20  ascertainment approaches.

21

*Figure 5. Simulated population and ascertainment fit to the accumulation of first-degree*

*relatedness in the DiscovEHR cohort. The real data was calculated at periodic "freezes"*

*indicated with the punctuation points connected by the faint red line. Most simulation parameters*

*were set based on information about the real population demographics and the DiscovEHR*

*ascertainment approach. However, two parameters were unknown and selected based on fit to*

*the real data: 1. the effective population size from which samples were ascertained and 2. the*

*increased chance that someone is ascertained given a first-degree relative previously ascertained,*

*which we call "clustered ascertainment". All panels show the same three simulated population*

1    *sizes spanning the estimated effective population size. We simulated clustered ascertainment by*

2    *randomly ascertaining an individual along with a Poisson-distributed random number of 1ˢᵗ*

3    *degree relatives (distributions' lambdas are indicated in the legends). (A) The accumulation of*

4    *pairs of first-degree relatives as additional samples are ascertained. (B) The proportion of the*

5    *ascertained participants that have one or more first-degree relatives that have also been*

6    *ascertained. (C) Simulated ascertainment projections with upper and lower bounds of the number*

7    *of first-degree relationships we expect with our current DiscovEHR ascertainment approach as*

8    *we scale to our goal of 250K participants. (D) Simulated projections with upper and lower*

9    *bounds of the proportion of the ascertained participants that have 1 or more first-degree relatives*

10    *that have also been ascertained.*

11

12    **Leveraging relatedness instead of treating it like a nuisance**

13    We have reconstructed pedigree structures for 12,574 first-degree family networks in

14    the DiscovEHR data set using the pedigree reconstruction tool PRIMUS[31], and found that

15    98.9% of these pedigrees reconstructed unambiguously to a single pedigree structure

16    when considering IBD estimates and reported participant ages. These pedigrees include

17    2,192 nuclear families (1,841 trios, 297 quartets, 50 quintets, 3 sextets, and 1 septet).

18    Table S9 shows a breakdown of the trios by ancestry. Figure 3C shows the largest first-

19    degree pedigrees, which contains 34 sequenced individuals. We have used these

20    relationships and pedigrees in several ways, and we highlight three main applications in

21    this section.

22

23

24

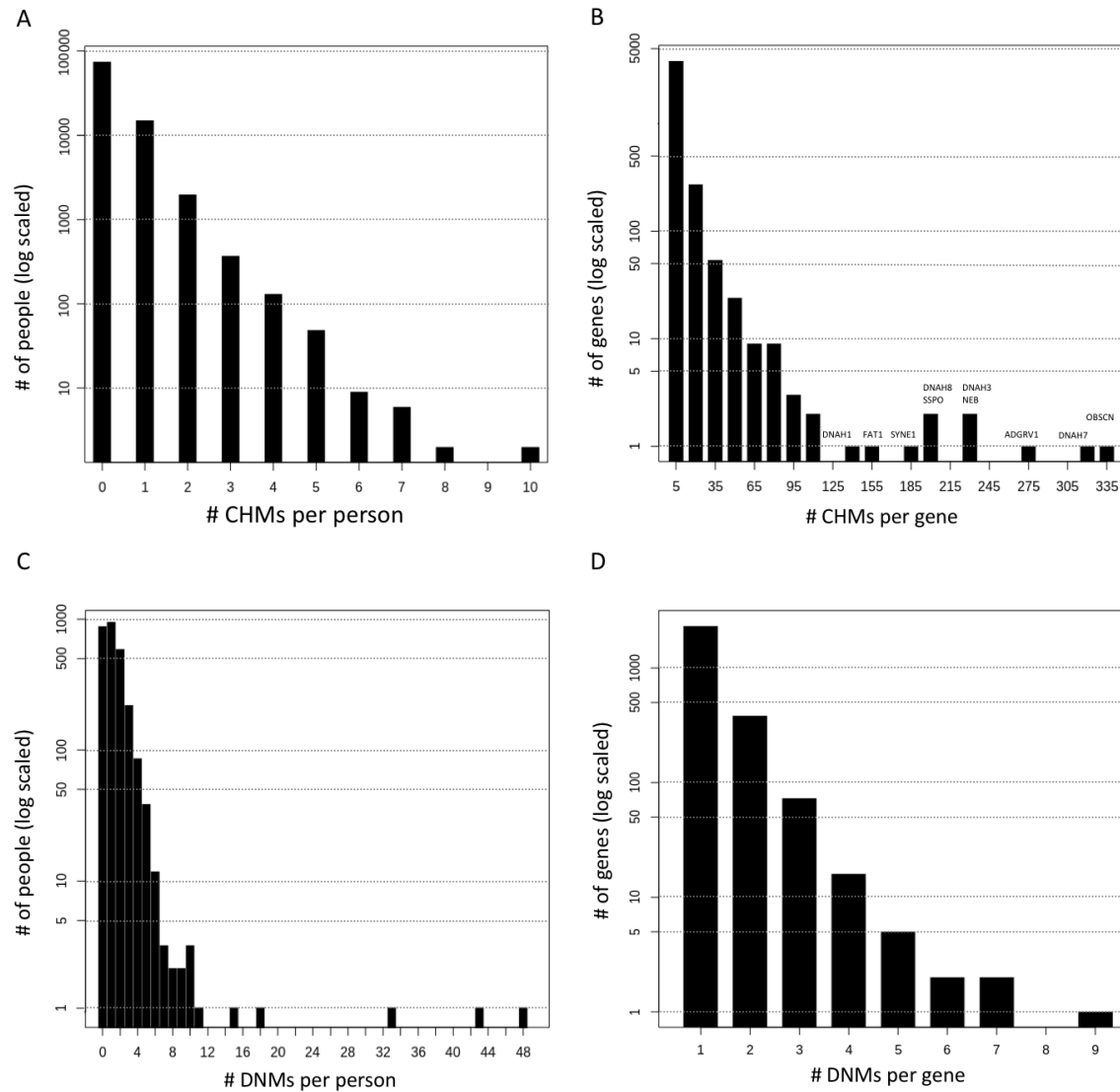1   *Compound Heterozygous mutations*

2       A major goal of human genetics is to better understand the function of every gene in

3   the human genome. Homozygous loss-of-function mutations (LoFs) are a powerful tool

4   to gain insight into gene function by analyzing the phenotypic effects of these "human

5   knockouts" (KOs). Rare (MAF < 1%) homozygous LoFs have been highlighted in recent

6   large-scale sequencing studies and have been critical in identifying many gene-phenotype

7   interactions[1,4,50,51]. While rare compound heterozygous mutations (CHMs) of two

8   heterozygous LoFs are functionally equivalent to rare homozygous KOs, they are more

9   difficult to identify (particularly with short-read sequencing) and are rarely interrogated

10  in large sequencing studies[1,4,50].

11      We performed a survey of rare CHMs in the DiscovEHR cohort. First, we identified

12  57,355 high-quality potential CHMs (pCHMs) consisting of pairs of rare heterozygous

13  variants that are either putative LoFs (pLoF, i.e., nonsense, frameshift, or splice-site

14  mutations) or missense variants with strong evidence of being deleterious (see Methods).

15  Second, we phased the pCHMs using a combination of allele-frequency-based phasing

16  using EAGLE and pedigree-based phasing using the reconstructed pedigrees and

17  relationship data (Figure 2). EAGLE phased the pCHMs with an average of 89.1%

18  accuracy based on trio validation (Table S1). However, because we had extensive

19  pedigree and relationship data within this cohort, we were able to use it to phase 25.2% of

20  the pCHMs and 33.8% of the *trans* CHMs with highly accurate trio and relationship

21  phasing data (≥98.0%; Table S1), reducing inaccurate phasing of *trans* CHMs by

22  approximately a third. The phased pCHMs spanned the entire frequency range from

23  singletons to 1% MAF (Table S10).

1    After processing, 40.3% of the pCHMs were phased in *trans*, yielding a high-

2    confidence set of 20,947 rare, deleterious CHMs distributed among 17,533 of the 92K

3    individuals (mean = 0.23 per person; max = 10 per person; Figure 6A). The median

4    genomic distance between pCHM variants in *cis* (5,955 bps) was a little more than half

5    the median distance between the pCHMs variants in *trans* (11,600 bps; Figure S6).

6    Nearly a third of the CHMs involved at least one pLoF and 8.9% of the of CHMs

7    consisted of two pLoF variants (Table S11). Over 4,216 of the 19,467 targeted genes

8    contain one or more CHM carriers (Table S12), and 2,468 have more than one carrier

9    (Figure 6B). The ten genes with more than 125 CHM carriers are estimated to be among

10   the most LoF tolerant in the genome based on ExAC pLI scores[4] (Table S13), so it is no

11   surprise that these genes would contain a higher number of CHMs.

12   In order to get a more robust set of human knockout genes and demonstrate the added

13   value of CHMs, we combined the CHMs with the 6,560 rare (MAF < 1%) homozygous

14   pLoFs found among the 92K DiscovEHR participants. pLoF-pLoF CHMs increased the

15   number of genes with $\geq 1$ and $\geq 20$ individuals with a putative KO by 15% and 61%,

16   respectively (Table S12). The benefit of including CHMs in a KO analysis is even more

17   significant when we consider missense variants that are predicted to disrupt protein

18   function. We found a combined 20,364 rare homozygous pLOF and deleterious missense

19   variants among the 92K participants. CHMs provided 26% more genes with $\geq 1$ carriers

20   and 397% more genes with $\geq 20$ carriers where both copies of the gene are predicted to be

21   completely knocked out or disrupted (Table S12).

22

1

*Figure 6. DiscovEHR results for compound heterozygous mutations (CHMs) and de novo*

*mutations (DNMs). (A) Distribution of the number of CHMs per individual in the DiscovEHR*

*cohort. (B) Distribution of the # of CHMs per gene. Names of genes with more than 125 CHMs*

*are listed. (C) Distribution of 3,415 exonic high and moderate confidence DNMs among the*

*children of trios in the DiscovEHR cohort. (D) The distribution non-synonymous DNMs across*

*the 2,802 genes with 1 or more.*

8

9

1    *De novo mutations*

2    *De novo* mutations (DNMs) are a class of rare variation that is more likely to produce

3    extreme phenotypes in humans due to sporadic occurrence and lack of purifying selection.

4    Many recent sequencing studies have shown that DNMs are a major driver in human

5    genetic disease[47,53,54], demonstrating that DNMs are a valuable tool to better understand

6    gene function.

7    We used the nuclear families reconstructed from the 92K DiscovEHR participants to

8    confidently call 3,415 moderate- and high-confidence exonic DNMs distributed among

9    1,783 of the 2,602 available children in trios (mean = 1.31; max = 48; Figure 6C).

10   PolyPhen2 predicts 29.1% (N=995) of the DNMs as "probably damaging" and an

11   additional 9.2% (N=316) as possibly damaging. The DNMs are distributed across 2,802

12   genes (Figure 6D) with *TTN* receiving the most at nine. The most common type of DNM

13   is nonsynonymous SNVs (58.5%) followed by synonymous SNVs (24.3%). Table 1

14   provides a complete breakdown of DNM types and shows that our proportions of DNMs

15   falling into the different functional classes generally match those found in a recent study

16   of DNMs in children with development disorders[53]. We also observed an increase in the

17   number of exonic DNMs with respect to both maternal (0.011 DNMs/year, $p=7.3 \times 10^{-4}$;

18   Poisson regression; Figure S7) and paternal age at birth (0.010 DNMs/year; $p=5.6 \times 10^{-4}$),

19   consistent with other reports[53,55-57]. Notably, maternal and paternal age at birth are highly

20   correlated in our dataset (rho=0.79; Figure S8), thus the rates are not additive and no

21   significant difference was identified to distinguish either as a driving factor (see

22   Supplemental Methods).

1    We attempted to perform visual validation of 23 high- and 30 moderate- and 47 low-

2    confidence DNMs spanning all functional classes. Eight moderate- and two low-

3    confidence variants could not be confidently called as true or false positive DNMs. Of

4    those remaining, 23/23 (100%) high-confidence, 19/22 (86%) moderate-confidence, and

5    12/43 (28%) low-confidence DNMs validated as true positives. Visual validation also

6    confirmed that the majority (40/49) of potential DNM in individuals with >10 DNMs are

7    likely false positive calls.

8

9    *Table 1. Breakdown by functional class of moderate- and high-confidence exonic DNMs found in*

10    *the DiscovEHR cohort compared to a recent developmental delay exome study of 4,293 trios.*

| Type of DNM | # of DNMs | % of DNMs | # in DDD study[a] | % in DDD study[a] |
|---|---|---|---|---|
| nonsynonymous SNV | 1,996 | 58.5% | 4,797 | 57.8% |
| synonymous SNV | 831 | 24.3% | 1,629 | 19.6% |
| splicing | 153 | 4.5% | 671 | 8.1% |
| non-frameshift deletion | 78 | 2.3% | 167 | 2.0% |
| non-frameshift insertion | 55 | 1.6% | 28 | 0.3% |
| frameshift | 187 | 5.5% | 603 | 7.3% |
| stop-gain SNV | 112 | 3.3% | 402 | 4.8% |
| stop-loss SNV | 3 | 0.1% | 7 | 0.1% |

11    *[a]The Deciphering Developmental Disorders Study (DDD)[53]. The DDD paper also reported 57*

12    *DNMs of other classes that were not included in our analysis nor in this table; percentages were*

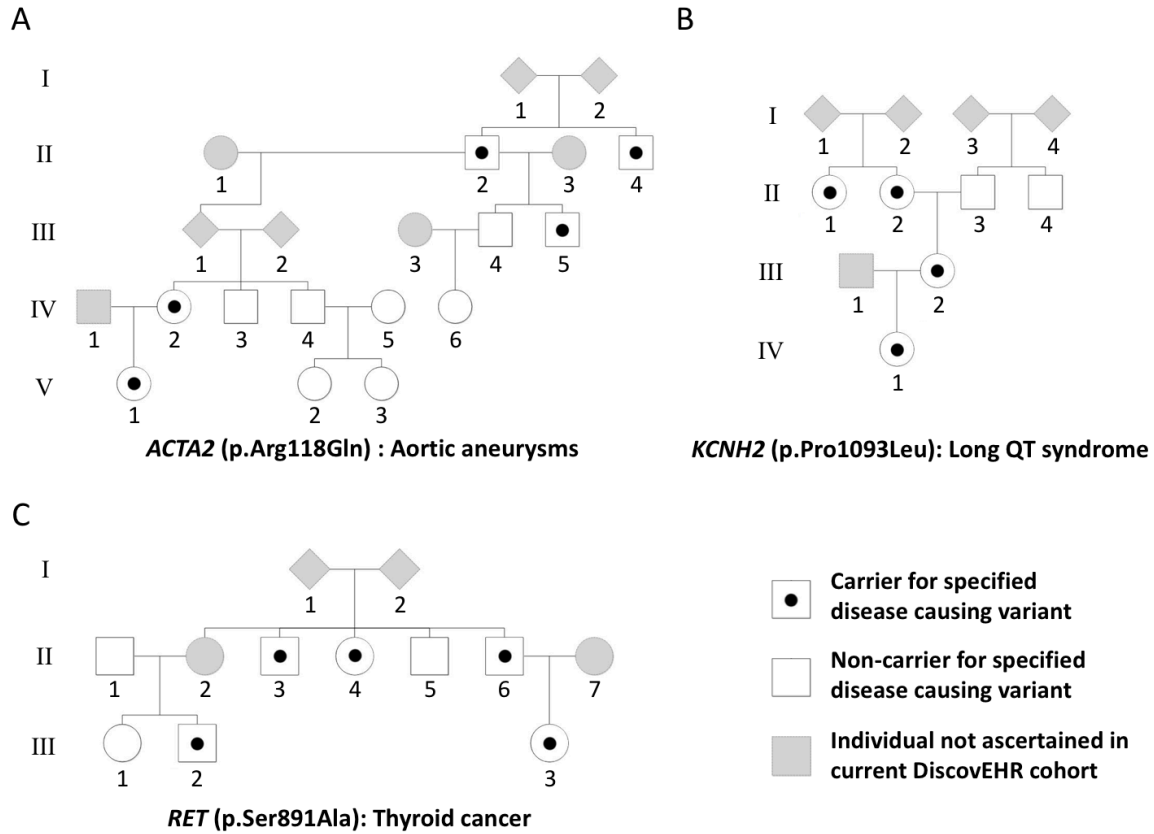13    *adjusted accordingly.*

14

15

1    *Variant and phenotype segregation in pedigrees*

2    We have used the reconstructed pedigree data from among the 92K DiscovEHR

3    participants to distinguish between novel/rare population variation and familial variants

4    and have leveraged it to identify highly penetrant disease variants segregating in families.

5    While this is not intended to be a survey of all known Mendelian disease-causing

6    variation transmitted through these pedigrees, we have identified a few illustrative

7    examples including familial aortic aneurysms (Figure 7A), long QT syndrome (Figure

8    7B), thyroid cancer (Figure 7C), and familial hypercholesterolemia (FH; Figure 8)[49]. The

9    FH example is particularly interesting as we previously reported a novel FH-causing

10    tandem duplication in *LDLR*[49]. We have updated the CNV calls and found 37 carriers of

11    the FH-causing tandem duplication among the 92K exomes, and we have reconstructed

12    30 out of the 37 carriers into a single extended pedigree. The carriers' shared ancestral

13    history provides evidence that they all inherited this duplication event from a common

14    ancestor approximately six generations back. While two of the seven remaining carriers

15    are second-degree relatives to each other, genotyping array data was not available to

16    confirm that the remaining seven carriers are also distantly related to the other carriers in
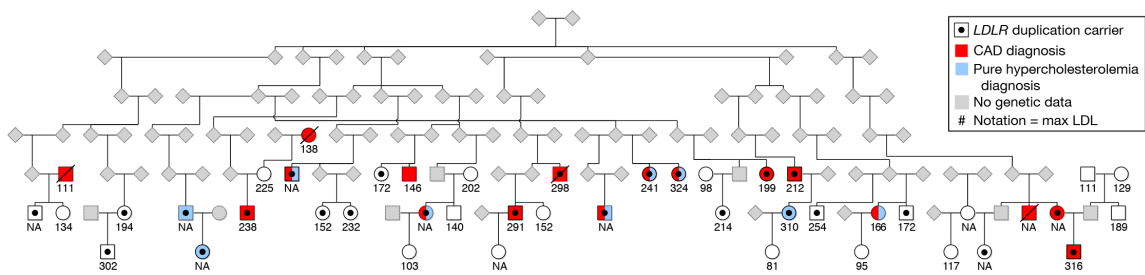
17    Figure 8.

18

*Figure 7. Reconstructed pedigree from DiscovEHR demonstrating the segregation of known disease-causing variants, including variants for (A) aortic aneurysms, (B) long QT syndrome, and (C) thyroid cancer.*



*Figure 8. Image of the reconstructed pedigree prediction containing 25/37 carriers of the novel FH-causing tandem duplication in LDLR and 20 non-carrier, related (first- or second-degree)*

1    *individuals from the sequenced cohort. Carrier and non-carrier status was determined from the*

2    *exome data from each individual. Elevated max LDL levels (value under symbols) as well as*

3    *increased prevalence of coronary artery disease (CAD, red fill) and pure hypercholesterolemia*

4    *(ICD 272.0; blue) segregate with duplication carriers. Five additional carriers (not drawn) were*

5    *found to be distant relatives (seventh- to ninth-degree relatives) of individuals in this pedigree.*

6

7

## Discussion

9      Sequencing studies continue to collect and sequence increasing proportions of

10   human populations and are uncovering the extremely complex, intertwined nature of

11   human relatedness. In the first 92K sequenced participants of the DiscovEHR cohort, we

12   have identified ~66K first- and second-degree relationships, reconstructed 12,574

13   pedigrees, and uncovered a second-degree family network of nearly 20,000 participants.

14   Studies in founder populations have already highlighted the complexity of relationships

15   (Old Order Amish[58], Hutterites[59], and Ashkenazi Jews[60]), and recent studies of non-

16   founder populations are reporting extensive levels of family structure (UK Biobank[61],

17   NHANES[62], and AncestryDNA[6]). We observed family structure (first- and second-degree

18   relationships) involving 55.6% of the DiscovEHR participants, and we expect family

19   structure will involve a large proportion, if not a majority, of individuals in other large

20   HPG studies. We have demonstrated through simulations and observations within our

21   own data that we can obtain a large number of close familial relationships, nuclear

22   families, and informative pedigrees within HPG studies. While underlying population

23   structure and depth of ascertainment will vary between studies, we do believe that our

1    observations in DiscovEHR will be applicable to other HPG studies since families tend to

2    visit the same healthcare system and have similar genetic and environmental disease risks.

3    The days of only having a handful of closely related individuals or samples in large

4    sequencing cohorts are over, and we can no longer simply remove closely related pairs of

5    individuals for our association studies knowing that it is only a small fraction of the

6    overall cohort. Instead, we need to continue developing methods that are capable of

7    leveraging the extensive relatedness of these rich cohorts and that can scale to

8    accommodate growing HPG study population sizes and phenotype diversity.

9         In this study, we have demonstrated several ways to leverage family structure.

10   First, we improved the phasing accuracy of rare compound heterozygous mutations

11   (CHMs). While we did obtain accurate phasing of CHMs with EAGLE, our pedigree-

12   and relationship-based phasing was far more accurate, reducing the pCHM phasing error

13   by approximately a third. We expect that the accuracy of the relationship-based phasing

14   of pCHMs will be lower for variants with >1% MAF because phasing using the pairwise

15   relationships assumes that if two variants appear together in two relatives, then they are

16   in *cis* and have segregated together from a common ancestor. There is a higher chance

17   that two independently segregating common variants will appear together in multiple

18   people, resulting in being incorrectly phased as *cis* by the algorithm. Therefore, common

19   variants may be better phased using population allele frequencies with programs like

20   EAGLE rather than phased using pairwise relationships.

21        Second, pedigree reconstruction within HPG studies provides trios and other

22   informative pedigree structures that can be leveraged for many use-cases. We used the

23   2,602 reconstructed trios to find 3,415 DNMs and tracked known disease-causing

1    mutations through extended pedigrees. Pedigrees and relationships are also particularly

2    useful for tracking transmission of rare variants, providing increased confidence in

3    variant calls and allowing for the use of more traditional Mendelian genetic analyses.

4    Pedigrees can be particularly useful when combined with follow-up chart reviews and the

5    ability to recontact patients and their family members.

6        We show that cryptic family structure in a large sequencing dataset presents an

7    opportunity to harness a valuable, untapped source of genetic insights rather than a

8    nuisance that must be managed during downstream analyses. As we enter the era of

9    genomic-based precision medicine, we see a critical need for additional innovative

10    methods and tools that are capable of effectively mining the familial structure and distant

11    relatedness contained within the ever-growing sequencing cohorts.

12

13

## 14    Description of Supplemental Data

15    Additional methods, thirteen tables, eight figures, and one excel file.

16

17

## 18    Conflict of interest

1    by Regeneron Pharmaceuticals (62/555,597), which discloses the simulation framework

2    and methods for identifying the familial relationships, phasing the pCHMs, and calling

3    the DNMs. Additional information for reproducing the results described in the article is

4    available upon reasonable request and subject to a data use agreement.

5

6

## Acknowledgments

8    We thank the MyCode Community Health Initiative participants for their permission to

9    use their health and genomics information in the DiscovEHR collaboration.

10

11

## Web resources

13    DeNovoCheck - https://sourceforge.net/projects/denovocheck

14    FBAT - https://www.hsph.harvard.edu/fbat/fbat.htm

15    GATK - https://software.broadinstitute.org/gatk/

16    QTDT - http://csg.sph.umich.edu/abecasis/qtdt/

17    SimProgeny - https://github.com/rgcgithub/SimProgeny

18    TopMed - https://www.nhlbiwgs.org

19

20

## Reference

22    1. Dewey, F.E., Murray, M.F., Overton, J.D., Habegger, L., Leader, J.B., Fetterolf, S.N.,

23    O'Dushlaine, C., Van Hout, C.V., Staples, J., Gonzaga-Jauregui, C., et al. (2016).

Distribution and clinical impact of functional variants in 50,726 whole-exome sequences from the DiscovEHR study. Science *354*, aaf6814–aaf6814.

2. Sudlow, C., Gallacher, J., Allen, N., Beral, V., Burton, P., Danesh, J., Downey, P., Elliott, P., Green, J., Landray, M., et al. (2015). UK biobank: an open access resource for identifying the causes of a wide range of complex diseases of middle and old age. PLoS Med. *12*, e1001779.

3. Collins, F.S., and Varmus, H. (2015). A new initiative on precision medicine. N. Engl. J. Med. *372*, 793–795.

4. Lek, M., Karczewski, K.J., Minikel, E.V., Samocha, K.E., Banks, E., Fennell, T., O'Donnell-Luria, A.H., Ware, J.S., Hill, A.J., Cummings, B.B., et al. (2016). Analysis of protein-coding genetic variation in 60,706 humans. Nature Publishing Group *536*, 285–291.

5. Henn, B.M., Hon, L., Macpherson, J.M., Eriksson, N., Saxonov, S., Pe'er, I., and Mountain, J.L. (2012). Cryptic distant relatives are common in both isolated and cosmopolitan genetic samples. PLoS ONE *7*, e34267.

6. Han, E., Carbonetto, P., Curtis, R.E., Wang, Y., Granka, J.M., Byrnes, J., Noto, K., Kermany, A.R., Myres, N.M., Barber, M.J., et al. (2017). Clustering of 770,000 genomes reveals post-colonial population structure of North America. Nat Commun *8*, 14238.

7. Agarwala, R., Biesecker, L.G., and Schäffer, A.A. (2003). Anabaptist genealogy database. Am J Med Genet C Semin Med Genet *121C*, 32–37.

1    8. Price, A.L., Butler, J., Patterson, N., Capelli, C., Pascali, V.L., Scarnicci, F., Ruiz-

2    Linares, A., Groop, L., Saetta, A.A., Korkolopoulou, P., et al. (2008). Discerning the

3    ancestry of European Americans in genetic association studies. PLoS Genet *4*, e236.

4    9. Rosenberg, N.A., Pritchard, J.K., Weber, J.L., Cann, H.M., Kidd, K.K., Zhivotovsky,

5    L.A., and Feldman, M.W. (2002). Genetic structure of human populations. Science *298*,

6    2381–2385.

7    10. Belbin, G.M., Ruderfer, D.M., Stahl, E.A., Jeff, J.M., Loos, R.J.F., Bottinger, E.P.,

8    Abul-Husn, N.S., Auton, A., and Kenny, E.E. (2015). Abstract: Reconstructing the

9    population history of New York City. Presented at American Society of Human Genetics

10    Annual Meeting. Baltimore, Md.

11    11. Stevens, E.L., Baugher, J.D., Shirley, M.D., Frelin, L.P., and Pevsner, J. (2012).

12    Unexpected relationships and inbreeding in HapMap phase III populations. PLoS ONE *7*,

13    e49575.

14    12. Huff, C.D., Witherspoon, D.J., Simonson, T.S., Xing, J., Watkins, W.S., Zhang, Y.,

15    Tuohy, T.M., Neklason, D.W., Burt, R.W., Guthery, S.L., et al. (2011). Maximum-

16    likelihood estimation of recent shared ancestry (ERSA). Genome Res. *21*, 768–774.

17    13. Fuchsberger, C., Flannick, J., Teslovich, T.M., Mahajan, A., Agarwala, V., Gaulton,

18    K.J., Ma, C., Fontanillas, P., Moutsianas, L., McCarthy, D.J., et al. (2016). The genetic

19    architecture of type 2 diabetes. Nature Publishing Group *536*, 41–47.

20    14. Locke, A.E., Kahali, B., Berndt, S.I., Justice, A.E., Pers, T.H., Day, F.R., Powell, C.,

21    Vedantam, S., Buchkovich, M.L., Yang, J., et al. (2015). Genetic studies of body mass

1    index yield new insights for obesity biology. Nature *518*, 197–206.

2    15. Surendran, P., Drenos, F., Young, R., Warren, H., Cook, J.P., Manning, A.K., Grarup,

3    N., Sim, X., Barnes, D.R., Witkowska, K., et al. (2016). Trans-ancestry meta-analyses

4    identify rare and common variants associated with blood pressure and hypertension. Nat

5    Genet *48*, 1151–1161.

6    16. Santorico, S.A., and Edwards, K.L. (2014). Challenges of linkage analysis in the era

7    of whole-genome sequencing. Genet. Epidemiol. *38 Suppl 1*, S92–S96.

8    17. Hu, H., Roach, J.C., Coon, H., Guthery, S.L., Voelkerding, K.V., Margraf, R.L.,

9    Durtschi, J.D., Tavtigian, S.V., Shankaracharya, Wu, W., et al. (2014). A unified test of

10    linkage analysis and rare-variant association for analysis of pedigree sequence data. Nat.

11    Biotechnol. *32*, 663–669.

12    18. Price, A.L., Zaitlen, N.A., Reich, D., and Patterson, N. (2010). New approaches to

13    population stratification in genome-wide association studies. Nat. Rev. Genet. *11*, 459–

14    463.

15    19. Kang, H.M., Sul, J.H., Service, S.K., Zaitlen, N.A., Kong, S.-Y., Freimer, N.B.,

16    Sabatti, C., and Eskin, E. (2010). Variance component model to account for sample

17    structure in genome-wide association studies. Nature Publishing Group *42*, 348–354.

18    20. Sun, L., and Dimitromanolakis, A. (2012). Identifying cryptic relationships. Methods

19    Mol. Biol. *850*, 47–57.

20    21. Devlin, B., and Roeder, K. (1999). Genomic Control for Association Studies.

1    Biometrics *55*, 997–104.

2    22. Voight, B.F., and Pritchard, J.K. (2005). Confounding from Cryptic Relatedness in

3    Case-Control Association Studies. PLoS Genet *1*, e32–10.

4    23. Staples, J., Nickerson, D.A., and Below, J.E. (2013). Utilizing graph theory to select

5    the largest set of unrelated individuals for genetic analysis. Genet. Epidemiol. *37*, 136–

6    141.

7    24. Chang, C.C., Chow, C.C., Tellier, L.C., Vattikuti, S., Purcell, S.M., and Lee, J.J.

8    (2015). Second-generation PLINK: rising to the challenge of larger and richer datasets.

9    Gigascience *4*, 7.

10   25. Zhang, Z., Ersoz, E., Lai, C.-Q., Todhunter, R.J., Tiwari, H.K., Gore, M.A., Bradbury,

11   P.J., Yu, J., Arnett, D.K., Ordovas, J.M., et al. (2010). Mixed linear model approach

12   adapted for genome-wide association studies. Nat Genet *42*, 355–360.

13   26. Yang, J., Zaitlen, N.A., Goddard, M.E., Visscher, P.M., and Price, A.L. (2014).

14   Advantages and pitfalls in the application of mixed-model association methods. Nat

15   Genet *46*, 100–106.

16   27. Kirkpatrick, B., and Bouchard-Côté, A. (2016). Correcting for Cryptic Relatedness in

17   Genome-Wide Association Studies. arXiv *q-bio.QM*.

18   28. Day-Williams, A.G., Blangero, J., Dyer, T.D., Lange, K., and Sobel, E.M. (2011).

19   Linkage analysis without defined pedigrees. Genet. Epidemiol. *35*, 360–370.

20   29. Conomos, M.P., Reiner, A.P., Weir, B.S., and Thornton, T.A. (2016). Model-free

1     Estimation of Recent Genetic Relatedness. Am. J. Hum. Genet. *98*, 127–148.

2     30. Conomos, M.P., Miller, M.B., and Thornton, T.A. (2015). Robust inference of

3     population structure for ancestry prediction and correction of stratification in the presence

4     of relatedness. Genet. Epidemiol. *39*, 276–293.

5     31. Staples, J., Qiao, D., Cho, M.H., Silverman, E.K., University of Washington Center

6     for Mendelian Genomics, Nickerson, D.A., and Below, J.E. (2014). PRIMUS: rapid

7     reconstruction of pedigrees from genome-wide estimates of identity by descent. Am. J.

8     Hum. Genet. *95*, 553–564.

9     32. Ko, A., and Nielsen, R. (2017). Composite likelihood method for inferring local

10     pedigrees. PLoS Genet *13*, e1006963.

11     33. Wijsman, E.M. (2012). The role of large pedigrees in an era of high-throughput

12     sequencing. Hum. Genet. *131*, 1555–1563.

13     34. Wijsman, E.M., Rothstein, J.H., and Thompson, E.A. (2006). Multipoint linkage

14     analysis with many multiallelic or dense diallelic markers: Markov chain-Monte Carlo

15     provides practical approaches for genome scans on general pedigrees. The American

16     Journal of Human Genetics *79*, 846–858.

17     35. Thornton, T., and McPeek, M.S. (2010). ROADTRIPS: case-control association

18     testing with partially or completely unknown population and pedigree structure. Am. J.

19     Hum. Genet. *86*, 172–184.

20     36. Sul, J.H., Cade, B.E., Cho, M.H., Qiao, D., Silverman, E.K., Redline, S., and Sunyaev,

S. (2016). Increasing Generality and Power of Rare-Variant Tests by Utilizing Extended Pedigrees. Am. J. Hum. Genet. *99*, 846–859.

37. He, Z., Zhang, D., Renton, A.E., Li, B., Zhao, L., Wang, G.T., Goate, A.M., Mayeux, R., and Leal, S.M. (2017). The Rare-Variant Generalized Disequilibrium Test for Association Analysis of Nuclear and Extended Pedigrees with Application to Alzheimer Disease WGS Data. Am. J. Hum. Genet. *100*, 371.

38. Ott, J., Kamatani, Y., and Lathrop, M. (2011). Family-based designs for genome-wide association studies. Nat. Rev. Genet. *12*, 465–474.

39. Abul-Husn, N.S., Manickam, K., Jones, L.K., Wright, E.A., Hartzel, D.N., Gonzaga-Jauregui, C., O'Dushlaine, C., Leader, J.B., Lester Kirchner, H., Lindbuchler, D.M., et al. (2016). Genetic identification of familial hypercholesterolemia within a single U.S. health care system. Science *354*, aaf7000–aaf7000.

40. International HapMap 3 Consortium, Altshuler, D.M., Gibbs, R.A., Peltonen, L., Dermitzakis, E., Schaffner, S.F., Yu, F., Bonnen, P.E., de Bakker, P.I.W., Deloukas, P., et al. (2010). Integrating common and rare genetic variation in diverse human populations. Nature Publishing Group *467*, 52–58.

41. Loh, P.-R., Danecek, P., Palamara, P.F., Fuchsberger, C., A Reshef, Y., K Finucane, H., Schoenherr, S., Forer, L., McCarthy, S., Abecasis, G.R., et al. (2016). Reference-based phasing using the Haplotype Reference Consortium panel. Nat Genet *48*, 1443–1448.

42. Adzhubei, I., Jordan, D.M., and Sunyaev, S.R. (2013). Predicting functional effect of

1   human missense mutations using PolyPhen-2. Curr Protoc Hum Genet *Chapter 7*,

2   Unit7.20–7.20.41.

3   43. Chun, S., and Fay, J.C. (2009). Identification of deleterious mutations within three

4   human genomes. Genome Res. *19*, 1553–1561.

5   44. Schwarz, J.M., Cooper, D.N., Schuelke, M., and Seelow, D. (2014). MutationTaster2:

6   mutation prediction for the deep-sequencing age. Nat. Methods *11*, 361–362.

7   45. Robinson, J.T., Thorvaldsdóttir, H., Winckler, W., Guttman, M., Lander, E.S., Getz,

8   G., and Mesirov, J.P. (2011). Integrative genomics viewer. Nat. Biotechnol. *29*, 24–26.

9   46. Wei, Q., Zhan, X., Zhong, X., Liu, Y., Han, Y., Chen, W., and Li, B. (2015). A

10  Bayesian framework for de novo mutation calling in parents-offspring trios.

11  Bioinformatics *31*, 1375–1381.

12  47. de Ligt, J., Willemsen, M.H., van Bon, B.W.M., Kleefstra, T., Yntema, H.G., Kroes,

13  T., Vulto-van Silfhout, A.T., Koolen, D.A., de Vries, P., Gilissen, C., et al. (2012).

14  Diagnostic exome sequencing in persons with severe intellectual disability. N. Engl. J.

15  Med. *367*, 1921–1929.

16  48. Staples, J., Witherspoon, D.J., Jorde, L.B., Nickerson, D.A., Below, J.E., Huff, C.D.,

17  the University of Washington Center for Mendelian Genomics1 (2016). PADRE:

18  Pedigree-Aware Distant-Relationship Estimation. The American Journal of Human

19  Genetics *99*, 154–162.

20  49. Maxwell, E.K., Packer, J.S., O'Dushlaine, C., McCarthy, S.E., Hare-Harris, A.,

1    Staples, J., Gonzaga-Jauregui, C., Fetterolf, S.N., Faucett, W.A., Leader, J.B., et al.

2    (2017). Profiling copy number variation and disease associations from 50,726

3    DiscovEHR Study exomes.

4    50. Saleheen, D., Natarajan, P., Armean, I.M., Zhao, W., Rasheed, A., Khetarpal, S.A.,

5    Won, H.-H., Karczewski, K.J., O'Donnell-Luria, A.H., Samocha, K.E., et al. (2017).

6    Human knockouts and phenotypic analysis in a cohort with a high rate of consanguinity.

7    Nature Publishing Group *544*, 235–239.

8    51. Narasimhan, V.M., Hunt, K.A., Mason, D., Baker, C.L., Karczewski, K.J., Barnes,

9    M.R., Barnett, A.H., Bates, C., Bellary, S., Bockett, N.A., et al. (2016). Health and

10   population effects of rare gene knockouts in adult humans with related parents. Science

11   *352*, 474–477.

12   52. Perdigoto, C. (2017). Mutations: Dawn of the Human Knockout Project. Nat. Rev.

13   Genet. *18*, 328–329.

14   53. Deciphering Developmental Disorders Study (2017). Prevalence and architecture of

15   de novo mutations in developmental disorders. Nature Publishing Group *542*, 433–438.

16   54. Fromer, M., Pocklington, A.J., Kavanagh, D.H., Williams, H.J., Dwyer, S., Gormley,

17   P., Georgieva, L., Rees, E., Palta, P., Ruderfer, D.M., et al. (2014). De novo mutations in

18   schizophrenia implicate synaptic networks. Nature Publishing Group *506*, 179–184.

19   55. Kong, A., Frigge, M.L., Masson, G., Besenbacher, S., Sulem, P., Magnusson, G.,

20   Gudjonsson, S.A., Sigurdsson, A., Jonasdottir, A., Jonasdottir, A., et al. (2012). Rate of

21   de novo mutations and the importance of father's age to disease risk. Nature Publishing

1    Group *488*, 471–475.

2    56. Rahbari, R., Wuster, A., Lindsay, S.J., Hardwick, R.J., Alexandrov, L.B., Turki, Al,

3    S., Dominiczak, A., Morris, A., Porteous, D., Smith, B., et al. (2016). Timing, rates and

4    spectra of human germline mutation. Nat Genet *48*, 126–133.

5    57. Wong, W.S.W., Solomon, B.D., Bodian, D.L., Kothiyal, P., Eley, G., Huddleston,

6    K.C., Baker, R., Thach, D.C., Iyer, R.K., Vockley, J.G., et al. (2016). New observations

7    on maternal age effect on germline de novo mutations. Nat Commun *7*, 10486.

8    58. McKusick, V.A., HOSTETLER, J.A., and EGELAND, J.A. (1964). GENETIC

9    STUDIES OF THE AMISH, BACKGROUND AND POTENTIALITIES. Bull Johns

10   Hopkins Hosp *115*, 203–222.

11   59. Ober, C., Abney, M., and McPeek, M.S. (2001). The genetic dissection of complex

12   traits in a founder population. The American Journal of Human Genetics *69*, 1068–1079.

13   60. Gusev, A., Palamara, P.F., Aponte, G., Zhuang, Z., Darvasi, A., Gregersen, P., and

14   Pe'er, I. (2012). The architecture of long-range haplotypes shared within and across

15   populations. Mol. Biol. Evol. *29*, 473–486.

16   61. Bycroft, C., Freeman, C., Petkova, D., Band, G., Elliott, L.T., Sharp, K., Motyer, A.,

17   Vukcevic, D., Delaneau, O., O'Connell, J., et al. (2017). Genome-wide genetic data on

18   ~500,000 UK Biobank participants.

19   62. Malinowski, J., Goodloe, R., Brown-Gentry, K., and Crawford, D.C. (2015). Cryptic

20   relatedness in epidemiologic collections accessed for genetic association studies:

1    experiences from the Epidemiologic Architecture for Genes Linked to Environment

2    (EAGLE) study and the National Health and Nutrition Examination Surveys (NHANES).

3    Front Genet *6*, 317.


4