

1 **A platform for case-control matching enables association studies without genotype sharing**

2

3 Mykyta Artomov<sup>1,2</sup>, Alexander A. Loboda<sup>2,3</sup>, Maxim N. Artyomov<sup>4,\*</sup>, Mark J. Daly<sup>1,2,5,\*</sup>

4

5 <sup>1</sup> Analytic and Translational Genetics Unit, Massachusetts General Hospital, Boston, MA, USA

6 <sup>2</sup> Broad Institute, Cambridge, MA, USA

7 <sup>3</sup> ITMO University, St. Petersburg, Russia

8 <sup>4</sup> Department of Immunology and Pathology, Washington University in St. Louis, St. Louis, MO, USA

9 <sup>5</sup> Institute for Molecular Medicine Finland, Helsinki, Finland

10 \* Correspondence: [mjdaly@atgu.mgh.harvard.edu](mailto:mjdaly@atgu.mgh.harvard.edu), [martyomov@wustl.edu](mailto:martyomov@wustl.edu)

11

12 Authors declare no competing financial interests.

13

## 1 **Abstract**

2

3           Acquiring a sufficiently powered cohort of control samples can be time consuming or,  
4 sometimes, impossible. Accordingly, an ability to leverage control samples that were already  
5 collected and sequenced elsewhere could dramatically improve power in all genetic association  
6 studies. However, since majority of the genotyped and sequenced human DNA samples to date  
7 are subject to strict data sharing regulations, large-scale sharing of, in particular, control samples  
8 is extremely challenging. Using insights from image recognition, we developed a method allowing  
9 selection of the best-matching controls in an external pool of samples that is compliant with  
10 personal genotype data protection restrictions. Our approach uses singular value decomposition  
11 of the matrix of case genotypes to rank controls in another study by similarity to cases. We  
12 demonstrate that this recovers an accurate case-control association analysis for both ultra-rare  
13 and common variants and implement and provide online access to a library of ~17,000 controls  
14 that enables association studies for case cohorts lacking control subjects.

15

## 16 **Introduction**

17

18           Traditionally, genetic association studies require construction of a dataset consisting of  
19 both case and control genotypes. With this, one tries to eliminate all potential biases (technical  
20 or ancestral) between case and control cohorts, ensuring that discovered associations are  
21 phenotype-driven. While technical biases could be significantly diminished without explicit  
22 sharing of sensitive individual-level data by using the same sequencing technology and data

1 processing standards for case and control cohorts, adjusting for population stratification requires  
2 researchers to either spend funds for sequencing study-specific controls (reducing the size of the  
3 case cohort) or file extensive paperwork to obtain what are sometimes a limited number of  
4 external control samples (roughly matched on ancestry and technology and consented for  
5 sharing) from public databases like dbGAP<sup>1</sup>. Moreover, re-using publicly available genetic data  
6 often requires each research team to process the data independently, performing redundant  
7 technical routines. The theoretical possibility of utilizing ancestry landscape for association  
8 studies without sharing genotypes was recently explored by UNICORN project<sup>2</sup>, proposing  
9 rigorous definition of a subject's ancestry in a numerical way, but practical implementation of  
10 such concepts is lacking. Similarly, Guo *et al*<sup>3</sup> recently proposed a methodology for usage of  
11 databases like gnomAD or ExAC as controls through calibration of association test statistics using  
12 benign variants, yet this approach is limited in ability to select matching control subjects for an  
13 association study.

14 In this work, we used insights from image recognition algorithms<sup>4</sup> and developed a  
15 rigorous methodology for matching background variation in independent datasets without  
16 explicit genotype sharing. Specifically, we used singular value decomposition (**Supplementary**  
17 **Figure S1**) which expresses original matrix of genotypes as multiplication of three matrices:  $G =$   
18  $USV^T$ ; where  $U$  – is an orthogonal basis of principal directions (with the 1<sup>st</sup> direction chosen along  
19 the largest variance in the original data,  $S$  – a diagonal matrix of singular values and  $V$  – is a matrix  
20 of principal components that provides coordinates of original samples in the basis of matrix  $U$ .  
21 Importantly, for a given genotype matrix  $G$  of cohort of cases, left singular vectors matrix -  $U$ , can  
22 be used alone to reconstruct original matrix of genotypes, and therefore can be shared

1 unrestrictedly with remote database/control server since it does not contain any information on  
2 the individual level data, rather only providing information about variance directions in case  
3 cohort as a whole. We show that given several first vectors of matrix  $U$  derived from the case  
4 cohort data, one can project control genotypes on such case-derived basis and subsequently rank  
5 controls by their similarity to case cohort. Accordingly, optimal control cohort can be chosen by  
6 taking maximum number of controls that is still able to maintain null distribution of the  
7 association test statistic<sup>5</sup>.

8         We performed tests both on genotyping array data and large-scale exome sequencing  
9 dataset confirming successful control selection. Additionally, we matched same exome samples  
10 from 1000 Genomes project processed through different variant discovery pipelines in two  
11 different ways to a set of controls showing that given exome capture concordance joint  
12 processing is not required for returning sufficient set of matched controls. Finally, we illustrate  
13 method performance with rare-variant association study in early onset breast cancer cohort,  
14 performed without genotype sharing. We use the proposed approach to make freely accessible  
15 a set of 16,946 controls that can be used by researchers worldwide without sharing data  
16 restrictions. Specifically, we provide an R-package for case cohort data transformation  
17 (SVDFunctions) at the local machine and accompanying online portal called SVD-based Control  
18 Repository (SCoRe, [www.dnascore.net](http://www.dnascore.net)), which selects a set of optimal controls and outputs  
19 corresponding control summary allele counts to enable association studies.

20

## 21 **Results**

22

## 1 Framework for control matching without genotype sharing

2 Following the standard workflow for genetic association studies – only autosomal LD-  
3 pruned DNA variants should be used for ancestral matching<sup>6,7</sup>. This ensures that downstream  
4 analysis reflects global LD structure rather than local LD. Case genotypes for common autosomal,  
5 LD-pruned variants could be represented as a genotype matrix (with entries “0”, “1”, “2” -  
6 corresponding to the number of alternative alleles in a genotype). With singular value  
7 decomposition (SVD) one can obtain an orthogonal basis of principal directions in cases ( $U$  matrix,  
8 **Figure 1**). Coordinates of principal directions (left singular vectors), especially a limited number  
9 of them, on their own cannot be used to recover individual level genotypes. They represent  
10 orthogonal directions within original data with the largest variance and do not contain any  
11 information of where original samples are located in this space. For reconstruction of the original  
12 matrix one needs to have some information about all three pieces of decomposition result<sup>8</sup>. From  
13 the matrix approximation properties of the SVD it is known that the first singular vector  
14 represents “dominating” direction of the data matrix. Thus, if folded back to case genotypes,  
15 vectors  $u_i$  should be similar in global LD structure to the case cohort.

16 It is established that these methods perform well if the following conditions are met: each  
17 sample could be well characterized by a few of the first singular vectors; an expansion in terms  
18 of the first few singular vectors discriminates well between the sample ancestries<sup>4</sup>; cases  
19 represent a relatively homogenous ancestral cluster (that is, if you are conducting a meta-analysis  
20 across two diverse ancestries, you should run the two case groups separately). Therefore, we can  
21 compute how well a prospective control sample can be represented in the basis of case cohort.  
22 This can be done by computing residual vector in the least squares problem of the type:

1

$$\min_{\alpha_i} \left\| z - \sum_{i=1}^k \alpha_i u_i \right\|$$

2 where  $z$  – a prospective control sample genotype vector,  $u_i$  – left singular vectors of the case  
3 genotype matrix SVD,  $k$  – first few singular vectors. Norm of the residual vector in this case could  
4 be computed as  $\|(I - UU^T)z\|_2$ . SVD requires access to raw genotypes so must be completed  
5 locally and limited basis of the left singular vectors  $\{u_k\}$  could be unrestrictedly shared. External  
6 samples that allow general research use could be stored on the remote control server (**Figure 1**)  
7 accepting anonymous SVD-processed data from local case clients. Remote control server then  
8 estimates residual norm for every prospective control and creates ranking for the representation  
9 quality in the basis of case cohort. Case allele frequencies need to be supplied to the control  
10 server for defining appropriate size of control set. For every control residual vector norm  
11 threshold an association test should be performed using allele frequencies for DNA variants used  
12 for matching and genomic inflation factor<sup>5</sup> estimated. Largest set of controls delivering  
13 acceptably null distribution of the test statistic ( $\lambda_{GC} \leq 1.05$ ) should be then taken as a matched set  
14 of controls. Allele frequencies for variants of interest should then be computed in matched  
15 control dataset and could be shared with client to run full-scale association test (Toy example of  
16 case-control study with 2 cases and 2 controls is explained at **Supplementary Figure S2**). Unlike  
17 individual level data, such summary statistic sharing from most consented resources is routinely  
18 allowed<sup>9,10</sup>.

19

20 **1000 genomes. OMNI array genotyping data**

1 To validate our approach, we used 1000 genomes genotyping data. Entire dataset was subjected  
2 to per genotype and per individual quality check procedures (**Supplementary Figure S3**). Final  
3 dataset consisted of 10,000 LD-pruned autosomal variants and 1,708 non-related samples (**Figure**  
4 **2A**). To simulate an association study we divided the dataset into 100 “cases” of European  
5 ancestry (as provided by 1000 genomes project annotation) and “control pool” of 1608 samples  
6 that included 359 European samples (**Figure 2A,B**). We simulated an association study without  
7 genotype sharing by separating every case cohort from control candidates (**Figure 2C**). For a case  
8 group SVD was computed and top 5 left singular vectors were used for control matching  
9 purposes. Residual vector norm was computed for every control, ranking controls by similarity to  
10 case cohort (**Figure 2D**) and separating ancestries (**Figure 2E**). Choosing different residual norm  
11 thresholds allows to create different pools of controls of variable matching quality (**Figure 2D,E**).  
12 Indeed, PCA plot built using shared genotypes confirmed that increase in residual vector norm  
13 threshold results in departure from original European case cluster delivering poor control  
14 matching quality (**Figure 2F**). To select optimal threshold value, allele frequency of variants in  
15 cases is used to estimate association test statistics and genomic control factor for different  
16 residual vector norm thresholds (**Figure 2G,H**). Largest control pool size with null distributed  
17 association test statistics (e.g.  $\lambda_{GC} \leq 1.05$ ) was selected as optimally matched set (**Figure 2I**). 100  
18 simulation instances with random 100 “case” cohorts (100 European samples each) yielded on  
19 average selection of 428 controls (SD=26.6) with estimated mean genomic inflation  $\lambda_{GC}=1.049$   
20 (SD= $8.76 \times 10^{-4}$ ) (**Figure 2J**). These control sets covered nearly all European samples in Control  
21 Pool and also selected a small number of controls annotated in 1000 Genomes as Latin-  
22 Americans (though still keeping genomic inflation factor below 1.05). In 100 random selections

1 of “case” cohort (100 European samples each), we identified a group of Latin-Americans  
2 recurrently selected for a matched control set (**Supplementary Figure S4A**). It appears to be the  
3 closest to Europeans on PCA (**Supplementary Figure S4B**). Selection of samples with non-  
4 matching ancestry annotation is bound to the statistical power of the association testing used for  
5 genomic inflation estimation, thus, the more cases are used for association testing the fewer  
6 controls of non-matching ancestry annotation are selected (**Supplementary Figure S4C**). Though,  
7 with any set of case cohort genomic inflation factor is below 1.05 implying neglectable effect on  
8 further phenotypic association test results (**Supplementary Figure S4D**).

9 Furthermore, such SVD-based approach can be used for creating an easy-to-use ancestry  
10 predictor that yield large confidence results (**Supplementary Methods, Supplementary Figure**  
11 **S5**, SVDFunctions R-package available at <https://github.com/alexloboda/SVDFunctions>).

12

### 13 **Exome sequencing data**

14 Further, we sought to test how this approach performs in exome sequencing data with  
15 relatively small number of LD-pruned DNA variants. We used an aggregated set of whole exome  
16 sequences consented for joint variant calling resulting in 37,607 samples to build a test dataset  
17 (**Supplementary Table S1**). Further, a subset of LD-pruned variants used for PCA in ExAC  
18 database<sup>9</sup> was subjected to quality check: only genotypes with DP>10, GQ>20 and variants with  
19 less than 500 missing genotypes were allowed, resulting in 4,561 DNA variant left for analysis.  
20 Samples were further analyzed for relatedness and related samples were removed: in pairs of  
21 samples with  $\hat{\pi} < 0.2$  only one sample was randomly kept (**Supplementary Methods**). Final  
22 dataset consisted of 32,677 whole exome sequences (**Figure 3A**). Similarly, to the experiment



1 with genotyping data that was described above, we randomly selected 3,000 European samples  
2 as a “case” set and leftover 29,677 samples became a disjointed pool of controls (**Figure 3B**). We  
3 used first 5 singular vectors to estimate residual norm for every sample in the control set and  
4 create a control ranking (**Figure 3C**), which distinctly separated ancestral groups (**Figure 3D**).  
5 Increase in residual norm results in departure from targeted European control cluster and  
6 inflation of the association test statistic (**Figure 3E,F**). Optimal control set size was defined as a  
7 largest set of controls with genomic inflation factor less or equal to 1.05 or the smallest inflation  
8 factor if no values below 1.05 are available (**Figure 3G**). In 100 random selections of “case” cohort  
9 (3000 Europeans each) on mean control set size was 17,435 (SD=921.97) with genomic inflation  
10 factor mean  $\lambda_{GC}=1.056$  (SD=0.012) (**Figure 3H**).

11 We analyzed how number of selected for analysis singular vectors affects size of matched  
12 controls cohort. Upon increase in number of singular vectors, more variance could be explained  
13 in limited reconstruction of the original matrix (e.g. if both  $U$ ,  $S$ ,  $V$  matrices are available for  
14 original genotype matrix, its reconstruction using first several vectors of  $U$  matrix would capture  
15 the variance better if more vectors are used, **Supplementary Figure S6A,B**). Though, since in our  
16 algorithm no matrix reconstruction is performed, there is no significant dependence of matched  
17 controls cohort size on the size of the case-derived basis (**Supplementary Figure S6C,D**). Thus, it  
18 is reasonable to use only first few singular vectors to increase computation speed.

19

## 20 **Controlling for technical bias**

21 In the above examples, genotyping or sequencing data was processed in the uniform way  
22 and in case of exomes – variants were called jointly, an action which is impossible without explicit

1 genotype sharing. Thus, we looked into effects of technical biases on the remote case-control  
2 matching. We used publicly available 1000 genomes Phase 3 data<sup>10</sup> to select 46 samples (CEU  
3 ancestry), that were also present in our exome control pool dataset and attempted to match  
4 them to a control group. For this analysis control pool dataset was modified to exclude all 1000  
5 genomes samples. Samples with known cancer phenotype were also excluded for further use of  
6 this data as a public database of controls (**Supplementary Figure S7A**). Both “external” and jointly  
7 processed 1000 genomes samples were acceptably matched to a group of controls  
8 (**Supplementary Figure S7B-D**), with exactly the same 5963 control candidates selected for both  
9 case groups and 1 additional control matched to “internal” cases. Important to note, that  
10 external 1000 genomes exomes were sequenced using similar to control pool’s Agilent exome  
11 capture kit. Therefore, we conclude that separate variant calling does not introduce major  
12 technical biases within the proposed approach.

13         One the level of experimental design, difference in DNA sequence capture used for exome  
14 sequencing are known to introduce major effects in variant calling and cannot be cross-used for  
15 association tests even if variants are jointly called<sup>9,11</sup>. The proposed approach will automatically  
16 treat platform-based bias as a kind of “ancestry”-matching and control candidates from the same  
17 sequencing platform would be prioritized over other options. However, in this initial release of  
18 17,000 we only provide the control samples sequenced with Whole Exome Agilent 1.1 RefSeq  
19 plus 3 boosters capture. Being the most represented capture in ExAC (~77%)<sup>9</sup> covers a major part  
20 of sequenced samples to date.

21         Additional data quality metrics: depth, missing data rate could be directly shared between  
22 control server and case client. Even with direct sharing of genotypes further work needs to be

1 done to get a gold-standard association study: lining up depth and accuracy at every site, gene,  
2 exon. While control platform described here will ensure ancestral and platform matching it is up  
3 to a researcher to match quality metrics and conduct careful statistical analysis.

4

## 5 **Breast Cancer Cohort Analysis**

6 We implemented presented methodology in combination of R-package (SVDFunctions)  
7 for data preprocessing and online SVD-based Control Repository (SCoRe) platform  
8 ([www.dnascore.net](http://www.dnascore.net)) and performed a gene-based association study for a cohort of early onset  
9 breast cancer patients (dbGAP: phs000822.v1.p1): 291 non-related cases matching quality  
10 standards were further used for analysis. Genotype matrix and summary genotype counts of  
11 cases were constructed for 4,561 LD-pruned DNA variants that are available for matching through  
12 SCoRe. 1,930 controls were matched to the case cohort with  $\lambda_{GC}=1.07$  (**Figure 4A**). To ensure that  
13 not only the targeted DNA variants are matched we inquired SCoRe to return summary allele  
14 counts for various variant classes found in cases: common synonymous variants (that were not  
15 used for matching), rare synonymous variants grouped by gene (singletons in cases, implying  
16  $MAF < 0.00172$ ). Null distribution of the test statistic in both inquiries confirmed case-control  
17 matching in both common and rare variant background variation (**Figure 4D,E**). Further, rare  
18 ( $MAF < 0.00172$ ) protein-truncating variants (PTV) gene burden was examined. 1930 genes had at  
19 least 1 PTV variant in cases (Bonferroni corrected significance threshold  $0.05/1930 = 2.56 \times 10^{-5}$ ).  
20 These genes were submitted to SCoRe to obtain summary gene-based counts of PTVs in matched  
21 controls and further, association test was performed locally using Fisher's exact test (**Figure 4F**)

1 re-discovering *BRCA1* as a susceptibility gene (with 12 carriers in 291 case and 7 carriers in 1930  
2 controls; Fisher test  $P=6.36 \times 10^{-7}$ ).

3 Strikingly, despite the complex ancestral structure of case cohort with substantial amount  
4 of Ashkenazi samples (**Figure 4B**), we identified a set of controls that were matched as a single  
5 batch with  $\lambda_{GC}=1.07$  (**Figure 4C**). This was very different when compared to conventional  
6 genotype sharing association study: due to complex ancestral structure in cases, matching of  
7 controls as a single batch is largely impeded. We used K-means clustering to identify matching  
8 batches (**Supplementary Figure S8A, B**) so that all of them return null-distributed test statistic  
9 (using same variants as were used for SCoRe-based matching, **Supplementary Figure S8C-F**).  
10 Using batch-based matching with fully shared genotypes, we identified a total of 6,415 controls  
11 (out of 16,946 candidate controls in SCoRe) to the same case group in jointly called dataset.

12 Statistical power was estimated for Fisher's exact test using multiple odds ratios for  
13 SCORE and conventional genotype sharing tests (**Figure 4G**), implying 291 case cohort and single-  
14 batch matched control cohorts of 1930 and 6415 samples, respectively (with disease prevalence  
15 12.4%<sup>12</sup>). We also estimated power for separate batches in a non-stratified analysis set up (**Figure**  
16 **4H**). The results show that SCoRe performs significantly better in a single-batch analysis than a  
17 shared genotype study and is producing sufficient number of controls to effectively saturate a  
18 power of the study in terms of size of control cohort.

19

## 20 **Web-portal for Case-Control Matching**

21 We established a public online platform with 16,946 non-cancer controls available for  
22 matching (**Supplementary Table S2**). We provide a list of variants available for matching in

1 control pool, and a R-package(SVDFunctions) for local generating of the case genotype matrix  
2 from standard PLINK<sup>13</sup> binary file format and a code to generate SVD and allele counts data that  
3 could be uploaded to the matching portal. Available at [www.dnascore.net](http://www.dnascore.net).

4 We set a hard threshold on minimal number of matched controls (500 samples) to deliver  
5 allele frequencies to a user and controls are selected in batches of 10. This ensures only  
6 anonymous summary allele frequencies are disclosed from control cohort. Quires to the  
7 database could be asking for variant based frequencies or cumulative counts of rare variants per  
8 gene. For variant-based quires we limit return to common variants (MAF>1%) to avoid disclosing  
9 unique DNA variants that would readily disclose a genotype.

10

## 11 **Discussion**

12 Statistical power is a key to successful association study and control set size is often a  
13 limiting factor. Despite potential availability of control sets through public repositories, large  
14 efforts should be put into processing case and control datasets jointly before even preliminary  
15 results of an association study could emerge. Practically, this often becomes infeasible for the  
16 small cohort studies limited to data access or computational power. Assembly of large case-  
17 control datasets is generally done by international consortia (ExAC, PGC, IBD Genetics  
18 Consortium, etc.) as this requires a lot of effort and generous funding. We provide a large pool  
19 of exome sequences and a tool enabling rapid selection of matched control sets without  
20 genotype sharing that ultimately provides allele frequency statistics required for performing  
21 association tests. Nearly no effort is required from the user side to get all information needed for  
22 association study, facilitating future discovery of associated genes and DNA variants.

1           Local cohorts assembled at hospitals as a part of clinical screening procedures or genetic  
2 counselling often have very modestly sized or none control sets and stringent sharing regulation.  
3 Our platform enables case-control study design and boosts statistical power for such patient  
4 cohorts. Especially, for rare mendelian phenotypes where assembly of well-powered case-control  
5 cohort is impeded by low disease prevalence.

6           Finally, hundreds of thousands samples were subjected to exome or genome sequencing  
7 to date in the world. However, all this data exists in isolated pieces limiting potential benefit for  
8 genetic studies. We provide a repository of the software codes used for running the matching  
9 platform so that it could readily be implemented by large data holders – National biobank  
10 initiatives and international disease consortia to let community benefit from large-scale genetic  
11 resources. This is also critically important for advancing genetic association studies in situations  
12 when explicit data sharing is not permitted or very challenging in the international settings, thus  
13 potentially providing insights into rare sample collections that were not available so far. Finally,  
14 approach developed in this work creates a path to creating unified central repository that would  
15 encompass all studies published in dbGAP and make it accessible to association studies run in  
16 any design and cohort.

17

## 1 **Methods**

2

### 3 **1000 genomes OMNI-array genotyping data**

4 Raw genotypes were downloaded from 1000 genomes FTP site at:

5 [ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/supporting/hd\\_genotype\\_chip/](ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/supporting/hd_genotype_chip/)

6 Dataset was subjected to variant and samples quality check, according to the standard

7 genotyping data handling protocol (**Supplementary Figure S2**)<sup>6,7</sup>. Quality check was performed

8 with PLINK package<sup>13</sup>. Code for converting binary PLINK format files into genotype matrix used

9 for matching is available through a control matching web-portal at

10

### 11 **Exome Sequencing Data**

12 Whole exome libraries were prepared using Whole Exome Agilent 1.1 RefSeq plus 3 boosters

13 capture kit and protocol, automated on the Agilent Bravo and Hamilton Starlet. Libraries were

14 then prepared for sequencing using a modified version of the manufacturer's suggested

15 protocol, automated on the Agilent Bravo and Hamilton Starlet, followed by sequencing on the

16 Illumina HiSeq 2000. We used an aggregated set of samples consented for joint variant calling

17 resulting in 37,607 samples (**Supplementary Table S1**). All samples were sequenced using the

18 same capture reagents at the Broad Institute and aligned on the reference genome with BWA<sup>14</sup>

19 and the best-practices GATK/Picard Pipeline, followed by joint variant calling with all samples

20 processed as a single batch using GATK v 3.1-144 Haplotype Caller<sup>15-17</sup>. The resulting dataset

21 had 7,094,027 distinct variants. Variant effect predictor was used for variant annotation<sup>18</sup>.

1 Early onset breast cancer cohort used for association study is available through dbGAP  
2 (phs000822.v1.p1).

3 1000 genomes sequencing data was downloaded from

4 <ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/>

5

## 6 **Software**

7 Software pipeline was developed using R-language and libraries<sup>19–21</sup>. PLINK/SEQ library was  
8 used for operations with sequencing data<sup>22</sup>.

9 Source code, usage examples and manual are available at supplemental materials and online:

10 [www.dnascore.net](http://www.dnascore.net) and <https://github.com/alexloboda/SVDFunctions>.

11

12



1 **Acknowledgements**

2 Authors would like to thank Konstantin Zaitsev (Washington University in St. Louis) for helpful  
3 advice on singular value decomposition methodology. M.A. would like to thank Dr. Andrey  
4 Shaw (Genentech) for inspiration to pursue research in human genetics.

5

6 **Author Contributions**

7 Conceptualization: M.A., M.N.A, M.J.D. Investigation: M.A., A.A.L., M.N.A., M.J.D. Software:  
8 M.A., A.A.L. Writing Original Draft: M.A., M.N.A., M.J.D.

9

10 **Competing Interests**

11 Authors declare no competing interests.

## 1 References

- 2 1. Mailman, M. D. *et al.* The NCBI dbGaP database of genotypes and phenotypes. *Nat.*  
3 *Genet.* **39**, 1181–1186 (2007).
- 4 2. Bodea, C. A. *et al.* A Method to Exploit the Structure of Genetic Ancestry Space to  
5 Enhance Case-Control Studies. *Am. J. Hum. Genet.* **98**, 857–868 (2016).
- 6 3. Guo, M. H., Plummer, L., Chan, Y.-M., Hirschhorn, J. N. & Lippincott, M. F. Burden Testing  
7 of Rare Variants Identified through Exome Sequencing via Publicly Available Control Data.  
8 *Am. J. Hum. Genet.* **103**, 522–534 (2018).
- 9 4. Elden, L. *Matrix methods in data mining and pattern recognition. Society of Industrial and*  
10 *Applied Mathematics.* (2007).
- 11 5. Devlin, B. & Roeder, K. Genomic control for association studies. *Biometrics* **55**, 997–1004  
12 (1999).
- 13 6. Anderson, C. A. *et al.* Data quality control in genetic case-control association studies. *Nat.*  
14 *Protoc.* **5**, 1564–1573 (2010).
- 15 7. Reed, E. *et al.* A guide to genome-wide association analysis and post-analytic  
16 interrogation. *Stat. Med.* **34**, 3769–3792 (2015).
- 17 8. Abdi, H. Singular Value Decomposition (SVD) and Generalized Singular Value  
18 Decomposition (GSVD). in *Encyclopedia of Measurement and Statistics* 907–912  
19 (Thousand Oaks (CA): Sage, 2007).
- 20 9. Lek, M. *et al.* Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536**,  
21 285–291 (2016).
- 22 10. Consortium, T. 1000 G. P. A global reference for human genetic variation. *Nature* **526**,

- 1           68–74 (2015).
- 2   11.   Shigemizu, D. *et al.* Performance comparison of four commercial human whole-exome  
3           capture platforms. *Sci. Rep.* **5**, 12742 (2015).
- 4   12.   *Cancer Statistics Facts. National Cancer Institute.* (2018).
- 5   13.   Purcell, S. *et al.* PLINK: A Tool Set for Whole-Genome Association and Population-Based  
6           Linkage Analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).
- 7   14.   Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler  
8           transform. *Bioinformatics* **25**, 1754–1760 (2009).
- 9   15.   DePristo, M. A. *et al.* A framework for variation discovery and genotyping using next-  
10          generation DNA sequencing data. *Nat. Genet.* **43**, 491–8 (2011).
- 11   16.   McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for analyzing  
12          next-generation DNA sequencing data. *Genome Res.* **20**, 1297–303 (2010).
- 13   17.   Van der Auwera, G. A. *et al.* From FastQ data to high confidence variant calls: the  
14          Genome Analysis Toolkit best practices pipeline. *Curr. Protoc. Bioinforma.* **43**, 11.10.1-33  
15          (2013).
- 16   18.   McLaren, W. *et al.* The Ensembl Variant Effect Predictor. *Genome Biol.* **17**, 122 (2016).
- 17   19.   Team, R. C. A language and environment for statistical computing. (2013).
- 18   20.   Wickham, H. *Elegant Graphics for Data Analysis.* (Springer-Verlag, 2016).
- 19   21.   Clayton, D. snpStats: SnpMatrix and XSnMatrix classes and methods. R package version  
20          1.30.0. (2017).
- 21   22.   <https://atgu.mgh.harvard.edu/plinkseq/>. PLINK/SEQ.

22

23

1 **Figure Legends**

2 **Figure 1. Scheme of an association study without genotype sharing.** Individual level genotype  
3 data is subject to data sharing restrictions. SVD-based processing creates anonymous data  
4 describing variation in case genotypes without storing individual data that could be shared with  
5 no restrictions. Remote server with a pool of controls selects a set of controls genotype  
6 variation matching cases, estimates allele frequency for sites to be used for association study  
7 and delivers results to the user.

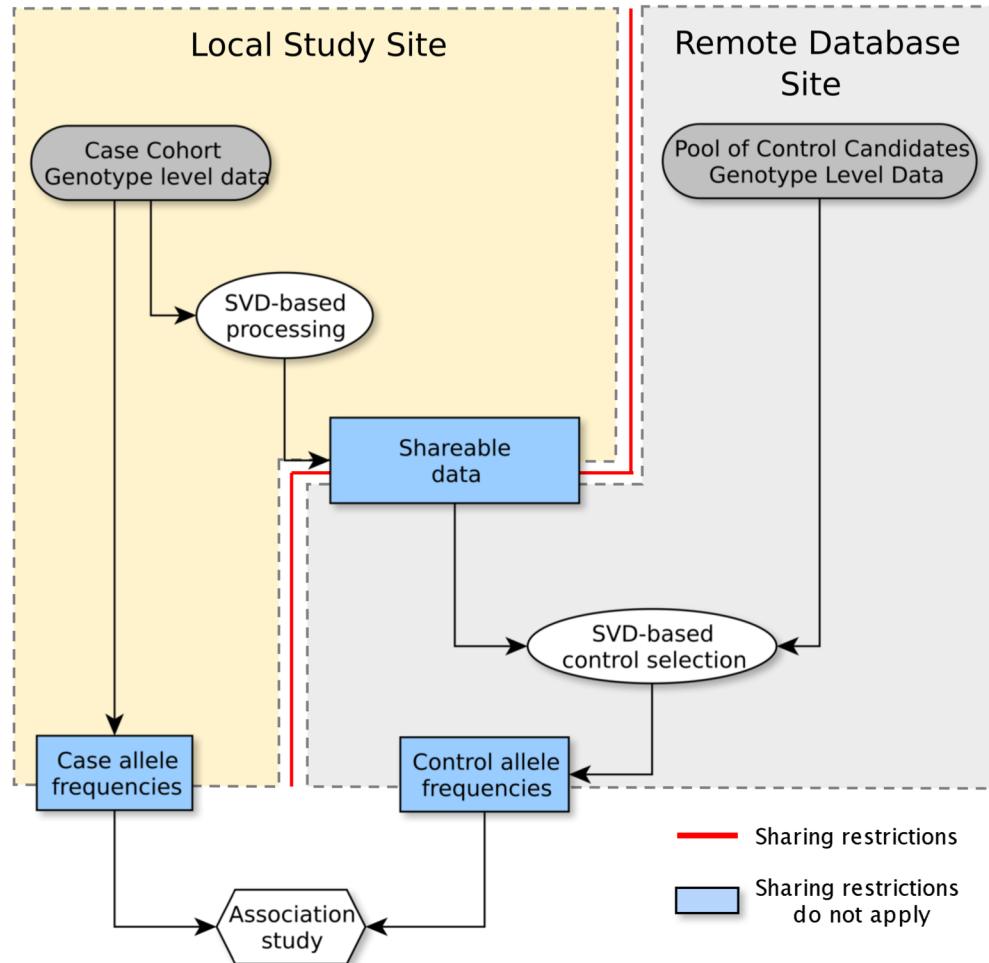
8  
9 **Figure 2. Case Study 1: Simulated Case-Control Study with 1000 Genomes OMNI Genotyping**  
10 **Array Data.** (A) Breakdown of ancestries present in the dataset and case-control study set up:  
11 random 100 European samples are selected as “cases” and the rest of the data is tested as  
12 prospective controls; (B) Conventional PCA showing European samples selected as case group;  
13 (C) Scheme of data handling simulating association study without genotype sharing; (D) Ranking  
14 of controls by quality of representation in case basis; (E) Distribution of residual vector norms  
15 for prospective controls; (F) Conventional PCA shows that greater values of residual vector  
16 norms correspond to greater departure from European cluster; (G) Quantile-Quantile plots for  
17 multiple thresholds demonstrate increasing inflation; (H) Optimal residual vector norm  
18 threshold should deliver  $\lambda_{GC} < 1.05$ ; (I) Optimal threshold selection scheme; (J) Matching  
19 experiment summary results.

20  
21 **Figure 3. Case Study 2: Simulated Case-Control Study with 32,677 Exomes.** (A) Breakdown of  
22 ancestries present in the dataset and case-control study set up: random 3000 European

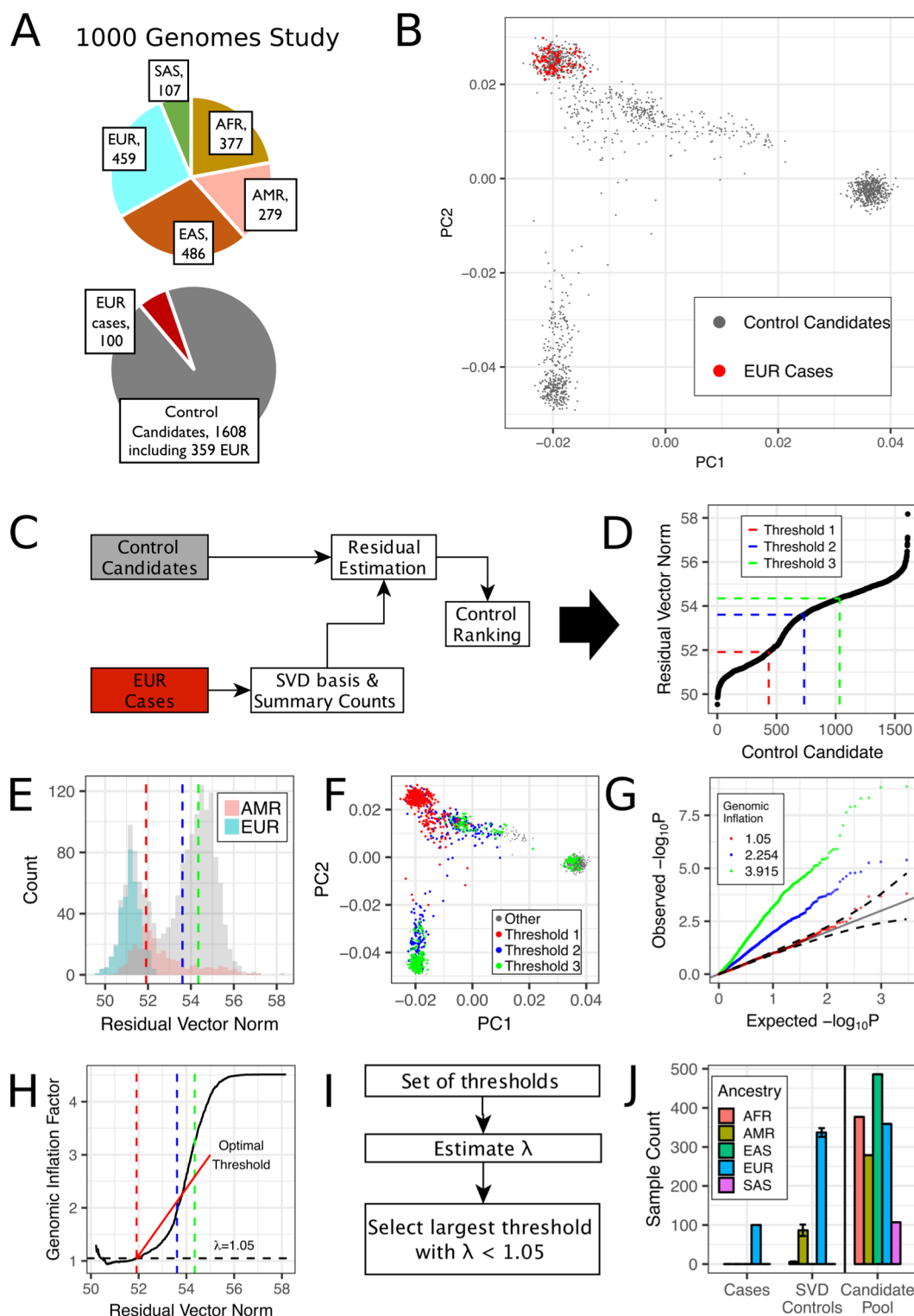
1 samples are selected as “cases” and the rest of the data is tested as prospective controls; (B)  
2 Conventional PCA showing European samples selected as case group; (C) Ranking of controls by  
3 quality of representation in case basis; (D) Distribution of residual vector norms for prospective  
4 controls; (E) Conventional PCA shows that greater values of residual vector norms correspond  
5 to greater departure from European cluster; (F) Quantile-Quantile plots for multiple thresholds  
6 demonstrate increasing inflation; (G) Optimal residual vector norm threshold should deliver  
7  $\lambda_{GC} < 1.05$ ; (H) Matching experiment summary results.

8  
9 **Figure 4. Case Study 3: Early-onset Breast Cancer Association Study.** (A) Experimental  
10 workflow for obtaining matched controls data from SCoRe online platform; (B) Principal  
11 component analysis of breast cancer cohort (EOBC) performed jointly with 1000 genomes  
12 European samples. Multiple European sub-ancestries are present in case cohort, including  
13 substantial Ashkenazi group; (C) QQ-plot obtained from SCoRe platform for 1930 matched  
14 controls; (D) QQ-plot of association study performed on common ( $MAF > 5\%$ ) synonymous  
15 variants found in cases with control allele frequencies obtained from SCoRe; (E) QQ-plot for  
16 gene-burden test for rare synonymous (singletons in cases, corresponds to  $MAF < 0.00172$ ) with  
17 gene-level summary allele counts obtained from SCoRe (F) QQ-plot for rare protein-truncating  
18 variants (singletons in cases, corresponds to  $MAF < 0.00172$ ) with gene-level summary allele  
19 counts obtained from SCoRe. BRCA1 is the top associated gene; (G) Fisher test statistical power  
20 estimate for 291 cases. Control cohorts matched with SCoRe platform and conventional  
21 stratified analysis with genotype sharing (**Supplementary Figure S8**) are labeled with dashed  
22 lines; (H) Fisher test statistical power for each case-control cohort: separate batches from

- 1 shared genotypes analysis, stratified analysis with shared genotypes and SCoRe platform
- 2 **(Supplementary Figure S8).**
- 3

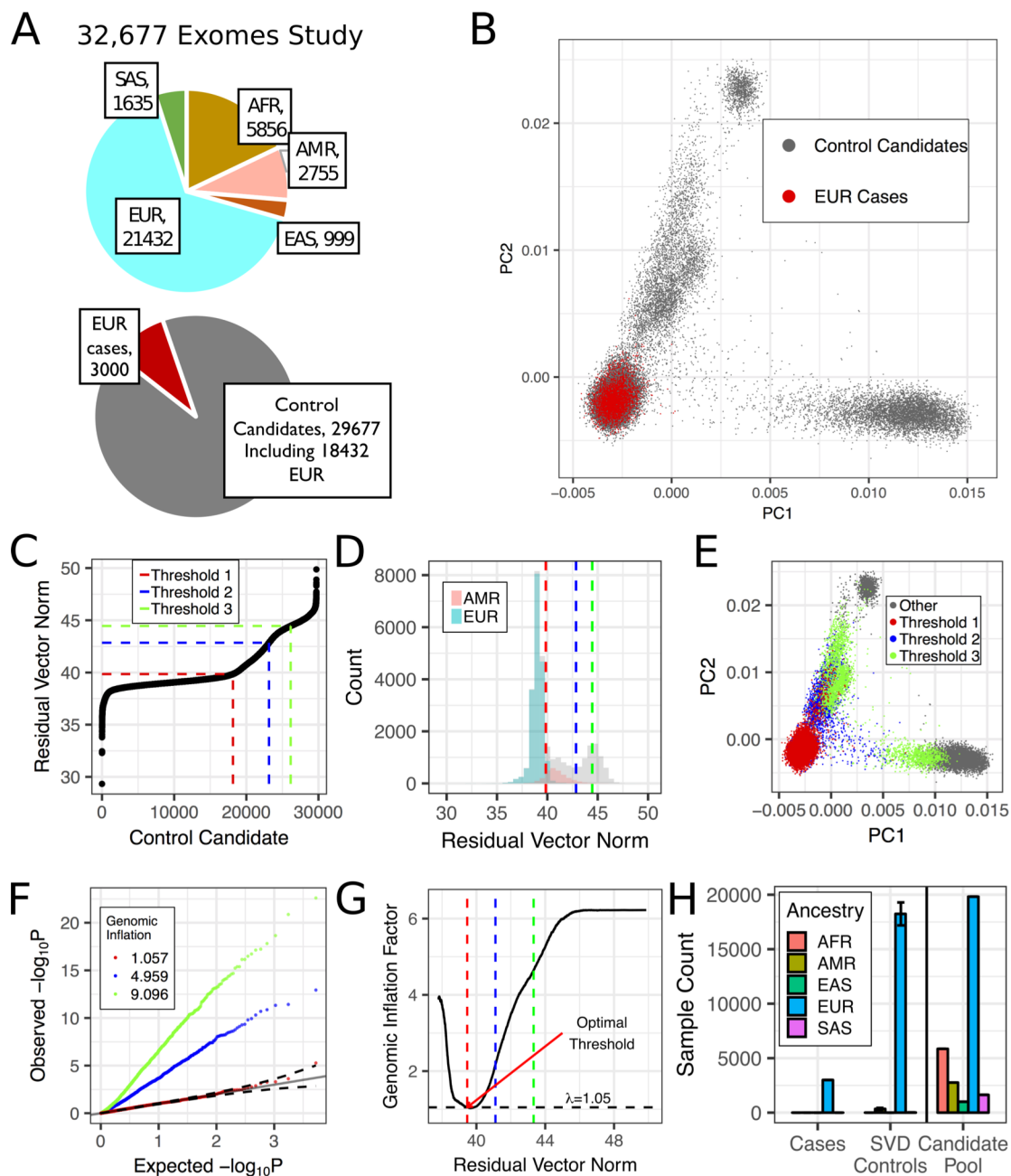


**Fig. 1. Scheme of an association study without genotype sharing.** Individual level genotype data is subject to data sharing restrictions. SVD-based processing creates anonymous data describing variation in case genotypes without storing individual data that could be shared with no restrictions. Remote server with a pool of controls selects a set of controls genotype variation matching cases, estimates allele frequency for sites to be used for association study and delivers results to the user.

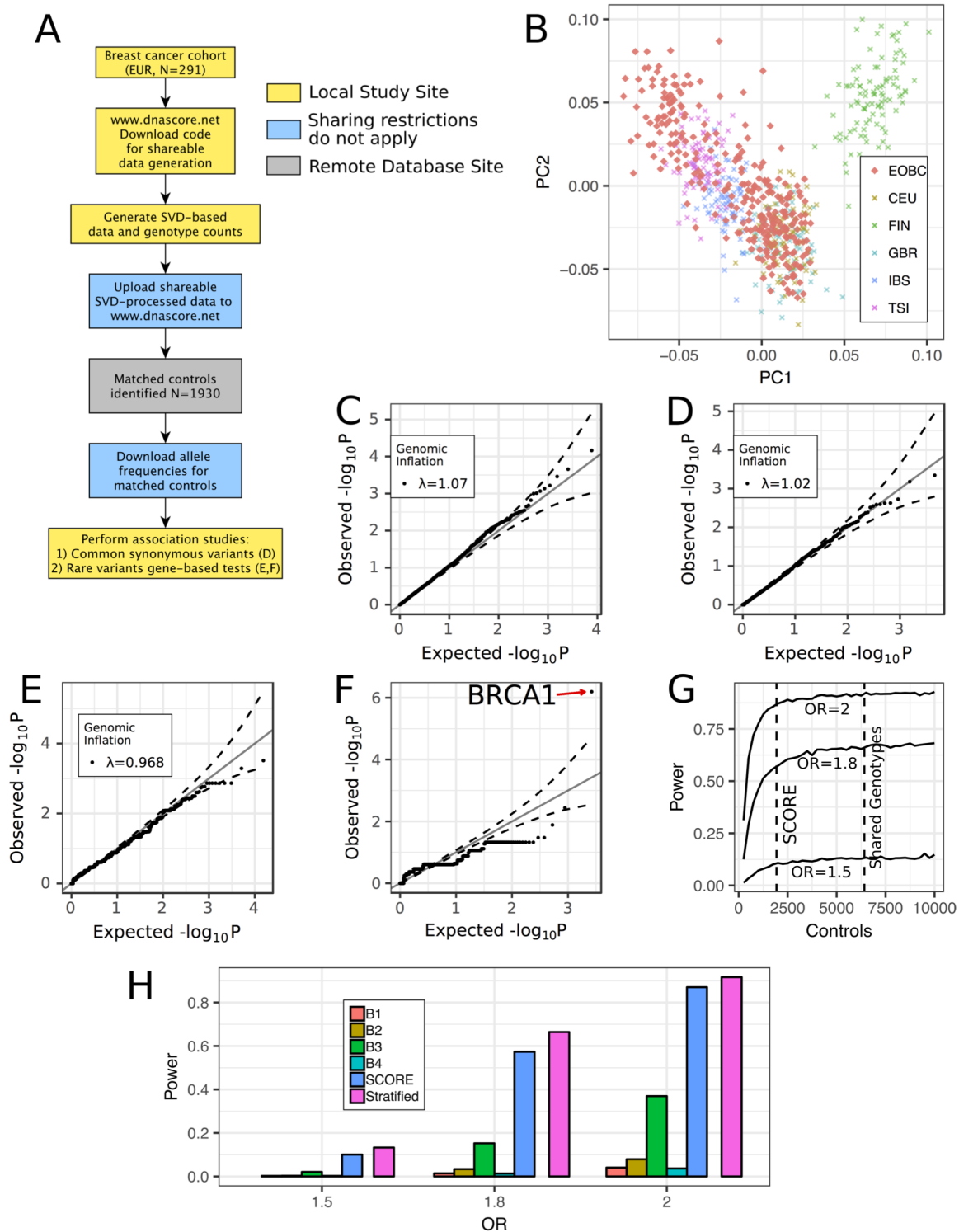


**Fig. 2. Case Study 1: Simulated Case-Control Study with 1000 Genomes OMNI Genotyping Array Data.** (A) Breakdown of ancestries present in the dataset and case-control study set up: random 100 European samples are selected as “cases” and the rest of the data is tested as prospective controls; (B) Conventional PCA showing European samples selected as case group; (C) Scheme of data handling simulating association study without genotype sharing; (D) Ranking of controls by quality of representation in case basis; (E) Distribution of residual vector norms for prospective controls; (F) Conventional PCA shows that greater values of residual vector norms correspond to greater departure from European cluster; (G) Quantile-Quantile plots for multiple thresholds demonstrate increasing inflation; (H) Optimal residual vector norm threshold should deliver  $\lambda_{GC} < 1.05$ ; (I) Optimal threshold selection scheme; (J) Matching experiment summary results





**Fig. 3. Case Study 2: Simulated Case-Control Study with 32,677 Exomes.** (A) Breakdown of ancestries present in the dataset and case-control study set up: random 3000 European samples are selected as “cases” and the rest of the data is tested as prospective controls; (B) Conventional PCA showing European samples selected as case group; (C) Ranking of controls by quality of representation in case basis; (D) Distribution of residual vector norms for prospective controls; (E) Conventional PCA shows that greater values of residual vector norms correspond to greater departure from European cluster; (F) Quantile-Quantile plots for multiple thresholds demonstrate increasing inflation; (G) Optimal residual vector norm threshold should deliver  $\lambda_{GC} < 1.05$ ; (H) Matching experiment summary results



**Fig. 4. Case Study 3: Early-onset Breast Cancer Association Study.** (A) Experimental workflow for obtaining matched controls data from SCORE online platform; (B) Principal component analysis of breast cancer cohort (EOBC) performed jointly with 1000 genomes European samples. Multiple European sub-ancestries are present in case cohort, including substantial Ashkenazi group; (C) QQ-plot obtained from SCORE platform for 1930 matched controls; (D) QQ-plot of association study performed on common (MAF>5%) synonymous variants found in cases with control allele frequencies obtained from SCORE; (E) QQ-plot for gene-burden test for rare synonymous (singletons in cases, corresponds to MAF<0.00172) with gene-level summary allele counts obtained from SCORE (F) QQ-plot for rare protein-truncating variants (singletons in cases, corresponds to MAF<0.00172) with gene-level summary allele counts obtained from SCORE. *BRCA1* is the top associated gene; (G) Fisher test statistical power estimate for 291 cases. Control cohorts matched with SCORE platform and conventional stratified analysis with genotype sharing (Fig. S8) are labeled with dashed lines; (H) Fisher test statistical power for each case-control cohort: separate batches from shared genotypes analysis, stratified analysis with shared genotypes and SCORE platform (Fig. S8).