1  **Identification of putative causal loci in whole-genome sequencing data via knockoff statistics**

2  Zihuai He[1,2#], Linxi Liu[3], Chen Wang[4], Yann Le Guen[1], Justin Lee[2], Stephanie Gogarten[5], Fred Lu[6], Stephen
3  Montgomery[7,8], Hua Tang[6,7], Edwin K. Silverman[9], Michael H. Cho[9], Michael Greicius[1], Iuliana Ionita-Laza[4#]
4
5  [1]Department of Neurology and Neurological Sciences, Stanford University, Stanford, CA 94305, USA
6  [2]Quantitative Sciences Unit, Department of Medicine, Stanford University, Stanford, CA, 94305, USA
7  [3]Department of Statistics, Columbia University, New York, NY 10027, USA
8  [4]Department of Biostatistics, Columbia University, New York, NY 10032, USA
9  [5]Department of Biostatistics, University of Washington, Seattle, WA, USA
10  [6]Department of Statistics, Stanford University, Stanford, CA, 94305, USA
11  [7]Department of Genetics, Stanford University, Stanford, CA, 94305, USA
12  [8]Department of Pathology, Stanford University, Stanford, CA, 94305, USA
13  [9]Channing Division of Network Medicine and Division of Pulmonary and Critical Care Medicine Division, Brigham
14  and Women's Hospital, Harvard Medical School, Boston, MA, 02215, USA
15
16  #Correspondence to: Zihuai He (zihuai@stanford.edu) and Iuliana Ionita-Laza (ii2135@cumc.columbia.edu)
17

18  **Abstract**

19  The analysis of whole-genome sequencing studies is challenging due to the large number of rare variants
20  in noncoding regions and the lack of natural units for testing. We propose a statistical method to detect and
21  localize rare and common risk variants in whole-genome sequencing studies based on a recently developed
22  knockoff framework. It can (1) prioritize causal variants over associations due to linkage disequilibrium
23  thereby improving interpretability; (2) help distinguish the signal due to rare variants from shadow effects
24  of significant common variants nearby; (3) integrate multiple knockoffs for improved power, stability and
25  reproducibility; and (4) flexibly incorporate state-of-the-art and future association tests to achieve the
26  benefits proposed here. In applications to whole-genome sequencing data from the Alzheimer's Disease
27  Sequencing Project (ADSP) and COPDGene samples from NHLBI Trans-Omics for Precision Medicine
28  (TOPMed) Program we show that our method compared with conventional association tests can lead to
29  substantially more discoveries.

30

31  **Introduction**

32  The rapid development of whole-genome sequencing technology allows for a comprehensive
33  characterization of the genetic variation in the human genome in both coding and noncoding regions. The
34  noncoding genome covers ~98% of the human genome, and includes regulatory elements that control when,
35  where, and to what degree genes will be expressed. Understanding the role of noncoding variation could
36  provide important insights into the molecular mechanisms underlying different traits.

37  Despite the increasing availability of whole-genome sequencing datasets including those from moderate to
38  large scale projects such as the Alzheimer's Disease Sequencing Project (ADSP), the Trans-Omics for
39  Precision Medicine (TOPMed) program etc., our ability to analyze and extract useful information from
40  these datasets remains limited at this point and many studies still focus on the coding regions and regions
41  proximal to genes[1,2]. The main challenges for analyzing the noncoding regions include the large number of
42  rare variants, the limited knowledge of their functional effects, and the lack of natural units for testing (such
43  as genes in the coding regions). To date, most studies have relied on association testing methods such as
44  single variant tests for common variants, gene-based tests for rare variants in coding regions, or a heuristic
45  sliding window strategy to apply gene-based tests to rare variants in the noncoding genome[3,4]. Only few
46  methods have been developed to systematically analyze both common and rare variants across the genome,
47  owing to difficulties such as an increased burden of the multiple testing problem, more complex correlations
48  and increased computational cost. Moreover, a common feature of the existing association tests is that they
49  often identify proxy variants that are correlated with the causal ones, rather than the causal variants that

1  directly affect the traits of interest. Identification of putative causal variants usually requires a separate fine-
2  mapping step. Fine-mapping methods such as CAVIAR[5] and SUSIE[6] were developed for single, common
3  variant analysis in GWAS studies, and are not directly applicable to window-based analysis of rare variants
4  in sequencing studies.

5  Methods that control the family-wise error rate (FWER) have been commonly used to correct for multiple
6  testing in genetic associations studies, e.g. a p-value threshold of $5 \times 10^{-8}$ based on a Bonferroni correction
7  is commonly used for genome-wide significance in GWAS corresponding to a FWER at 0.05. The number
8  of genetic variants being considered in the analysis of whole-genome sequencing data increases
9  substantially to more than 400 million in TOPMed[2], and FWER-controlling methods become highly
10 conservative[7]. As more individuals are being sequenced, the number of variants increases accordingly. The
11 false discovery rate (FDR), which quantifies the expected proportion of discoveries which are falsely
12 rejected, is an alternative metric to the FWER in multiple testing control, and can have greater power to
13 detect true positives while controlling FDR at a specified level. This metric has been popular in the
14 discovery of eQTLs and Bayesian association tests for rare variation in autism spectrum disorder studies[8-11]. Given the limited power of conventional association tests for whole-genome sequencing data and the
16 potential for many true discoveries to be made in studies for highly polygenic traits, controlling FDR can
17 be a more appealing strategy. However, the conventional FDR-controlling methods, such as the Benjamini-
18 Hochberg procedure[12], often do not appropriately account for correlations among tests and therefore cannot
19 guarantee FDR control at the target level, which can limit the widespread application of FDR control to
20 whole-genome sequencing data.

21 The knockoff framework is a recent breakthrough in statistics to control the FDR under arbitrary correlation
22 structure and to improve power over methods controlling the FWER[13,14]. The main idea behind it is to first
23 construct synthetic features, i.e. knockoff features, that resemble the true features in terms of the correlation
24 structure but are conditionally independent of the outcome given the true features. The knockoff features
25 serve as negative controls and help us select the truly important features, while controlling the FDR.
26 Compared to the well-known Benjamini-Hochberg procedure[12], which controls the FDR under
27 independence or a type of positive-dependence, the knockoff framework appropriately accounts for
28 arbitrary correlations between the original variables while guaranteeing control of the FDR. Moreover, it is
29 not limited to using calibrated p-values, and can be flexibly applied to feature importance scores computed
30 based on a variety of modern machine learning methods, with rigorous finite-sample statistical guarantees.
31 Several knockoff constructions have been proposed in the literature including the second-order knockoff
32 generator proposed by Candès et al.[14] and the knockoff generator for Hidden Markov Models (HMMs)
33 proposed by Sesia et al.[15,16]. The HMM construction has been applied to phased GWAS data in the UK
34 biobank. However, these constructions can fail for rare variants in whole-genome sequencing data whose
35 distribution is highly skewed and zero-inflated, leading to inflated FDR. Romano et al.[17] proposed deep
36 generative models for arbitrary and unspecified data distributions, but such an approach is computationally
37 intensive, and therefore not scalable to whole-genome sequencing data.

38 Our contributions in this paper include a sequential knockoff generator, a powerful genome-wide screening
39 method, and a robust inference procedure integrating multiple knockoffs. The sequential knockoff generator
40 is more than 50 times faster than state-of-the-art knockoff generation methods, and additionally allows for
41 the efficient generation of multiple knockoffs. The genome-wide screening method builds upon our recently
42 proposed scan statistic framework, WGScan[18], to localize association signals at genome-wide scale. We
43 adopt the same screening strategy, but incorporate several recent advances for rare-variant analysis in
44 sequencing studies, including the aggregated Cauchy association test to combine single variant tests, burden
45 and dispersion (SKAT) tests, the saddlepoint approximation for unbalanced case-control data, the
46 functional score test that allows incorporation of functional annotations, and a modified variant threshold
47 test that accumulates extremely rare variants such as singletons and doubletons[19-26]. We compute statistics
48 measuring the importance of the original and knockoff features using an ensemble of these tests. Feature
49 statistics that contrast the original and knockoff statistics are computed for each feature, and can be used

1  by the knockoff filter to select the important features, i.e. those significant at a fixed FDR threshold. The
2  integration of multiple knockoffs further helps improve the power, stability and reproducibility of the results
3  compared with state-of-the-art alternatives. Using simulations and applications to two whole-genome
4  sequencing studies, we show that the proposed method is powerful in detecting signals across the genome
5  with guaranteed FDR control.

6  Our knockoff method can be considered a synthetic alternative to knockout functional experiments designed
7  to identify functional variation implicated in a trait of interest. For each individual in the original cohort,
8  the proposed method generates a synthetic sequence where each genetic variant is being randomized,
9  making it silent and not directly affecting the trait of interest while preserving the sequence correlation
10  structure. Then the proposed method compares the original cohort where the variants are potentially
11  functional with the synthetic cohort where the variants are silenced. The randomization utilizes the knockoff
12  framework that ensures that the original sequence and the synthetic sequence are "exchangeable". That is,
13  if one replaces any part of the original sequence by its synthetic, silenced sequence, the joint distribution of
14  genetic variants (the LD structure etc.) remains the same. This leads to an important feature of our proposed
15  screening procedure that is similar to real functional experiments, namely the ability to prioritize causal
16  variants over associations due to linkage disequilibrium and other unadjusted confounding effects (e.g.
17  shadow effects of nearby significant variants and unadjusted population stratification) as we show below.

18  In this paper, we present a statistical approach that addresses the challenges described above, and leads to
19  increased power to detect and localize common and rare risk variants at genome-wide scale. The framework
20  appropriately accounts for arbitrary correlations while guaranteeing FDR control at a desired level, and
21  therefore has higher power than existing association tests that control FWER. Furthermore, the proposed
22  method has additional important advantages over the standard approaches due to some intrinsic properties
23  of the underlying framework. Specifically, it allows for the prioritization of causal variants over
24  associations due to linkage disequilibrium. For analyses specifically focusing on rare variants, the method
25  naturally distinguishes the signal due to rare variants from shadow effects of nearby significant (common
26  or rare) variants. Additionally, it naturally reduces false positives due to unadjusted population stratification.

27

28  **Results**

29  **Overview of the screening procedure with multiple knockoffs (*KnockoffScreen*).** We describe here the
30  main ideas behind our method, *KnockoffScreen*. We assume a study population of $n$ subjects, with $Y_i$ being
31  the quantitative/dichotomous outcome value; $\mathbf{X_i} = (X_{i1}, \ldots, X_{id})^T$ being the $d$ covariates which can include
32  age, gender, principal components of genetic variation etc.; $\{G_{ij}\}_{1 \leq j \leq p}$ being the $p$ genetic variants in the
33  genome. For each target window $\Phi_{kl} = \{j : k \leq j \leq l\}$, we are interested in determining whether $\Phi_{kl}$
34  contains any variants associated with the outcome of interest while adjusting for covariates.

35  The idea of the proposed method is to augment the original cohort with a synthetic cohort with genetic
36  variants, $\{\tilde{G}_{ij}\}_{1 \leq j \leq p}$, referred to as knockoff features. $\{\tilde{G}_{ij}\}_{1 \leq j \leq p}$ are generated by a data driven algorithm
37  such that they are exchangeable with $\{G_{ij}\}_{1 \leq j \leq p}$, yet they do not directly affect $Y_i$ (i.e. are "silenced", and
38  therefore not causal). More precisely, $\{\tilde{G}_{ij}\}_{1 \leq j \leq p}$ is independent of $Y_i$ conditional on $\{G_{ij}\}_{1 \leq j \leq p}$. Note that
39  the knockoff generation procedure is different from the well-known permutation procedure which generates
40  control features by permuting the samples; for such a permutation procedure, the exchangeability property
41  between the original genetic variants and the synthetic ones does not hold and hence the FDR control cannot
42  be guaranteed[13,14].

43  The screening procedure examines every target window $\Phi_{kl}$ in the genome and performs hypothesis testing
44  in both the original cohort and the synthetic cohort, to test for association of $G_{\Phi_{kl}}$ and $\tilde{G}_{\Phi_{kl}}$ with $Y$

1 respectively. As explained below, the knockoff procedure is amenable to any form of association test within
2 the window. Let $p_{\Phi_{kl}}, \tilde{p}_{\Phi_{kl}}$ be the resulting p-values. We define a feature statistic as

$$W_{\Phi_{kl}} = T_{\Phi_{kl}} - \tilde{T}_{\Phi_{kl}}, \qquad (1)$$

4 where $T_{\Phi_{kl}} = -\log_{10} p_{\Phi_{kl}}$ and $\tilde{T}_{\Phi_{kl}} = -\log_{10} \tilde{p}_{\Phi_{kl}}$. Essentially, the observed p-value for each window is
5 compared to its control counterpart in the synthetic cohort. A threshold $\tau$ for $W_{\Phi_{kl}}$ can be determined by
6 the knockoff filter so that the FDR is controlled at the nominal level. We select all windows with $W_{\Phi_{kl}} \geq \tau$.
7 We additionally derived the corresponding Q-value for a window, $q_{\Phi_{kl}}$, that unifies the feature statistic
8 $W_{\Phi_{kl}}$ and the threshold $\tau$. More details are given in the Methods section.

9 The knockoff construction ensures exchangeability of features, namely that $\{G_{ij}\}_{1 \leq j \leq p}$ and $\{\tilde{G}_{ij}\}_{1 \leq j \leq p}$ are
10 exchangeable. Hence if one swaps any subset of variants with their synthetic counterpart, the joint
11 distribution remains the same. For instance, suppose that $G_{i1}$ and $G_{i2}$ are two genetic variants, then the
12 knockoff generator will generate their knockoff counterparts $\tilde{G}_{i1}$ and $\tilde{G}_{i2}$ such that
13 $(G_{i1}, G_{i2}, \tilde{G}_{i1}, \tilde{G}_{i2}) \sim (G_{i1}, \tilde{G}_{i2}, \tilde{G}_{i1}, G_{i2})$, where "$\sim$" denotes equality in distribution. More generally, for any
14 subset $S \subset \{1, \dots, p\}$,

$$(G_i, \tilde{G}_i)_{\text{swap}(S)} \sim (G_i, \tilde{G}_i), \qquad (2)$$

16 where $(G_i, \tilde{G}_i)_{\text{swap}(S)}$ is obtained from $(G_i, \tilde{G}_i)$ by swapping the variants $G_{ij}$ and $\tilde{G}_{ij}$ for each $j \in S$. This
17 feature exchangeability implies the exchangeability of the importance scores $T_{\Phi_{kl}}$ and $\tilde{T}_{\Phi_{kl}}$ under the null
18 hypothesis, i.e. $(T_{\Phi_{kl}}, \tilde{T}_{\Phi_{kl}}) \sim (\tilde{T}_{\Phi_{kl}}, T_{\Phi_{kl}})$ if $\Phi_{kl}$ does not contain any causal variant. Thus $\tilde{T}_{\Phi_{kl}}$ can be used
19 as the negative control, and we reject the null when $W_{\Phi_{kl}} = T_{\Phi_{kl}} - \tilde{T}_{\Phi_{kl}}$ is sufficiently large. This
20 exchangeability property leads to several interesting properties of our proposed screening procedure relative
21 to conventional association tests as mentioned in the Introduction, and which will be discussed in detail in
22 later sections.

23 Once the knockoff generation is completed, we apply a genome-wide screening procedure. Our screening
24 procedure considers windows with different sizes (1bp, 1kb, 5kb, 10kb) across the genome, with half of
25 each window overlapping with adjacent windows of the same size. To calculate the importance score for
26 each window $\Phi_{kl}$, we incorporate several recent advances for association tests for sequencing studies to
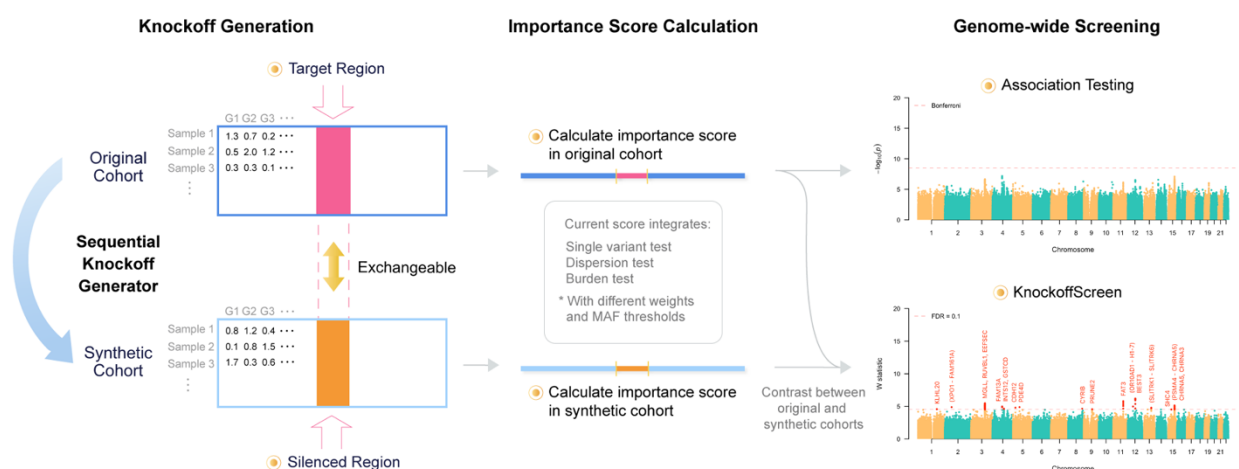27 compute $p_{\Phi_{kl}}$.

28 • For each 1bp window (i.e. single variant): we only consider common (minor allele frequency
29   (MAF)>0.05) and low frequency (0.01<MAF<0.05) variants and compute $p_{\Phi_{kl}}$ from single variant
30   score test. For binary traits, we implement the saddlepoint approximation for unbalanced case-control
31   data.
32 • For each 1kb/5kb/10kb window, we perform:
33   a. Burden and dispersion tests for common and low frequency variants with Beta (MAF, 1, 25)
34      weights, where Beta (.) is the probability density function of the beta distribution with shape
35      parameters 1 and 25[26]. These tests aim to detect the combined effects of common and low frequency
36      variants.
37   b. Burden and dispersion tests for rare variants (MAF<0.01 & minor allele count (MAC)>=5) with
38      Beta (MAF,1, 25) weights. These tests aim to detect the combined effects of rare variants.
39   c. Burden and dispersion tests for rare variants, weighted by functional annotations[23]. Current
40      implementation includes CADD[27] and tissue/cell type specific GenoNet scores[28]. These tests aim
41      to utilize functional annotations for improved power.
42   d. Burden test for aggregation of ultra-rare variants (MAC<5). These tests aim to aggregate effects
43      from extremely rare variants such as singletons, doubletons etc.
44   e. Single variant score tests for common, low frequency and rare variants in the window.

1    f.   The aggregated Cauchy association test[29] to combine a-e to compute $p_{\Phi_{kl}}$.

2    We also extend the single knockoff described above to the setting with multiple knockoffs to improve the
3    power, stability and reproducibility of the findings. Let $q$ be the FDR threshold. The inference based on a
4    single knockoff is limited by a detection threshold of $1/q$, defined as the minimum number of independent
5    signals required for making any discovery. It has no power at the target FDR level $q$ if there are fewer than
6    $1/q$ discoveries to be made. The multiple knockoffs improve the detection threshold from $1/q$ to $1/(Mq)$,
7    where $M$ is the number of knockoffs[30]. For example, the detection threshold is 10 when the target FDR=0.1.
8    In scenarios where the signal is sparse (<10 independent causal variants) in the target region or across the
9    genome, inference based on a single knockoff can have very low power to detect any of the causal variants.
10   In such a setting, *KnockoffScreen* with $M$ knockoffs reduces the detection threshold from 10 to $10/M$,
11   which allows *KnockoffScreen* to detect sparse signals in a target region or across the genome. Furthermore,
12   integrating multiple knockoffs leads to improvements in the stability and reproducibility of the knockoff
13   procedure. Specifically, the results of the *KnockoffScreen* procedure depend to some extent on the sampling
14   of knockoff features $\left\{\tilde{G}_{ij}\right\}_{1 \le j \le p}$, which is random. Therefore, running the analysis twice on the same dataset
15   may lead to the selection of slightly different subsets of features. In particular, for weak causal effects, there
16   is a chance that the causal variant is selected in only one of the analyses. We demonstrate in the Methods
17   section that our choice of multiple knockoff statistics helps improve the stability of the results compared
18   with state-of-the-art alternatives.

19   In the Methods section, we describe in detail our computationally efficient method to generate the knockoff
20   features, and our multiple knockoffs method. A flowchart of our approach is shown in Figure 1.

21   **Figure 1: Overview of *KnockoffScreen*.** The left panel illustrates the knockoff generation based on the original genotype matrix.
22   Each row in the matrix corresponds to an individual and each column corresponds to a genetic variant. Each cell presents the
23   genotype value/dosage. The mid panel illustrates the calculation of the importance score for each 1bp, 1kb, 5kb or 10kb window.
24   The right panel presents a typical example of genome-wide screening results using conventional association testing (top) and
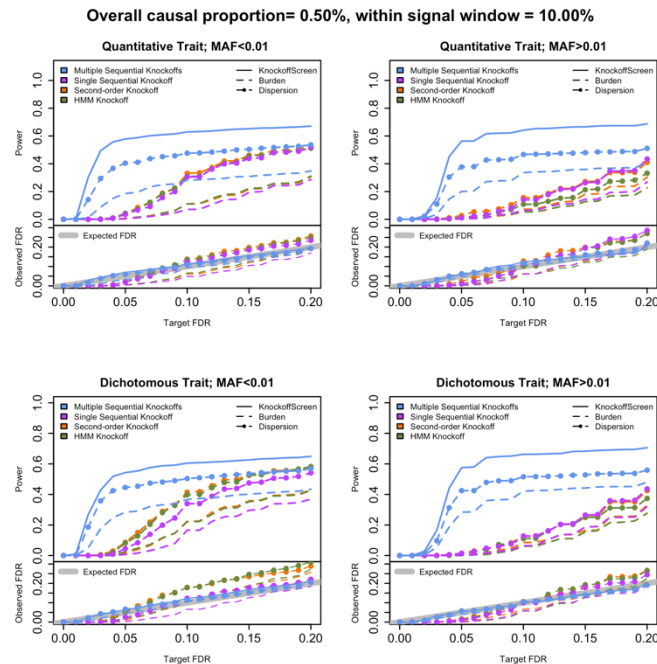25   *KnockoffScreen* (bottom).



26

27   ***KnockoffScreen* improves power and guarantees FDR control in single-region simulation studies.** We
28   performed empirical power and FDR simulations to evaluate the performance of *KnockoffScreen* in a single
29   region. We compared it with existing alternatives for sequence-based association testing, including the
30   burden and dispersion (SKAT) tests with Beta(MAF;1,25) weights. For a fair and simplified comparison,
31   we did not include additional functional annotations in our method for these simulations. Note that burden
32   and SKAT are also applied within the knockoff framework, and therefore we still aim at controlling the
33   FDR. We also compared with state-of-the-art methods for generating knockoff features, including the
34   second-order knockoff generator proposed by Candès et al.[14], referred to as SecondOrder, and the knockoff

1 generator for Hidden Markov Models (HMMs) proposed by Sesia et al.[15,16] with number of states S=50.
2 For simulating the sequence data, each replicate consists of 10,000 individuals with genetic data on 1,000
3 genetic variants from a 200kb region, simulated using the haplotype dataset in the SKAT package. The
4 SKAT haplotype dataset was generated using a coalescent model (COSI), mimicking the linkage
5 disequilibrium structure of European ancestry samples. Simulation details are provided in the Methods
6 section. We compared the methods in different scenarios for common and rare variants, quantitative traits
7 and dichotomous traits. For each replicate, the empirical power is defined as the proportion of detected
8 windows among all causal windows (windows that contain at least one causal variant); the empirical FDR
9 is defined as the proportion of non-causal windows among all detected windows. We present the average
10 power and FDR over 1,000 replicates in Figure 2. We additionally present the distribution of power and
11 FDP (the false discovery proportion) at target FDR level 0.1 over 1,000 replicates in Figure S1.

12 The comparisons of the different knockoff generators show that *KnockoffScreen* has significantly improved
13 power with a better FDR control. For single knockoff generators, SecondOrder and HMMs have inflated
14 FDR for rare variants. We also observed that the HMM based knockoff has inflated FDR for common
15 variants for the window-based screening procedure considered in this paper. *KnockoffScreen* has well
16 controlled FDR, and significantly higher power compared with single knockoff, especially when the target
17 FDR $q$ is small. This is due to the high detection threshold ($1/q$) needed for the single knockoff. Our
18 multiple knockoff method *KnockoffScreen* incorporates five knockoffs, and as a consequence the detection
19 threshold is reduced from $1/q$ to $1/(5q)$, which helps improve power. We note that the power of methods
20 with single knockoff and multiple knockoffs may be comparable in settings where the detection threshold
21 is not a primary factor that limits the power, such as for higher target FDR values. Furthermore, we observed
22 that the additional tests included in *KnockoffScreen* improve its power, compared to the burden and SKAT
23 tests with the same number of knockoffs. In summary, the simulation results show that the screening
24 procedure and multiple knockoffs help improve power while controlling FDR at the nominal level.

25 **Figure 2: Power and false discovery rate (FDR) simulation studies in a single region.** The four panels show power and FDR
26 base on 500 replicates for different types of traits (quantitative and dichotomous) and different types of variants (rare and common),
27 with different target FDR varying from 0 to 0.2. The different colors indicate different knockoff generators. The different types of
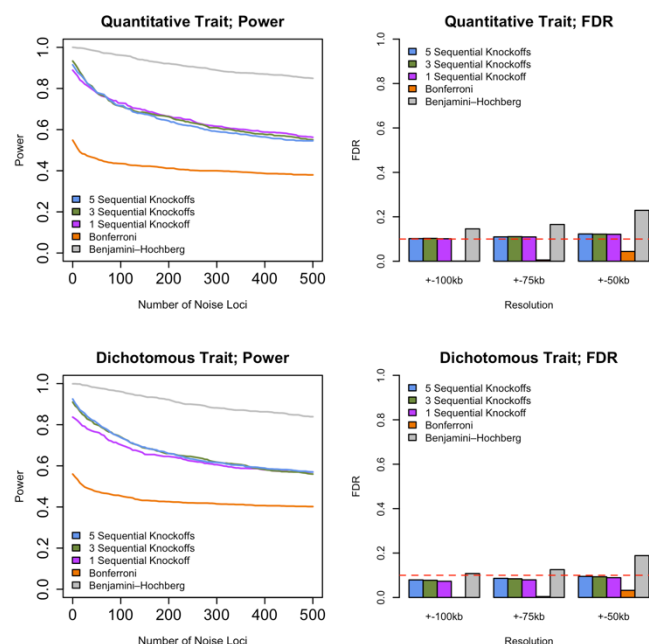28 lines indicate different tests to define the importance score.



Overall causal proportion= 0.50%, within signal window = 10.00%

29

30

1  ***KnockoffScreen* improves genome-wide locus discovery for polygenic traits.** We conducted genome-
2  wide empirical FDR and power simulations using ADSP whole-genome sequencing data to evaluate the
3  performance of *KnockoffScreen* in the presence of multiple causal loci. Specifically, we randomly choose
4  10 causal loci and 500 noise loci across the whole genome, each of size 200kb. Each causal locus contains
5  a 10kb causal window. For each replicate, we randomly set 10% variants in each 10kb causal window to
6  be causal. In total, there are approximately 335 causal variants on average across the genome. Simulation
7  details are provided in the Methods section. We compared the proposed *KnockoffScreen* method to
8  conventional p-value based methods including Bonferroni correction for FWER control, and BH procedure
9  for FDR control. For *KnockoffScreen* we also evaluated the effect of different numbers of knockoffs. We
10 evaluated the empirical power and FDR at target FDR 0.10. For each replicate, the power is defined as
11 proportion of the 200kb causal loci detected by each method; the empirical FDR is defined as the proportion
12 of significant windows +/- 100/75/50kb away from the causal windows. We report the average power and
13 FDR over 100 replicates in Figure 3.

14 The simulation results show that *KnockoffScreen* exhibits substantially higher power than using Bonferroni
15 correction. Additionally, using the conventional Benjamini-Hochberg FDR control may have higher power
16 than *KnockoffScreen*, but fails to control FDR at higher resolution (e.g. +/-75kb). Statistically, the knockoff
17 filter is expected to have similar or higher power for independent tests compared with the BH procedure[13].
18 For correlated genetic variants/windows, the higher empirical power of the BH procedure in our simulation
19 studies is subject to false-positive inflation. Therefore, we do not recommend directly using the
20 conventional BH FDR control in whole genome sequencing studies. In the presence of multiple causal loci
21 and at a moderate target FDR, we observe that the power is similar for different number of knockoffs
22 because the aforementioned detection threshold is no longer an issue. Thus, multiple knockoffs are
23 particularly useful when the number of causal loci is small, and the target FDR is stringent. Regardless of
24 the effect on power, an important advantage of using multiple knockoffs is that it can significantly improve
25 the stability and reproducibility of knockoff-based inference. Since the knockoff sampling is random, each
26 run of the knockoff procedure may lead to different selected sets of features. In practice, strong signals will
27 always be selected but weak signals may be missed at random with a single knockoff. The proposed multiple
28 knockoff procedure has significantly smaller variation in feature statistic in our simulation study based on
29 real data from ADSP. We discuss the details in the Methods section (Figure 9).

30

**Figure 3: Genome-wide power and false discovery rate (FDR) simulations studies in the presence of multiple causal loci.**
The two left panels show power for different types of traits (quantitative and dichotomous), defined as the average proportion of 200kb causal loci being identified at target FDR 0.1. The two right panels show empirical FDR for different types of traits (quantitative and dichotomous) at different resolutions, defined as the proportion of significant windows (target FDR 0.1) +/- 100/75/50kb away from the causal windows. The empirical power and FDR are averaged over 100 replicates.



***KnockoffScreen* prioritizes causal variants/loci over associations due to linkage disequilibrium.** The exchangeability properties for the features help the inference based on the feature statistic $W_{\Phi_{kl}} = T_{\Phi_{kl}} - \tilde{T}_{\Phi_{kl}}$ to prioritize causal variants/loci over associations due to LD. For example, suppose $G_{i1}$ is causal and $G_{i2}$ is a null variant correlated with $G_{i1}$; $(\tilde{G}_{i1}, \tilde{G}_{i2})$ are exchangeable with $(G_{i1}, G_{i2})$, therefore $\text{cor}(G_{i1}, \tilde{G}_{i2}) \approx \text{cor}(G_{i1}, G_{i2})$. Thus, the resulting p-values $p_2 \sim \tilde{p}_2$, and hence $W_2 = -\log p_2 - (-\log \tilde{p}_2)$ follows a distribution that is symmetric around 0. This way, by comparing the p-value of $G_{i2}$ (a null variant) to that of its control counterpart, the method no longer identifies the proxy variant $G_{i2}$ as significant. On the other hand, the knockoff generation minimizes the correlation between feature $G_{i1}$ and its knockoff counterpart $\tilde{G}_{i1}$, such that $W_1 = -\log p_1 - (-\log \tilde{p}_1)$ takes positive value with higher probability and therefore can identify the causal variant $G_{i1}$ as significant.

We compared *KnockoffScreen* with state-of-the-art methods which perform association tests in each window and apply a hard threshold (e.g. Bonferroni correction) to control for family wise error rate (FWER). For a fair comparison, for the conventional association testing we adopted the same combination of tests (i.e. we combined the same single variant and region-based tests) implemented in *KnockoffScreen* to calculate the p-value. As a proof of concept, we show first the results from an analysis of common and rare variants within a 200kb region near the apolipoprotein E (*APOE*) gene for Alzheimer's Disease (AD), using data on 3,894 individuals from the Alzheimer's Disease Sequencing Project (ADSP). More details on the data analysis for ADSP are described in a later section. *APOE* is a major genetic determinant of AD risk, containing AD risk/protective alleles. *APOE* comes in three forms (*APOE* $\varepsilon 2/\varepsilon 3/\varepsilon 4$). Among them, $\varepsilon 2$ is the least common and confers reduced risk to AD, $\varepsilon 4$ is the most common and increases risk to AD, while $\varepsilon 3$ appears neutral. We found that the conventional association test using a Bonferroni correction identifies a large number of significant associations ($p < 0.05$/number of tested windows), but most of these windows are presumably false positives due to LD since they are no longer significant after adjusting for the *APOE* alleles (Figure 4A). In contrast, *KnockoffScreen* filtered out a considerable number of associations that are

8

1   likely due to LD, and identified more refined windows that reside in *APOE* and *APOC1* at target FDR=0.1
2   (Figure 4B). A recent study identified AD risk variants and haplotypes in the *APOC1* region, and showed
3   that these signals are independent of the *APOE-ε4* coding change, consistent with our findings[31].

4   We conducted additional simulation studies to further investigate this property. We randomly drew a subset
5   of variants (1,000 variants) from the 200kb region near *APOE*, set a 5kb window (similar to the size of
6   *APOE*) as the causal window and then simulated disease phenotypes. More details on these simulations are
7   provided in the Methods section. With target FDR=0.1, we evaluated the proportion of selected windows
8   overlapping the true causal window, and the maximum distance between the selected windows and the
9   causal window. Figures 4C-D show the results over 500 replicates. We found that windows selected by
10  *KnockoffScreen* have a significantly better chance to overlap with the causal window relative to the
11  conventional association test. We also found that the maximum distance between the selected windows and
12  the causal window is significantly smaller for *KnockoffScreen*. Particularly, the distribution of the
13  maximum distance to the causal window is zero-inflated for *KnockoffScreen*; these are cases where all
14  windows detected by *KnockoffScreen* overlap/cover the causal window.

15  Overall, the real data example and these simulation results demonstrate that *KnockoffScreen* is able to
16  prioritize causal variants over associations due to linkage disequilibrium and produces more accurate results
17  in detecting disease risk variants/loci, thereby improving interpretation of the findings.

18  **Figure 4: *KnockoffScreen* prioritizes causal variants/loci and distinguishes the signal due to rare variants from shadow**
19  **effects of significant common variants nearby.** The top two panels present the results of the data analyses of the APOE+/-100kb
20  region from the ADSP data. Each dot represents a window. Windows selected by *KnockoffScreen* are highlighted in red. Windows
21  selected by conventional association testing but not by *KnockoffScreen* are shown in gray. The bottom three panels present
22  simulation results based on the APOE+/-100kb region, comparing the conventional association testing and *KnockoffScreen* methods
23  in terms of prioritizing causal regions, and distinguishing true signals from shadow effects of nearby variants. The target FDR is
24  0.1. The results are based on 500 replicates.



25

26  ***KnockoffScreen* distinguishes the signal due to rare variants from shadow effects of significant**
27  **common variants nearby.** Conventional sequence-based association tests focused on rare variants (MAF
28  below a certain threshold, e.g. 0.01) can lead to false positive findings by identifying rare variants that are
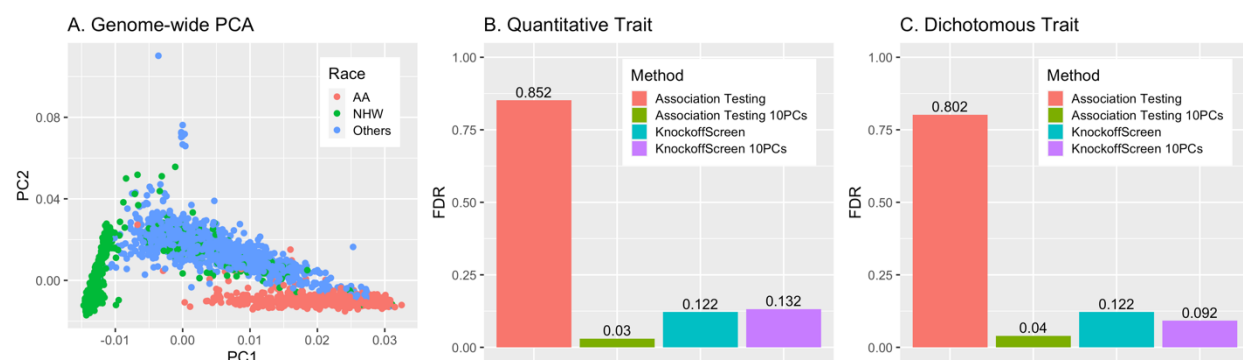
1  not causal but instead correlated with a known causal common variant at the same locus; this is referred to
2  as the shadow effects[32]. For illustration, we conducted simulation studies based on the same 200kb region
3  near the *APOE* gene as described above. We adopted the same simulation setting but set the causal variants
4  to be common (MAF>0.01) and apply the methods to rare variants only (MAF<0.01). More details on these
5  simulations are provided in the Methods section. Since all causal variants are common, all detected
6  windows (focusing on rare variants) are false positives due to the shadow effect. We compared
7  *KnockoffScreen* with conventional association testing by counting the number of false positives and show
8  the distribution over 500 replicates in Figure 4E. For a fair comparison, for the conventional association
9  testing we adopted the same ensemble of tests implemented in *KnockoffScreen* to calculate the p-value. We
10 observed that the conventional tests tend to identify a large number of false positives due to the shadow
11 effect. In contrast, *KnockoffScreen* has a significantly reduced number of false positives, demonstrating that
12 it is able to distinguish the effect of rare variants from that of common variants nearby. This feature is
13 particularly appealing in detecting novel rare association signals in whole-genome sequencing studies. The
14 same argument also holds if instead rare variants were causal; by construction, *KnockoffScreen* applied to
15 common variants only can distinguish effects attributable to common causal variants from those due to rare
16 causal variants nearby.

17 **Empirical evaluation of *KnockoffScreen* in the presence of population stratification.** Population
18 structure is an important confounder in genetic association studies. Standard methods to adjust for
19 population stratification, including principal component analysis or mixed effect models, help control for
20 global ancestry in conventional sequencing association tests. We performed an empirical evaluation of
21 *KnockoffScreen* in the presence of population stratification using sequencing data from the ADSP project.
22 We also evaluated whether, by regressing out the top principal components when computing the association
23 statistics (p-values), *KnockoffScreen* is able to control FDR. Specifically, we randomly drew a subset of
24 variants (1,000 variants) from the 200kb region near the *APOE* region in the ADSP study. The ADSP
25 includes three ethnic groups: African American (AA), Non-Hispanic White (NHW) and Others (of which,
26 98% are Caribbean Hispanic) (see genome-wide PCA results in Figure 5A). We set the mean/prevalence
27 for the quantitative/dichotomous trait to be a function of the subpopulation, but not directly affected by any
28 genetic variants. More details on these simulations are provided in the Methods section. We compared
29 *KnockoffScreen* with the conventional association test with no adjustment for population stratification. We
30 also included a modified version of *KnockoffScreen* that adjusts for the top 10 global PCs when computing
31 the p-values used to compute the window feature statistic, referred to as *KnockoffScreen+10PCs*. For
32 comparison, we also included the conventional association test based on Bonferroni correction, which
33 defines significant associations by p-value<0.05/number of tests.

34 Since in these simulations none of the genetic variants are causal, all detected windows are false positives
35 due to the confounding effects of population structure. With a target FDR=0.1, we calculated the observed
36 FDR, defined as the proportion of replicates where any window is detected, and present the results in
37 Figures 5B-C. We observed that both PC-adjusted *KnockoffScreen* and the conventional PC-adjusted
38 association test are able to control FDR at the target level.  This is further illustrated by our real data analysis
39 of ADSP where despite the combined analysis of three ethnicities there is no apparent inflation in false
40 positive signals. Interestingly, *KnockoffScreen* exhibits lower FDR than association test when they are both
41 unadjusted, indicating that the use of knockoffs naturally helps to prioritize causal variants over association
42 due to population stratifications. We additionally performed simulation studies to mimic population
43 stratification driven by rare variants and present the results in Table S1. As before, we found that both PC-
44 adjusted *KnockoffScreen* and association test are able to control FDR in the scenarios considered here, and
45 *KnockoffScreen* exhibits a lower FDR than the conventional association test for an unadjusted model. Since
46 the reduction of false positives for *KnockoffScreen* does not require observing/estimating the underlying
47 ancestry, the knockoff procedure can potentially complement existing tools for ancestry adjustment to better
48 reduce false positive findings due to population substructure. However, we clarify that *KnockoffScreen*
49 itself does not completely eliminate the confounding due to population stratification (Table S1) because the
50 current knockoff generator assumes the same LD structure across individuals and it only accounts for local

1   LD structure. Therefore, it does not capture heterogeneous LD structure across populations and strong long-
2   range LD due to population stratification. Development of new knockoff generators that explicitly account
3   for population structure will be of interest[33].

4
5   **Figure 5: Empirical evaluation of *KnockoffScreen* in the presence of population stratification**. The left panel presents the
6   principal component analysis of the ADSP data, which contains three ethnic groups: African American (AA), Non-Hispanic White
7   (NHW) and Others (of which, 98% are Caribbean Hispanic). Each dot represents an individual. The middle and right panels present
8   results from a simulation study that mimics the ADSP data, comparing *KnockoffScreen* with conventional association testing. Each
9   panel shows empirical FDR based on 500 replicates. *KnockoffScreen* 10PCs is a modified version of *KnockoffScreen* method that
10  includes adjustment for the top principal components while computing the association statistics (p-values). *KnockoffScreen* controls
11  FDR at 0.10; Association Testing is based on usual Bonferroni correction (0.05/number of tests), controlling FWER at 0.05.



12

13  ***KnockoffScreen* enables computationally efficient screening of whole-genome sequencing data**. One
14  obstacle for the widespread application of knockoffs to genetic data, particularly whole-genome sequencing
15  data, is their computational cost. The knockoff generation can be computationally intensive when the
16  number of genetic variants $p$ is large; depending on the method, it may require the calculation of the eigen
17  values of a $p \times p$ covariance matrix, or iteratively fitting a prediction model for every variant. The whole-
18  genome sequencing data from ADSP (~4000 individuals) contains ~85 million variants in total, much larger
19  than the number of variants in GWAS datasets. Similarly, in 53,581 TOPMed samples, more than 400
20  million single-nucleotide and insertion/deletion variants were detected[2]. As more individuals are being
21  sequenced, the number of variants will increase accordingly. We demonstrate that the proposed sequential
22  model to simultaneously generate multiple knockoffs is significantly more computationally efficient than
23  existing knockoff generation methods, making it scalable to whole-genome sequencing data. We compared
24  the computing time of our proposed knockoff generator with two existing alternatives: the second-order
25  knockoff generator proposed by Candès et al.[14], referred to as SecondOrder; and knockoffs for Hidden
26  Markov Models (HMMs) proposed by Sesia et al.[15,16] with varying number of states (S=12 and S=50). We
27  estimate the complexity of our proposed method as $O(np)$, where $n$ is the sample size and $p$ is the number
28  of genetic variants. The details of this calculation are described in the Methods section. The complexity of
29  the HMM method is also $O(np)$, as discussed in Sesia et al.[16]. However, it is significantly less efficient
30  than the proposed method for unphased genotype data as we show below. We note that the computing time
31  of the SecondOrder method is of order $O(np^2 + p^3)$ because it requires calculating the eigen values of a
32  $p \times p$ covariance matrix. Therefore, it is not a feasible approach for whole-genome analysis with a large
33  number of variants.

34  We performed simulations to empirically evaluate the computational time for the different methods. We
35  note that the proposed method focuses on the analysis of whole-genome sequencing data, and thus the
36  computational cost is reported on unphased genotype data, which is the usual format for sequencing data.
37  Since the HMM model assumes the availability of phased data, we report the computing time separately
38  for phasing with fastPhase and sampling with SNPknock as described in Sesia et al.[15]. We simulated genetic
39  data using the SKAT package, with varying sample sizes and number of genetic variants (Table 1). The
40  computing time was evaluated on a single CPU (Intel(R) Xeon(R) CPU E5-2640 v3 @ 2.60GHz). For the

1 simulation scenario considered in the previous section with 10,000 individuals and 1,000 genetic variants,
2 we observed that the proposed method takes 6.59 seconds to generate a single set of knockoff features,
3 which is ~130 times faster than the HMM model with S=12 states (881.43 seconds). The application of the
4 HMM model with the recommended S=50 states to unphased sequencing data (13681.53 seconds for 10,000
5 individuals and 1,000 genetic variants) is currently not practical at genome-wide scale. As shown, a
6 substantial fraction of the total computing time is taken by the phasing step, and therefore using more
7 computationally efficient phasing algorithms can further improve the computational cost of the HMM-
8 based knockoff generation.

9 **Table 1: Computing time of different knockoff generators.** Each cell shows the computing time in seconds to generate knockoffs
10 based for unphased genotype data. The multiple sequential knockoffs approach generates five knockoffs. The computing time was
11 measured on unphased genotype data using a single CPU (Intel(R) Xeon(R) CPU E5-2640 v3 @ 2.60GHz). Since the HMM model
12 was mainly proposed for phased data, we report the computing time separately for phasing with fastPhase, and sampling with
13 SNPknock.

| n | p | MSK (5 knockoffs) | SK | SecondOrder | HMM with S=12 | | HMM with S=50 | |
|---|---|---|---|---|---|---|---|---|
| | | | | | Phasing | Sampling | Phasing | Sampling |
| 1000 | 500 | 2.11 | 0.86 | 8.9 | 37.86 | 6.02 | 580.87 | 93.88 |
| 1000 | 1000 | 3.99 | 1.92 | 57.01 | 76 | 12.01 | 1147.66 | 188.74 |
| 1000 | 2000 | 8.89 | 4.06 | 491.19 | 161.94 | 24.76 | 2336.83 | 376.93 |
| 5000 | 500 | 4.66 | 1.63 | 8.51 | 188.5 | 30.45 | 2878.43 | 485.34 |
| 5000 | 1000 | 11.76 | 3.95 | 52.63 | 380.06 | 60.28 | 5914.19 | 996.11 |
| 5000 | 2000 | 31.58 | 11.09 | 479.01 | 811.61 | 129.6 | 11734.66 | 1865.11 |
| 10000 | 500 | 7.42 | 2.34 | 9.29 | 377.07 | 58.8 | 5784.24 | 957.49 |
| 10000 | 1000 | 20.57 | 6.59 | 54.66 | 757.49 | 123.94 | 11744.68 | 1936.85 |
| 10000 | 2000 | 52.86 | 16.92 | 445.05 | 1571.19 | 253.46 | 23584.8 | 3870.07 |

14

15 ***KnockoffScreen* detects more independent disease risk loci across the genome in two whole-genome**
16 **sequencing studies.** Here we show results from the application of *KnockoffScreen* to two whole-genome
17 sequencing datasets from two different studies, namely the Alzheimer's Disease Sequencing Project
18 (ADSP), and the COPDGene study from the NHLBI Trans-Omics for Precision Medicine (TOPMed)
19 Program. For each study, we considered windows with sizes (1bp, 1kb, 5kb, 10kb) across the genome as
20 described before. In addition to the different weighting and thresholding strategies, we include several
21 functional scores to improve the power of detecting rare functional variants. The functional scores include
22 non-tissue specific CADD score and 10 tissue/cell type specific GenoNet scores. The GenoNet scores were
23 trained using epigenetic annotations from the Roadmap Epigenomics Project across 127 tissues/cell types.
24 We partition the tissues/cell types into 10 groups (including Stem Cells, Blood, Connective Tissue, Brain,
25 Internal Organs, Fetal Brain, Muscle, Fetal Tissues, and Gastrointestinal; Table S2 has more details on
26 these tissue groupings) and we use the maximum GenoNet score per group.

27 We show results from conventional association tests (using the same combination of single variant and
28 region-based tests as implemented in *KnockoffScreen*) and using Bonferroni correction ($p < 0.05$/number
29 of tested windows) to control the family wise error rate. QQ-plots of all tests (Figure S5) show that the type
30 I error rate is well controlled. We also report results from *KnockoffScreen* at an FDR threshold of 0.1. We
31 assigned each significant window to its overlapping locus (gene or intergenic region). If the locus is a gene,
32 we report the gene's name; if the locus is intergenic, we report the upstream and downstream genes
33 (enclosed within parentheses and separated by "-"). To assess the degree of overlap with previously
34 described associations, we additionally searched if the loci have known associations with Alzheimer's
35 disease and lung related traits in the NHGRI-EBI GWAS Catalog[34], acknowledging that some of the studies

1  in the GWAS catalog included ADSP and COPDGene data. The details on gene annotations are described
2  in the Methods section.

3  ***Application to ADSP.*** We first applied *KnockoffScreen* to the whole-genome sequencing data from the
4  Alzheimer's Disease Sequencing Project (ADSP) for a genome-wide scan. The data includes 3,085 whole
5  genomes from the ADSP Discovery Extension Study and 809 whole genomes from the Alzheimer's Disease
6  Neuroimaging Initiative (ADNI), for a total of 3,894 whole genomes. More details on the ADSP data are
7  provided in the Methods section. We adjusted for age, age^2, gender, ethnic group, sequencing center, and
8  the leading 10 principal components of ancestry. We present the results in Figure 6.

9  The conventional association test with Bonferroni correction identified a region (~50kb long) at the known
10  *APOE* locus, containing a large number of significant associations (Figure 6), but, as discussed before, most
11  of them are presumably due to LD with the known *APOE* risk variants since they are no longer significant
12  after adjusting for the *APOE* alleles. Within the *APOE* region, *KnockoffScreen* identified fewer windows
13  that overlap with known AD genes, namely *APOE*, *APOC1, APOC1P1* and *TOMM40* at FDR<0.1, while
14  removing a considerable number of associations that are likely due to LD. Beyond the *APOE* locus,
15  *KnockoffScreen* identified several other loci that potentially affect AD risk, including *KAT8* and an
16  intergenic region on chromosome 18q22 between *DSEL* and *TMX3*. *KAT8* (lysine acetyltransferase 8) has
17  been recently identified in two large scale GWAS focused on clinically diagnosed AD and AD-by-proxy
18  individuals[35,36]. It is a promising candidate gene that affects multiple brain regions including the
19  hippocampus and plays a putative role in neurodegeneration in both AD and Parkinson's disease[37]. The
20  intergenic region identified by *KnockoffScreen* resides in a known linkage region for AD and bipolar
21  disorder on chromosome 18q22.1[38,39]. *DSEL* (dermatan sulfate epimerase-like) is implicated in D-
22  glucuronic acid metabolism and tumor rejection. A recent study has shown that glucuronic acid levels
23  increase with age and predict future healthspan-related outcomes[40]. Furthermore, *DSEL* is highly expressed
24  in the brain and has been found associated with AD in an imaging-wide association study[41]. SNPs upstream
25  of *DSEL* have also been associated with recurrent early-onset major depressive disorder[42]. Two other
26  intergenic loci, *ANKRD18A-FAM240B* and *TAFA5-BRD1* were reported in the GWAS catalog to have
27  suggestive associations ( $5 \times 10^{-8} < p < 1 \times 10^{-5}$ ) with late-onset Alzheimer's disease[43]. We
28  additionally present results when applying the Benjamini-Hochberg procedure for FDR control in Figure
29  S6; we observed that the associations identified by *KnockoffScreen* are largely replicated in the GWAS
30  catalog, while the new discoveries uniquely identified using the conventional BH FDR control do not
31  overlap with previous GWAS findings, suggesting they may be false positives.

32  ***Application to COPDGene study in TOPMed.*** The Genetic Epidemiology of COPD (COPDGene) study
33  includes current and former cigarette smokers aged > 45. All subjects underwent spirometry to measure
34  lung function. Cases were identified as those with moderate-to-severe chronic obstructive pulmonary
35  disease (COPD), controls were those with normal lung function, and a third set were neither cases nor
36  controls. These individuals have been whole-genome sequenced as part of the larger TOPMed project at an
37  average ~30X coverage depth, with joint-sample variant calling and variant level quality control in
38  TOPMed samples[2,44]. The COPDGene Freeze 5b dataset used for this analysis includes a total of 8,444
39  individuals, of which 5,713 are Non Hispanic White and 2,731 are African American. We tested lung
40  function measurements on all individuals: forced expiratory volume in one second ($FEV_1$), forced vital
41  capacity (FVC) and their ratio ($FEV_1/FVC$), as well as for case-control COPD status on a subset (NHW:
42  2366 cases/2084 controls, AA: 702 cases/1409 controls).

43  We applied *KnockoffScreen* separately to the two ethnic groups, and four phenotypes, while adjusting for
44  covariates as follows. In all analyses we adjusted for sequencing center, and the 10 leading principal
45  components of ancestry. Additionally, for $FEV_1$ and $FEV_1/FVC$ ratio, we adjusted for age, $age^2$, gender,
46  height, $height^2$, pack-years of smoking, and current smoking. For FVC, we adjusted for age, $age^2$, gender,
47  height, $height^2$, weight, pack-years of smoking, and current smoking. For COPD case/control status, we
48  adjusted for age, gender, and pack-years of smoking. Results for the NHW group for $FEV_1$ are shown in
49  Figure 7 and those for $FEV_1/FVC$ are shown in Figure S4.

1  Note that for $FEV_1$ and $FEV_1/FVC$, *KnockoffScreen* has been able to identify many more significant
2  associations compared with the application to Alzheimer's disease, a reflection of the larger sample size
3  but also the higher degree of polygenicity for lung function phenotypes relative to AD. Compared with the
4  conventional association test with Bonferroni correction, *KnockoffScreen* detected several known signals
5  for $FEV_1$, including the *PSMA4/CHRNA5/CHRNA3* locus on chromosome 15, the *INTS12/GSTCD* locus
6  on chromosome 4, and the *EEFSEC/RUVBL1* locus on chromosome 3. Overall, the majority of the single
7  variant signals that were found significant at FDR 0.1 have been associated with COPD-related phenotypes
8  in the GWAS catalog (81.8% for $FEV_1$ and 69.2% for $FEV_1/FVC$) (Figures 7 and S4) supporting the ability
9  of *KnockoffScreen* to identify previously discovered loci in GWAS studies with sample sizes much larger
10  than used here. *KnockoffScreen* additionally identified new loci by aggregating common/rare variants.
11  Although the new loci identified by *KnockoffScreen*, particularly those identified by rare variant methods,
12  will need to be validated in larger datasets, and the effector genes are not known, some of the genes in these
13  regions may be of interest. For FVC and COPD, as well as all traits for the African-Americans, we did not
14  identify any significant associations at FDR 0.1, likely a reflection of low power due to the smaller sample
15  size and possibly non-genetic covariates that might be associated with risk in AA and unaccounted for in
16  these analyses.

17  It is interesting to note that the significant loci identified by *KnockoffScreen* are markedly enriched for
18  windows (single bp or larger) overlapping protein coding genes despite an unbiased screen of the entire
19  genome. In particular, 40%, 80% and 56.4% of the loci significant for AD, $FEV_1$ and $FEV_1/FVC$
20  respectively overlap protein coding genes. Given the modest sample size of the datasets analyzed here, this
21  is perhaps expected; *KnockoffScreen* is able to identify the stronger effects closer to genes (e.g. coding and
22  promoter regions). As sample sizes for whole-genome sequencing studies continue to increase, we can
23  expect additional loci in noncoding regions to be identified.

24  In summary, these empirical results suggest that *KnockoffScreen* can identify additional signals that are
25  missed by conventional Bonferroni correction, while filtering out proxy associations that are likely due to
26  LD. Scatter plots comparing genome-wide W statistics vs. -log10(p-values) further illustrate this point
27  (Figure 8).

28  **Figure 6: *KnockoffScreen* application to the Alzheimer's Disease Sequencing Project (ADSP) data to identify variants**
29  **associated with the Alzheimer's Disease.** The top-left panel presents the Manhattan plot of p-values (truncated at $10^{-20}$ for clear
30  visualization) from the conventional association testing with Bonferroni adjustment ($p < 0.05$/number of tested windows) for
31  FWER control. The bottom-left panel presents the Manhattan plot of *KnockoffScreen* with target FDR at 0.1. The right panel
32  presents a heatmap that shows stratified p-values (truncated at $10^{-10}$ for clear visualization) of all loci passing the FDR=0.1
33  threshold, and the corresponding Q-values that already incorporate correction for multiple testing. The loci are shown in descending
34  order of the knockoff statistics. For each locus, the p-values of the top associated single variant and/or window are shown indicating
35  whether the signal comes from a single variant, a combined effect of common variants or a combined effect of rare variants. The
36  names of those genes previously implicated by GWAS studies are shown in bold (names were just used to label the region and may
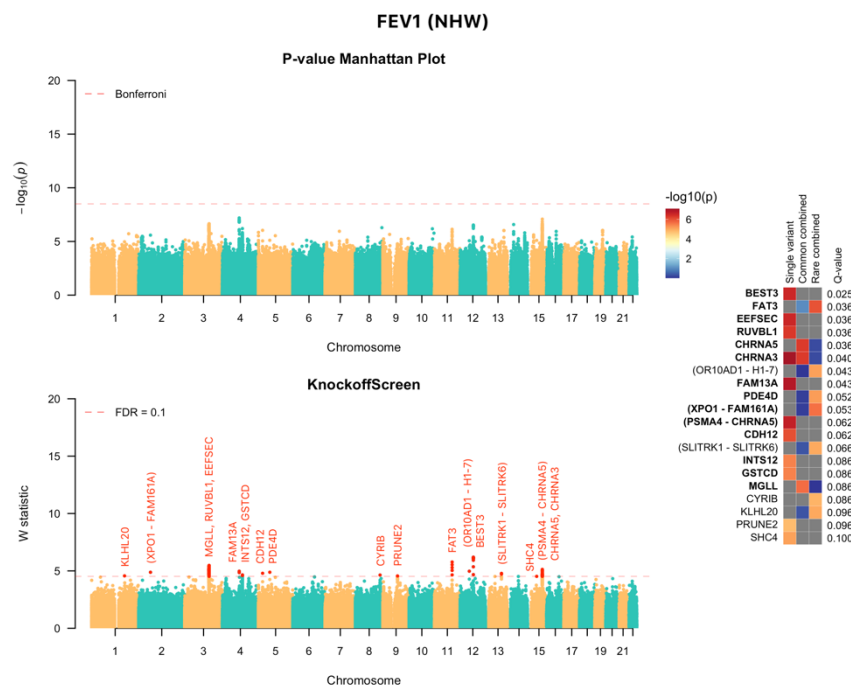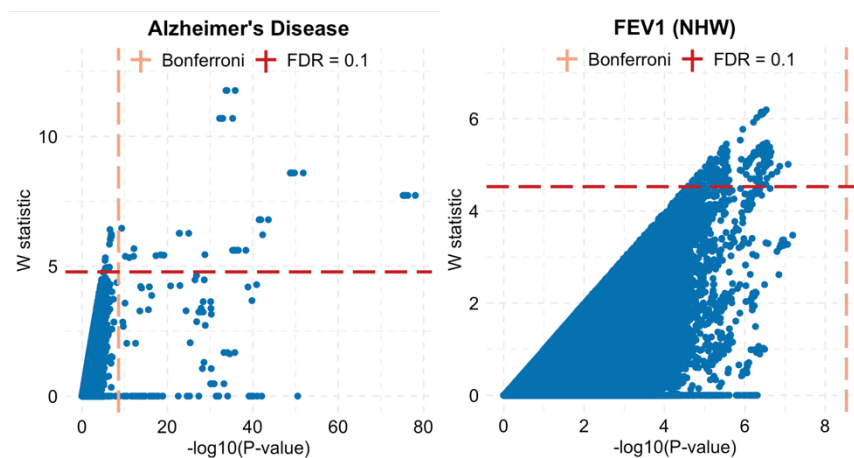37  not represent causative gene in the region).

**Figure 7:** *KnockoffScreen* **application to the COPDGene study in TOPMed to identify variants associated with FEV$_1$ in Non Hispanic White (NHW).** The top-left panel presents the Manhattan plot of p-values from the conventional association testing with Bonferroni adjustment ($p < 0.05$/number of tested windows) for FWER control. The bottom-left panel presents the Manhattan plot of *KnockoffScreen* with target FDR at 0.1. The right panel presents a heatmap that shows stratified p-values of all loci passing the FDR=0.1 threshold, and the corresponding Q-values that already incorporate correction for multiple testing. The loci are shown in descending order of the knockoff statistics. For each locus, the p-values of the top associated single variant and/or window are shown indicating whether the signal comes from a single variant, a combined effect of common variants or a combined effect of rare variants. The names of those genes previously implicated by GWAS studies are shown in bold (names were just used to label the region and may not represent causative gene in the region).

1 **Figure 8: Scatter plot of genome-wide W statistic vs. -log10(p-value).** Each dot represents one variant/window. The dashed lines
2 show the significance thresholds defined by Bonferroni correction (for p-values) and by false discovery rate (FDR; for W statistic).
3 The p-values are from the conventional association testing described in the main text.



4

## 5 Discussion

6 In summary, we propose a computationally efficient algorithm, *KnockoffScreen*, for the identification of
7 putative causal loci in whole-genome sequencing studies based on the knockoff framework. This
8 framework guarantees the FDR control at a desired level under general dependence structure, and has
9 appealing properties relative to conventional association tests, including a reduction in LD-contaminated
10 associations and false positive associations due to unadjusted population stratification. Through
11 applications to two whole-genome sequencing studies for Alzheimer's disease, COPD and lung function
12 phenotypes we demonstrate the ability of the approach to identify more significant associations, many of
13 which have been identified in previous GWAS studies, with sample sizes orders of magnitude larger than
14 the ones considered here. As sample sizes for whole-genome sequencing studies continue to increase,
15 *KnockoffScreen* can help discover more risk loci with even more stringent FDR thresholds.

16 In *KnockoffScreen*, we choose to control FDR at the nominal level. Our analyses of data from ADSP and
17 COPDGene show that our method compared with conventional association tests leads to significantly more
18 discoveries. The majority of the single variant signals that were found significant at FDR 0.1 have been
19 associated with AD or COPD-related phenotypes respectively in the GWAS catalog (87.5% for AD, 81.8%
20 for $FEV_1$ and 69.2% for $FEV_1$/FVC), supporting our claim that the FDR control in *KnockoffScreen* is able
21 to replicate previously discovered loci in GWAS studies with sample sizes much larger than those used
22 here. Furthermore, *KnockoffScreen* identified a set of new discoveries driven by the combined effects of
23 multiple common/rare variants. The results demonstrate that controlling FDR is an appealing strategy when
24 there are potentially many discoveries to be made as in genetic association studies for highly polygenic
25 traits, the dependence structure is local, and the investigators are willing to accept a rigorously defined
26 small fraction of false positives in order to substantially increase the total number of true discoveries. We
27 note that the choice of target FDR should be defined rigorously and interpreted appropriately. For example,
28 loci identified at a liberal FDR threshold (e.g. 0.3 as in Iossifoy et al.[9]) can be useful for enrichment and
29 pathway analyses; our analyses of data from ADSP and COPDGene used FDR=0.1 for identifying putative
30 causal loci. As large-scale whole-genome sequencing data become increasingly available, one will be able
31 to apply *KnockoffScreen* with a lower, more stringent FDR threshold (e.g. 0.01 or 0.05).

32 The model-X knockoff framework underlying *KnockoffScreen* makes our approach robust to violations of
33 model assumptions. Specifically, by imposing a model on genetic variants ($G_i$) instead of on the conditional
34 distribution of the outcome given the variants (distribution of $Y_i | G_i$), the FDR control is guaranteed even
35 when the model for $Y_i | G_i$ is mis-specified. We do however need to construct a valid synthetic cohort $\tilde{G}_i$'s
36 such that the exchangeability conditions are satisfied, and define a test statistic with the sign-flip property

1 (i.e. the effect of swapping a variant with its knockoff is only a sign flip of the corresponding test statistic).
2 This robustness feature is particularly useful for genetic studies of complex traits, as the underlying genetic
3 model is unknown, and it is difficult to evaluate whether a model is appropriate for describing the
4 relationship between the trait and the variants.

5 There is some limited work on controlling the FWER within the knockoff framework using a single
6 knockoff[45]. One obstacle for its application is that it only allows controlling for $k$-FWER at significance
7 level $\alpha$ (the probability of making at least $k$ false rejections) where $k$ or $\alpha$ has to be relatively large in order
8 to detect any association. Therefore, it cannot be directly applied to control the conventional FWER ($k =$
9 $1, \alpha = 0.05$) without further modifications. Although our proposed multiple knockoffs method has the
10 potential to be extended to control the FWER, we estimated that about 20 knockoffs are necessary to achieve
11 the conventional FWER control. This leads to additional computational burden that will need to be
12 overcome in order to become scalable to the large-scale genetic data.

13 In addition to controlling FDR, our approach contrasts to conventional association testing methods in that
14 it naturally helps prioritize the underlying causal variants, a property that usually requires a second stage
15 conditional analysis or statistical fine mapping[46]. It also helps separate causal effects from shadow effects
16 of significant variants nearby. This property can help distinguish effects due to common causal variants or
17 rare causal variants at the same locus due to LD, by applying *KnockoffScreen* to common/rare variants
18 separately. Overall, *KnockoffScreen* serves as a powerful and efficient method that attempts to unify
19 association testing and statistical fine mapping. However, similar to statistical fine-mapping methods that
20 only leverage LD to fine-map a complex trait, it remains challenging to fully distinguish highly correlated
21 variants. As we discussed in the Methods section, *KnockoffScreen* currently detects clusters of tightly linked
22 variants, without removing any variants that are potentially causal. In the future, we may consider using
23 functional genomics data to further improve the ability of *KnockoffScreen* to identify causal variants among
24 highly correlated ones.

25 Unlike existing knockoff methods for genetic data that define coefficients in a LASSO regression as the
26 importance score[15,16], *KnockoffScreen* directly uses transformed p-values as importance score. This leads
27 to another appealing property of *KnockoffScreen*, namely it can serve as a wrapper method that can flexibly
28 utilize p-values from any existing or future association testing methods to achieve the benefits proposed
29 here. For example, the current implementation of *KnockoffScreen* calculates importance score using an
30 ACAT type test to aggregate several recent advances for rare-variant analysis. To extend its application to
31 studies with large unbalanced case-control ratios or sample relatedness, one can apply methods like
32 SAIGE[47] to calculate p-values for the original cohort and the synthetic cohort generated by *KnockoffScreen*,
33 and then apply the same knockoff filter for variable selection. Moreover, recent studies have demonstrated
34 that multivariate models have many advantages over marginal association testing, including improved
35 power by reducing the residual variation and better control of population stratification[15]. *KnockoffScreen* is
36 able to integrate tests from multivariate models (e.g. BOLT-LMM and its extension to window-based
37 analysis of sequencing data).

38 Meta-analyses are important in allowing the integration of results from multiple whole-genome sequencing
39 studies without sharing individual level data. Several methods have been proposed for meta-analysis of
40 single variant tests for common variants or "set" based (e.g. window based) tests for rare variants[48-50]. Those
41 methods integrate summary statistics from each individual cohort, such as p-values or score statistics, and
42 then compute a combined p-value for each genetic variant or each window for a meta-analysis. As we
43 discussed, *KnockoffScreen* can also directly utilize p-values from existing methods for meta-analysis. We
44 have discussed the detailed procedure in the Methods section.

45 Variable selection based on knockoff procedure depends on the random sampling of knockoff features
46 $\{\tilde{G}_{ij}\}_{1 \le j \le p}$. Although FDR control is guaranteed, the randomness may lead to slightly different feature
47 statistics and selection of slightly different subsets of variants. We propose a stable inference procedure

1  integrating multiple knockoffs that significantly improves the stability and reproducibility of the results
2  compared with state-of-the-art knockoff methods as discussed in the Method section.

3  We have demonstrated that the proposed sequential knockoff generator is significantly faster than existing
4  alternatives. Besides the generation of knockoff features, another source of computational burden is the
5  calculation of the importance score (p-value for a window). The total CPU time is 7,616 hours for the ADSP
6  data analysis (15.2 hours with 500 cores) and 14,274 hours for the COPDGene data analysis (28.5 hours
7  with 500 cores). The calculation of p-values in the current analysis is time consuming because of the
8  comprehensive inclusion of many different functional annotations. Specifically, for each window, there are
9  in total 29 tests being implemented for the original genetic variants and each of their five knockoffs, leading
10  to a total of 29*6=174 p-value calculations per window. If computational resources are limited, using a
11  limited number of functional annotations can substantially reduce the computing time. In addition, several
12  methods have been proposed in recent years to use state-of-the-art optimization strategies for scalable
13  association testing for large scale datasets with thousands of phenotypes in large biobanks.[51-53] By directly
14  utilizing p-values from those association testing methods, *KnockoffScreen* can scale up to biobank sized
15  datasets at a comparable computational efficiency.

16  Despite the aforementioned advantages, *KnockoffScreen* has some limitations related to underlying
17  modeling assumptions needed to improve the computational efficiency of the multiple knockoff generation
18  and calculation of the feature importance scores. In particular, the implemented feature importance scores
19  rely on computing p-values from a marginal model (e.g. single variant score test, burden test or SKAT) or
20  a partly multivariate model (BOLT-LMM and its extension to window-based analysis of sequencing data).
21  We made this choice of feature importance score due to its flexibility to integrate state-of-the-art tests for
22  sequencing studies, but we recognize that a fully multivariate model as implemented in Sesia et al.[15] can be
23  more powerful. In addition, the knockoff generator used in *KnockoffScreen* assumes a linear approximation
24  model based on unphased genotype dosage data. This model is well motivated based on the sequential
25  model to generate knockoff features, and the approximate multivariate normal model for the genotype data
26  commonly used in the genetic literature. Additionally, it is computationally efficient relative to existing
27  knockoff generation methods. We acknowledge that relative to a generative model like HMM it is less
28  interpretable. More complex models for discrete genotype values that can also account for non-linear effects
29  among genetic variants could be of interest in future work.

30

18

## Methods

**Sequential model to generate model-X knockoff features.** We propose a computationally efficient sequential model to generate knockoff features $\widetilde{G}$ that leverages local linkage disequilibrium structure. Our method is an extension of the general sequential conditional independent pairs (SCIP) approach in Candès et al. (2018)[14].

---

**Algorithm 1** Sequential Conditional Independent Pairs (Single Knockoff)

---

$j = 1$
while $j \leq p$ do
    Sample $\tilde{G}_j$ independently from $\mathcal{L}\left(G_j | \boldsymbol{G}_{-j}, \widetilde{\boldsymbol{G}}_{\boldsymbol{1}:(\boldsymbol{j-1})}\right)$
    $j = j + 1$
**end**

---

where $\boldsymbol{G}_{-j}$ denotes all genetic variants except for the $j$-th variant; $\mathcal{L}\left(G_j | \boldsymbol{G}_{-j}, \widetilde{\boldsymbol{G}}_{\boldsymbol{1}:(\boldsymbol{j-1})}\right)$ is the conditional distribution of $G_j$ given $\boldsymbol{G}_{-j}$ and $\widetilde{\boldsymbol{G}}_{\boldsymbol{1}:(\boldsymbol{j-1})}$. Candès et al. showed that knockoffs generated by this algorithm satisfy the exchangeability condition, and they lead to a guaranteed FDR control[14]. Intuitively, the exchangeability condition can be described as follows: if one swaps any subset of variants and their synthetic counterpart, the joint distribution (LD structure etc.) does not change. They also noted that the ordering in which knockoffs are created does not affect the exchangeability property and equally valid constructions may be obtained by looping through an arbitrary ordering of the variants. Although the SCIP method represents a general knockoff generator, the conditional distribution at each iteration depends on all genetic variants in the study, which can be very difficult or impossible to compute in practice. We draw inspiration from Markov models for sequence data to consider the genetic sequence as a Markov chain with memory, such that

$$\mathcal{L}\left(G_j | \boldsymbol{G}_{-j}\right) = \mathcal{L}\left(G_j | \boldsymbol{G}_{\boldsymbol{k} \in B_j}\right), \quad (3)$$

where the index set $B_j$ defines a subset of genetic variants "near" the $j$-th variant, which we will define later. Furthermore, by noting that the correlation among genetic variables approximately exhibits a block diagonal structure[54], under certain model assumptions which will be specified in the Appendix, we have

$$\mathcal{L}\left(G_j | \boldsymbol{G}_{-j}, \widetilde{\boldsymbol{G}}_{\boldsymbol{1}:j-1}\right) = \mathcal{L}\left(G_j | \boldsymbol{G}_{\boldsymbol{k} \in B_j}, \widetilde{\boldsymbol{G}}_{\boldsymbol{1} \leq k \leq j-1, k \in B_j}\right). \quad (4)$$

To generate knockoff features from $\mathcal{L}\left(G_j | \boldsymbol{G}_{\boldsymbol{k} \in B_j}, \widetilde{\boldsymbol{G}}_{\boldsymbol{1} \leq k \leq j-1, k \in B_j}\right)$, we assume a semiparametric model

$$G_j = g\left(\boldsymbol{G}_{\boldsymbol{k} \in B_j}, \widetilde{\boldsymbol{G}}_{\boldsymbol{1} \leq k \leq j-1, k \in B_j}\right) + \varepsilon_j, \quad (5)$$

where $\varepsilon_j$ is a random error term, $E\left(\varepsilon_j \big| \boldsymbol{G}_{\boldsymbol{k} \in B_j}, \widetilde{\boldsymbol{G}}_{\boldsymbol{1} \leq k \leq j-1, k \in B_j}\right) = 0$. We consider $g(\cdot)$ to be parametric as follows,

$$g\left(G_{ij} | \boldsymbol{G}_{\boldsymbol{k} \in B_j}, \widetilde{\boldsymbol{G}}_{\boldsymbol{1} \leq k \leq j-1, k \in B_j}\right) = \alpha + \sum_{k \neq j, k \in B_j} \beta_k G_{ik} + \sum_{k \leq j-1, k \in B_j} \gamma_k \tilde{G}_{ik}, \quad (6)$$

and will explain in detail when such a linear form is an appropriate model in the Appendix. We estimate $(\alpha, \boldsymbol{\beta}, \boldsymbol{\gamma})$ by minimizing the mean squared loss. Let $\widehat{\boldsymbol{G}}_j = \hat{\alpha} + \sum_{k \neq j, k \in B_j} \hat{\beta}_k \boldsymbol{G}_k + \sum_{k \leq j-1, k \in B_j} \hat{\gamma}_k \widetilde{\boldsymbol{G}}_k$. We calculate the residual $\widehat{\boldsymbol{\varepsilon}}_j = \boldsymbol{G}_j - \widehat{\boldsymbol{G}}_j$ and its permutation $\widehat{\boldsymbol{\varepsilon}}_j^*$, and then define the knockoff feature for $\boldsymbol{G}_j$ to be $\widetilde{\boldsymbol{G}}_j = \widehat{\boldsymbol{G}}_j + \widehat{\boldsymbol{\varepsilon}}_j^*$. This permutation-based algorithm is particularly designed to generate knockoff features for rare genetic variants in sequencing studies, whose distribution is highly skewed and zero-inflated. We note that the algorithm does not generate categorical variables in {0,1,2}. Instead, it generates continuous variables to mimic genotype dosage value, making it more robust for rare variants. In addition, we evaluated

19

1    a multinomial logistic regression model for generating categorical knockoffs. We found that the conditional
2    mean of a rare variant can be extremely small, and it is very likely to generate knockoffs with all 0 values
3    where statistical inference cannot be applied. We show in simulation studies that existing knockoff
4    generators, such as the second-order model-X knockoffs proposed by Candès et al.[14] and knockoffs for
5    HMM proposed by Sesia et al.[15,16], do not control FDR for rare variant analysis based on the feature score
6    considered in this paper (Figure 2). In Figures S7 and S8, we present an additional comparison between the
7    proposed method and HMM-based knockoff generators (S=12 and S=50), stratified by allele frequency. As
8    shown, the proposed method generates knockoff versions for rare variants with better exchangeability with
9    the original variants compared with the HMM model. That is, the correlation coefficients are closer to those
10    for the original variants for *KnockoffScreen* compared to HMM (bottom panel, the dots are mostly above
11    the diagonal line). One plausible explanation is that the application of HMM to whole genome sequencing
12    data requires accurate phased data for rare variants, which itself is a challenging task and also an active
13    research area.

14    We discuss now in detail how we define $B_j$ while taking into account the linkage disequilibrium (LD)
15    structure in the neighborhood of $j$. Let $r_{jk}$ be the sample correlation coefficient between variants $j$ and $k$.
16    We define $B_j$ to include "$K$-nearest" genetic variants within a 200kb window (+/-100kb from the target
17    variant)[55] using $|r_{jk}|$ as a similarity measure. The choice of the window size aims to balance accurate
18    modeling of local LD structure and computational efficiency. The choice of $K$ is to ensure that
19    $P\left(G_j | \boldsymbol{G}_{k \in B_j}, \widetilde{\boldsymbol{G}}_{1 \le k \le j-1, k \in B_j}\right)$ accurately mimics the joint distribution $P\left(G_j | \boldsymbol{G}_{-j}, \widetilde{\boldsymbol{G}}_{1:j-1}\right)$ and to avoid
20    overfitting. We adopt the theoretical result for regression analysis with diverging number of covariates and
21    choose to include top $K$ variants with $|r_{jk}| > 0.05$ up to $K = n^{1/3}$, which ensures that the coefficient
22    estimations achieve asymptotic normality[56].

23    We note that the sequential model is flexible enough and we could consider other supervised learning
24    techniques like Lasso, support vector regression and artificial neural networks. However, since the auto-
25    regressive model is fitted iteratively for every variant in the genome, these methods require cross-validation
26    at each variant level which is computationally not applicable at genome-wide scale.

27    **Multiple sequential knockoffs to improve power and stability.** Inference based on single knockoff is
28    limited by the detection threshold $[\frac{1}{q}]$, which is the minimum number of independent rejections needed in
29    order to detect any association. For example, in scenarios where the signal is sparse (<10 independent true
30    associations) in the target region or across the genome, inference based on a single knockoff has very low
31    power to detect any association with target FDR 0.1. Another limitation of the single knockoff is its
32    instability. Since the knockoff sample is random, running the knockoff procedure multiple times may lead
33    to different selected sets of features. The idea of constructing multiple knockoffs was first discussed by
34    Barber and Candès[13] and Candès et. al.[14], and further studied in detail by Gimenez and Zou[30]. However,
35    current methods are not applicable to rare variants and not scalable to whole genome sequencing data.

36    We extend the above SCIP based knockoff generator procedure to multiple knockoffs ($M$ is the total number
37    of knockoffs), as follows.

---

**Algorithm 2** Sequential Conditional Independent Tuples (Multiple Knockoffs)

$j = 1$
while $j \le p$ do
    Sample $\tilde{G}_j^1, \cdots, \tilde{G}_j^M$ independently from $\mathcal{L}\left(G_j | \boldsymbol{G}_{-j}, \widetilde{\boldsymbol{G}}_{1:j-1}^1, \cdots, \widetilde{\boldsymbol{G}}_{1:j-1}^M\right)$
    $j = j + 1$
**End**

---

38    Gimenez and Zou[30] proposed this general algorithm and proved that the knockoffs generated by this
39    algorithm satisfy the extended exchangeability condition (see Appendix for precise definition and proof).

20

1    Based on this general algorithm, we extend our previous sequential model to this setting to estimate $\widehat{G}_j =$

2    $\hat{\alpha} + \sum_{k \neq j, k \in B_j} \hat{\beta}_k G_k + \sum_{1 \leq m \leq M} \sum_{k \leq j-1, k \in B_j} \hat{\gamma}_k^m \widetilde{G}_k^m$. We calculate the residual $\hat{\varepsilon}_j = G_j - \widehat{G}_j$ and its $M$

3    permutations $\hat{\varepsilon}_j^{*1}, \dots, \hat{\varepsilon}_j^{*M}$, and then define the knockoff feature for $G_j$ to be $\widetilde{G}_j^m = \widehat{G}_j + \hat{\varepsilon}_j^{*m}$.

4    **Knockoff filter to define the threshold $\tau$ for FDR control.** For single knockoff, we follow the result

5    derived by Candès et al.[14] to define the feature statistic as $W_{\Phi_{kl}} = T_{\Phi_{kl}} - \widetilde{T}_{\Phi_{kl}}$ where $T_{\Phi_{kl}} = -\log_{10} p_{\Phi_{kl}}$

6    and $\widetilde{T}_{\Phi_{kl}} = -\log_{10} \tilde{p}_{\Phi_{kl}}$ and

7
$$\tau = \min\left\{t > 0: \frac{1 + \#\{\Phi_{kl}: W_{\Phi_{kl}} \leq -t\}}{\#\{\Phi_{kl}: W_{\Phi_{kl}} \geq t\}} \leq q\right\}, \quad (7)$$

8    where "#" denote the number of elements in the set; $q$ is the target FDR level. We select all windows with

9    $W_{\Phi_{kl}} > \tau$. For multiple knockoffs, we modify the result in Gimenez and Zou[30] and define

10
$$W_{\Phi_{kl}} = \left(T_{\Phi_{kl}} - \underset{1 \leq m \leq M}{\text{median}} T_{\Phi_{kl}}^m\right) I_{T_{\Phi_{kl}} \geq \underset{1 \leq m \leq M}{\max} T_{\Phi_{kl}}^m}, \quad (8)$$

11    and

12
$$\tau = \min\left\{t > 0: \frac{\frac{1}{M} + \frac{1}{M}\#\{\Phi_{kl}: \kappa_{\Phi_{kl}} \geq 1, \tau_{\Phi_{kl}} \geq t\}}{\#\{\Phi_{kl}: \kappa_{\Phi_{kl}} = 0, \tau_{\Phi_{kl}} \geq t\}} \leq q\right\}, \quad (9)$$

13    where $T_{\Phi_{kl}}^m = -\log p_{\Phi_{kl}}^m$; $I.$ is an indicator function, $I_{T_{\Phi_{kl}} \geq \underset{1 \leq m \leq M}{\max} T_{\Phi_{kl}}^m} = 1$ if $T_{\Phi_{kl}} \geq \underset{1 \leq m \leq M}{\max} T_{\Phi_{kl}}^m$ and 0

14    otherwise; $\kappa_{\Phi_{kl}} = \underset{0 \leq m \leq M}{\arg\max} T_{\Phi_{kl}}^m$ denote the index of the original (denoted as 0) or knockoff feature that has

15    the largest importance score; $\tau_{\Phi_{kl}} = T_{\Phi_{kl}}^{(0)} - \underset{1 \leq m \leq M}{\text{median}} T_{\Phi_{kl}}^{(m)}$ denote the difference between the largest

16    importance score and the median of the remaining importance scores. It reduces to the knockoff filter for

17    single knockoff when $M = 1$. Essentially, $W_{\Phi_{kl}} > \tau$ selects windows where the original feature has higher

18    importance score than any of the $M$ knockoffs (i.e. $\kappa_{\Phi_{kl}} = 0$), and the gap with the median of knockoff

19    importance score is above some threshold.

20    We note that this definition of feature statistic and knockoff filter is a modified version of that proposed by

21    Gimenez and Zou[30], where they considered the maximum instead of the median of the knockoff importance

22    scores,    i.e.    $\kappa_{\Phi_{kl}} = \underset{0 \leq m \leq M}{\arg\max} T_{\Phi_{kl}}^m$ ,    $\tau_{\Phi_{kl}} = T_{\Phi_{kl}}^{(0)} - \underset{1 \leq m \leq M}{\max} T_{\Phi_{kl}}^{(m)}$    and    $W_{\Phi_{kl}} = \big(T_{\Phi_{kl}} - $

23    $\underset{1 \leq m \leq M}{\max} T_{\Phi_{kl}}^m\big) I_{T_{\Phi_{kl}} \geq \underset{1 \leq m \leq M}{\max} T_{\Phi_{kl}}^m}$. To improve stability and reproducibility of knockoff based inference, we

24    change $\tau_{\Phi_{kl}}$ from $T_{\Phi_{kl}}^{(0)} - \underset{1 \leq m \leq M}{\max} T_{\Phi_{kl}}^{(m)}$ to $T_{\Phi_{kl}}^{(0)} - \underset{1 \leq m \leq M}{\text{median}} T_{\Phi_{kl}}^{(m)}$ . The modified method reduces the

25    randomness coming from sampling knockoff features given the fact that sample median has much smaller

26    variation than each individual sample or the sample maximum.

27    **Knockoff Q-value.** The Q-value in statistics is similar to the well-known p-value, except that it measures

28    significance in terms of the FDR[57] rather than the FWER and already incorporates correction for multiple

29    testing. For multiple hypothesis testing, a general mathematical definition of the Q-value for a null

30    hypothesis is the minimum FDR that can be attained when all tests showing evidence against the null

31    hypothesis at least as strong as the current one are declared as significant[58]. For example, the Q-value for

32    usual FDR control based on ordered p-values can be estimated by,

33
$$q = \min_{t \geq p} \widehat{\text{FDR}}(t), \quad (10)$$

34    where $p$ is the p-value of the hypothesis under consideration and $\widehat{\text{FDR}}(t)$ is the estimated FDR if we are to

35    reject all tests with p-values less than $t$. In order to introduce a more informative and interpretable measure

1    of significance for the top signals, we extend the Q-value framework for the usual FDR control to the
2    knockoffs based case. The proposed Q-value combines the information from both feature importance
3    statistics $W_{\Phi_{kl}}$ and the threshold $\tau$. It also makes results comparable even we choose different feature
4    importance statistics across multiple runs. By definition, we shall see that selecting windows with $q_\Phi < q$,
5    where $q$ is the target FDR, is equivalent to the aforementioned knockoff filter which selects those with
6    $W_\Phi > \tau$.

7    For single knockoff, we define the Q-value for window $\Phi$ with feature statistic $W_\Phi > 0$ as,

$$q_\Phi = \min_{t \leq W_\Phi} \frac{1 + \#\{\Phi_{kl} : W_{\Phi_{kl}} \leq -t\}}{\#\{\Phi_{kl} : W_{\Phi_{kl}} \geq t\}}, \qquad (11)$$

9    where $\frac{1 + \#\{\Phi_{kl} : W_{\Phi_{kl}} \leq -t\}}{\#\{\Phi_{kl} : W_{\Phi_{kl}} \geq t\}}$ is an estimate of the proportion of false discoveries if we are to select all windows
10    with feature statistic greater than $t > 0$, referred to as the knockoff estimate of FDR[13]. For window $\Phi$ with
11    feature statistic $W_\Phi \leq 0$, we define $q_\Phi = 1$ and the window will never be selected. For multiple knockoffs,
12    we define the Q-value for window $\Phi$ with statistics $\kappa_\Phi = 0$ and $\tau_\Phi$ as

$$q_\Phi = \min_{t \leq \tau_\Phi} \frac{\frac{1}{M} + \frac{1}{M}\#\{\Phi_{kl} : \kappa_{\Phi_{kl}} \geq 1, \tau_{\Phi_{kl}} \geq t\}}{\#\{\Phi_{kl} : \kappa_{\Phi_{kl}} = 0, \tau_{\Phi_{kl}} \geq t\}}, \qquad (12)$$

14    where $\frac{\frac{1}{M} + \frac{1}{M}\#\{\Phi_{kl} : \kappa_{\Phi_{kl}} \geq 1, \tau_{\Phi_{kl}} \geq t\}}{\#\{\Phi_{kl} : \kappa_{\Phi_{kl}} = 0, \tau_{\Phi_{kl}} \geq t\}}$ is an estimate of the proportion of false discoveries if we are to select all
15    windows with feature statistic $\kappa_{\Phi_{kl}} = 0, \tau_{\Phi_{kl}} \geq t$, which is our extension of the knockoff estimate of FDR
16    to multiple knockoffs. For window $\Phi$ with $\kappa_\Phi \neq 0$, we again define $q_\Phi = 1$ and the window will never be
17    selected.

18    **Choice of windows for genome-wide screening.** *KnockoffScreen* considers windows with different sizes
19    (1bp, 1kb, 5kb, 10kb) across the genome, with half of each window overlapping with adjacent windows at
20    the same window size. This choice of windows is similar to the scan statistic framework, WGScan, for
21    whole-genome sequencing data[18]. It is also similar to that in *KnockoffZoom* proposed by Sesia et al.[15] for
22    GWAS data where they also consider windows of different sizes; for each fixed window size the windows
23    are non-overlapping but smaller windows are fully nested within larger windows. We theoretically prove
24    the FDR control using the proposed statistic in the Appendix for nonoverlapping windows; however, the
25    theoretical justification for the more general setting of overlapping windows remains an open question. For
26    the proposed choice of overlapping windows, we demonstrate via empirical simulation studies that the FDR
27    is well controlled (Figure 2) as window overlapping is a local phenomenon.

28    ***KnockoffScreen* improves stability and reproducibility of knockoff-based inference.** We conducted
29    simulation studies to compare *KnockoffScreen* with single knockoff approach, and the multiple knockoffs
30    approach proposed by Gimenez and Zou, referred to as MK-Maximum[30].

31    We designed these simulations to mimic the real data analysis of ADSP. For each replicate, we randomly
32    drew 1,000 variants, including both common and rare variants, from the 200kb region near gene *APOE*
33    (chr19: 44905796-44909393). We set 1.25% variants to be causal, all within a 5kb signal window (similar
34    to the size of *APOE*) and then simulated a dichotomous trait as follows

$$g(\mu_i) = \beta_0 + X_{i1} + \beta_1 g_1 + \cdots + \beta_s g_s,$$

36    where $g(x) = \log(\frac{x}{1-x})$ and $\mu_i$ is the conditional mean of $Y_i$; $\beta_0$ is chosen such that the prevalence is 10%.
37    We set the effect $\beta_j = 0.7|\log_{10} m_j|$, where $m_j$ is the MAF for the $j$-th variant. Given the same genotype
38    and phenotype data, we first generated 100 knockoffs. Then we repeatedly drew five knockoffs randomly
39    among them for 100 replicates. For each replicate, we scanned the regions with candidate window sizes

22

1  (1bp, 1kb, 5kb, 1kb) using *KnockoffScreen*, the multiple knockoffs feature statistic based on sample
2  maximum by Gimenez and Zou[30], and the single knockoff method. For a fair comparison, we adopted the
3  same tests implemented in *KnockoffScreen* to calculate the p-value for all comparison methods. We
4  calculated the variation of feature statistic $W_{\Phi_{kl}}$ for each window (stability) and the frequency with which
5  each causal window is selected (reproducibility) over 100 replicates. We present the results in Figure 9.

6  In the left panel, we observed that *KnockoffScreen* has significantly smaller variation in feature statistic
7  $W_{\Phi_{kl}}$ than the other two comparison methods. We note that the method based on sample maximum, MK-
8  Maximum, exhibits comparable and sometimes even larger variation than the method based on single
9  knockoff. In the mid panel, we observed that *KnockoffScreen* has a higher chance (~0.94) to replicate
10 findings across different knockoff replicates compared to MK-Maximum (~0.74-0.83) and single knockoff
11 (~0.43). This improvement is further demonstrated in the right panel, where we show that *KnockoffScreen*
12 exhibits smaller variation in feature statistics for the causal windows, resulting in higher reproducibility.
13 The significantly lower reproducibility rate for single knockoffs relative to MK-Maximum is presumably
14 due to its higher detection threshold because it exhibits similar level of variation as MK-Maximum for the
15 causal windows.

16 **Figure 9: Simulation studies to evaluate the stability and reproducibility of different knockoff procedures.** Different colors
17 indicate different knockoff procedures: *KnockoffScreen,* single knockoff and MK – Maximum (the multiple knockoff method based
18 on the maximum statistic proposed by Gimenez and Zou[30]). All three methods are based on the same knockoff generator proposed
19 in this paper for a fair comparison. The stability is quantified as the variation of $\tau_{\Phi_{kl}}$ across 100 replicates due to randomly sampling
20 knockoffs for a given data (left and right panels). The reproducibility is quantified as the frequency of a causal window being
21 selected across 100 replicates.



22

23 **Practical strategy for tightly linked variants.** Variants residing in short genetic regions can be in
24 moderate to high LD. Although the knockoff method helps to prioritize causal variants over associations
25 due to low/moderate LD, strong correlations can make it difficult or impossible to distinguish the causal
26 genetic variants from their highly correlated variants (see also Sesia et al.[16]). In fact, the knockoff method
27 will rank all those highly correlated variants lower, which diminishes the power if causal variants exist (see
28 below for a concrete example). We are primarily interested in the identification of relevant clusters of tightly
29 linked variants, rather than individual variants. To address this issue, we propose a practical solution by
30 slightly modifying $B_j$. The resulting algorithm improves the power to detect clusters of tightly linked
31 variants, without removing any variants that are potentially causal.

32 Specifically, we create a hierarchical clustering dendrogram using $|r_{jk}|$ as a similarity measure and define
33 clusters by $|r_{jk}| > 0.75$, such that variants from two different clusters do not have a correlation greater than
34 0.75. To generate the knockoff feature for the $j$-th variant, we exclude variants from $B_j$ that are in the same
35 cluster. For example, let $G_1$, $G_2$ and $G_3$ be three genetic variants; $G_1$ and $G_2$ are tightly linked with $|r_{12}| >$
36 0.75 . The standard knockoff procedure will generate $\tilde{G}_1$ based on $P(G_1|G_2, G_3)$ , $\tilde{G}_2$ based on
37 $P(G_2|G_1, G_3, \tilde{G}_1)$. Since $G_1$ and $G_2$ are highly correlated, $\tilde{G}_1 \approx G_1, \tilde{G}_2 \approx G_2$ and there will be no power to
38 detect $G_1$ or $G_2$ even if one of them is causal. To improve the power, our modified algorithm simultaneously

23

1 generates $\tilde{G}_1$ and $\tilde{G}_2$ based on a joint distribution $P(G_1, G_2|G_3)$ by first estimating the conditional means
2 and then permuting the residuals jointly. This avoids the situation of $\tilde{G}_1$ and $\tilde{G}_2$ being identical to $G_1$ and
3 $G_2$ because $G_1$ ($G_2$) is excluded from the generation of $\tilde{G}_2(\tilde{G}_1)$ . Thus both $G_1$ and $G_2$ can be detected as a
4 cluster. The idea is similar to that of group-wise exchangeable knockoffs proposed by Sesia et al.[15]. We
5 further discuss limitations and some alternative approaches in the Discussion section.

6 **Computational efficiency of the knockoff generator.** We estimate the computational complexity of our
7 proposed method for each variant $j$ as $O(nL) + O(LlogL) + O(n(K + MK)^2 + (K + MK)^3) = O(n)$,
8 where $n$ is the sample size; $L$ is a predefined constant for the length of the nearby region; $K$ is the number
9 of variants in the defined set $B_j$, which is bounded by the predefined constant $L$; $M$ is a predetermined
10 constant for the number of knockoffs. $O(nL)$ is for calculating the correlation between variant $j$ and
11 variants in the nearby region; $O(LlogL)$ is for the hierarchical clustering; $O(n(K + MK)^2 + (K + MK)^3)$
12 is for fitting the conditional auto-regressive model. Since we iteratively generate the knockoff for every
13 variant, we estimate the complexity of our proposed method for all variants as $O(np)$, where $p$ is the
14 number of genetic variants. We note that the genotype matrix $G$ is sparse for rare variants. Therefore, the
15 cost for calculation of correlation and hierarchical clustering can be drastically reduced. In addition, the
16 approach that we proposed to define $B_j$ ensures that $K$ is relatively small and this further reduces the
17 computational cost.

18 *KnockoffScreen* **allows meta-analysis of multiple cohorts.** Meta-analysis is a powerful approach that
19 enables integration of multiple cohorts for a larger sample size without sharing individual level data. Several
20 methods have been proposed for meta-analysis of single variant tests for common variants or set-based (e.g.
21 window based) tests for rare variants [48-50]. Those methods integrate summary statistics from each individual
22 cohort, such as p-values or score statistics, and then compute a combined p-value for each genetic variant
23 or each window for a meta-analysis. Since *KnockoffScreen* directly uses p-value as importance score, it can
24 flexibly incorporate the aforementioned methods for a meta-analysis. The meta-analysis procedure is
25 described as follows:

26 1. Generate knockoff features for each individual cohort.
27 2. Calculate summary statistics within each individual cohort for original data and knockoff data.
28 3. Apply existing meta-analysis methods to aggregate summary statistics to compute combined p-values
29     $p_{\Phi_{kl},combined}$ and $\tilde{p}_{\Phi_{kl},combined}$, for original data and knockoff data respectively.
30 4. Define $W_{\Phi_{kl}} = T_{\Phi_{kl}} - \tilde{T}_{\Phi_{kl}}$ where $T_{\Phi_{kl}} = -\log_{10} p_{\Phi_{kl},combined}$ and $\tilde{T}_{\Phi_{kl}} = -\log_{10} \tilde{p}_{\Phi_{kl},combined}$,
31     and apply *KnockoffScreen* to select putative causal variants. It naturally extends to multiple knockoffs
32     as described above.

33 **Single-region empirical power and FDR simulations.** We conducted empirical FDR and power
34 simulations. Each replicate consists of 10,000 individuals with genetic data on 1,000 genetic variants from
35 a 200kb region, simulated using the SKAT package. The SKAT haplotype dataset was generated using a
36 coalescent model (COSI), mimicking the linkage disequilibrium structure of European ancestry samples.
37 The simulations focus on both rare and common variants with minor allele frequency (MAF) <0.01
38 and >0.01 respectively. It has been discussed in Sesia et al.[16] that the false discovery proportion is difficult
39 to define if the method identifies a variant that is tightly linked with the causal variant. The analysis of
40 sequencing data targets different test units (set-based vs. single variant-based), further complicating the
41 FDR comparisons. We note that the simulations here focus on method comparison for locus discovery to
42 identify relevant clusters of tightly linked variants. Therefore, we simplify the simulation design in this
43 particular section to avoid difficulties in defining the FDR in the presence of strong correlations by keeping
44 one representative variant from each tightly linked cluster. Specifically, we applied hierarchical clustering
45 such that no two clusters have cross-correlations above a threshold value of 0.75 and then randomly choose
46 one representative variant from each cluster to be included in the simulation study.

1    We set 0.5% variants in the 200kb region to be causal, all within a 10kb signal window. Then we generated
2    the quantitative/dichotomous trait as follows:

3    $$\text{Quantitative trait: } Y_i = X_{i1} + \beta_1 g_1 + \cdots + \beta_s g_s + \varepsilon_i,$$

4    $$\text{Dichotomous trait: } g(\mu_i) = \beta_0 + X_{i1} + \beta_1 g_1 + \cdots + \beta_s g_s,$$

5    where $X_{i1} \sim N(0,1)$, $\varepsilon_i \sim N(0,3)$ and they are all independent; $(g_1, \ldots, g_s)$ are selected risk variants; $g(x) =$
6    $\log(\frac{x}{1-x})$ and $\mu_i$ is the conditional mean of $Y_i$; for dichotomous trait, $\beta_0$ is chosen such that the prevalence
7    is 10%. We set the effect $\beta_j = \frac{a}{\sqrt{2m_j(1-m_j)}}$, where $m_j$ is the MAF for the $j$-th variant. We define $a$ such
8    that the variance due to the risk variants, $\beta_1^2 var(g_1) + \cdots + \beta_s^2 var(g_s)$, is 0.05 for the simulations focusing
9    on common variants and 0.1 for the simulations focusing on rare variants. We scan the regions with
10   candidate window sizes (1bp, 1kb, 5kb, 10kb), and we consider several tests including the burden test,
11   dispersion test, and Cauchy combination test to aggregate burden, dispersion, and individual variant test
12   results (as discussed in the main text). This combined test is the method implemented in the *KnockoffScreen*
13   method. A window is considered causal if it contains at least one causal variant. For each replicate, the
14   empirical power is defined as the proportion of detected windows among all causal windows; the empirical
15   FDR is defined as the proportion of non-causal windows among all detected windows. We simulated 500
16   replicates and calculated the average empirical power and FDR.

17   **Genome-wide empirical power and FDR simulations in the presence of multiple causal loci.** We
18   conducted empirical FDR and power simulations using ADSP whole genome sequencing data, and
19   compared the proposed method with state-of-the-art tests for sequencing data analysis adjusted by
20   Bonferroni correction and Benjamini-Hochberg procedure for FDR control. We randomly choose 10 causal
21   loci and 500 noise loci across the genome, each spanning 200kb. Each causal locus contains a 10kb causal
22   window. For each replicate, we randomly set 10% variants in each 10kb causal window to be causal. In
23   total, there are approximately 335 causal variants on average across the genome. We generated the
24   quantitative/dichotomous trait as follows:

25   $$\text{Quantitative trait: } Y_i = X_{i1} + \sum_{k=1}^{10}(\beta_{k1} g_{k1} + \cdots + \beta_{k,k_s} g_{k,k_s}) + \varepsilon_i,$$

26   $$\text{Dichotomous trait: } g(\mu_i) = \beta_0 + X_{i1} + \sum_{k=1}^{10}(\beta_{k1} g_{k1} + \cdots + \beta_{k,k_s} g_{k,k_s}) + \varepsilon_i,$$

27   where $X_{i1} \sim N(0,1)$, $\varepsilon_i \sim N(0,3)$ and they are all independent; $(g_1, \ldots, g_s)$ are selected risk variants; $g(x) =$
28   $\log(\frac{x}{1-x})$ and $\mu_i$ is the conditional mean of $Y_i$; for dichotomous trait, $\beta_0$ is chosen such that the prevalence
29   is 10%. We set the effect $\beta_{kj} = \frac{a_k}{\sqrt{2m_{kj}(1-m_{kj})}}$, where $m_{kj}$ is the MAF for the $j$-th variant in causal window
30   $k$. We define $a_k$ such that the phenotypic variance due to the risk variants for each causal locus, $\beta_{k1} g_{k1} +$
31   $\cdots + \beta_{k,k_s} g_{k,k_s}$, is 1. We scan the regions with candidate window sizes (1bp, 1kb, 5kb, 10kb), and we
32   consider several tests including the burden test, dispersion test, and Cauchy combination test to aggregate
33   burden, dispersion, and individual variant test results (as discussed in the main text). This combined test is
34   the method implemented in the *KnockoffScreen* method. For each replicate, the empirical power is defined
35   as the proportion of causal loci (the 200kb regions) being identified; the empirical FDR is defined as the
36   proportion of detected windows not overlapping with the causal window +/- 50kb/75kb/100kb, which
37   evaluates FDR at different resolutions. The empirical power and FDR are averaged over 100 replicates.

38   **Simulations for investigating various properties of the *KnockoffScreen* method (the prioritization of
39   causal variants, the influence of shadow effects from common variants, and robustness to population
40   stratification).** We design these simulations to mimic the real data analysis of ADSP. For each replicate,
41   we randomly drew 1,000 variants, including both common and rare variants, from the 200kb region near
42   gene *APOE* (chr19: 44905796-44909393). We scanned the regions with candidate window sizes (1bp, 1kb,
43   5kb, 1kb) using the conventional association test and *KnockoffScreen*. For a fair comparison, we adopted

25

1  the same tests implemented in *KnockoffScreen* to calculate the p-value for the conventional association
2  testing method.

3  *Prioritization of causal variants*. We set 0.25% variants to be causal, all within a 5kb signal window (similar
4  to the size of *APOE*), and then simulated a dichotomous trait by

$$g(\mu_i) = \beta_0 + X_{i1} + \beta_1 g_1 + \cdots + \beta_s g_s,$$

6  where $g(x) = \log(\frac{x}{1-x})$ and $\mu_i$ is the conditional mean of $Y_i$; for dichotomous trait, $\beta_0$ is chosen such that
7  the prevalence is 10%. We set the effect $\beta_j = a|\log_{10} m_j|$, where $m_j$ is the MAF for the $j$-th variant. We
8  defined $a = 1.4$ such that the risk variant has a similar odds ratio as *APOE-ε4* (~3.1) given a similar MAF
9  (~0.137)[59,60]. For each replicate, we compared the two methods in terms of (1) the proportion of selected
10  windows that overlaps with the causal window; and (2) the maximum distance between selected windows
11  and the causal window.

12  *Shadow effect.* We adopted the same simulation setting but set the causal variants to be common
13  (MAF>0.01) and apply the methods to rare variants only (MAF<0.01). Since all causal variants are common,
14  all detected windows are false positives due to the shadow effect. We counted the number of false positives
15  and show the distribution over 500 replicates.

16  *Population stratification.* The ADSP includes three ethnic groups: African American (AA), Non Hispanic
17  White (NHW) and Others (98% of which are Caribbean Hispanic). Let $Z_i$ denote the ethnic group ($Z_i = 0$:
18  AA; $Z_i = 1$: NHW; $Z_i = 2$: Others). We simulated quantitative and dichotomous traits by

19  Quantitative trait: $Y_i = X_{i1} + Z_i + \varepsilon_i$

20  Dichotomous trait: $g(\mu_i) = \beta_0 + X_{i1} + Z_i$

21  where $X_{i1} \sim N(0,1)$, $\varepsilon_i \sim N(0,3)$ and they are all independent; $g(x) = \log(\frac{x}{1-x})$ and $\mu_i$ is the conditional
22  mean of $Y_i$ ; for dichotomous trait, $\beta_0$ is chosen such that the prevalence is 10%. This way, the
23  mean/prevalence for the quantitative/dichotomous trait is a function of the subpopulation, but not directly
24  affected by the genetic variants. We counted the number of false positives and show the distribution over
25  500 replicates. We also calculated an estimate of the FDR, defined as the proportion of replicates where
26  any window is detected.

27  *Population stratification driven by rare variants.* We carried out additional simulation studies to simulate
28  population stratification driven by rare variants using the ADSP data. Specifically, we randomly choose
29  100 regions across the whole genome but outside chromosome 19 with each region of size 200kb. Each
30  region contains a 10kb causal window. We randomly set 10% rare variants (MAF<0.01; MAC>10) in each
31  causal window to exhibit small effects on the trait of interest, Thus the allele frequency differences across
32  ethnic groups will lead to different disease prevalence, reflecting a population stratification driven by rare
33  variants. Then we evaluate the FDR for the selected 200kb region near gene *APOE* (chr19: 44905796-
34  44909393). Since the causal variants are independent of the target region, the confounding effect will be
35  due to population stratification. Specifically, we generated the quantitative/dichotomous trait as follows:

36  Quantitative trait: $Y_i = X_{i1} + \gamma \sum_{k=1}^{100}(\beta_{k1} g_{k1} + \cdots + \beta_{k,k_s} g_{k,k_s}) + \varepsilon_i,$

37  Dichotomous trait: $g(\mu_i) = \beta_0 + X_{i1} + \gamma \sum_{k=1}^{100}(\beta_{k1} g_{k1} + \cdots + \beta_{k,k_s} g_{k,k_s}) + \varepsilon_i,$

38  where $X_{i1} \sim N(0,1)$, $\varepsilon_i \sim N(0,3)$ and they are all independent; $(g_1, \ldots, g_s)$ are selected risk variants; $g(x) =$
39  $\log(\frac{x}{1-x})$ and $\mu_i$ is the conditional mean of $Y_i$; for dichotomous trait, $\beta_0$ is chosen such that the prevalence
40  is 10%. We set the effect $\beta_{kj} = \frac{a_k}{\sqrt{2m_{kj}(1-m_{kj})}}$, where $m_{kj}$ is the MAF for the $j$-th variant in causal window
41  $k$. We define $a_k$ such that the variance due to the risk variants for each causal locus, $\beta_{k1} g_{k1} + \cdots +$
42  $\beta_{k,k_s} g_{k,k_s}$, is 0.01; we set $\gamma = 0, 0.25, 0.5, 0.75$ which quantifies the magnitude of population stratification.

26

1  **The Alzheimer's Disease Sequencing Project.** We first applied *KnockoffScreen* to whole-genome
2  sequencing (WGS) data from the Alzheimer's Disease Sequencing Project (ADSP)[61]. The data include
3  3,085 whole genomes from the ADSP Discovery Extension Study including 1,096 Non-Hispanic White
4  (NHW), 977 African American (AA) descent and 1,012 Caribbean Hispanic (CH). Sequencing for these
5  samples was conducted through three National Human Genome Research Institute (NHGRI) funded Large
6  Scale Sequencing and Analysis Centers (LSACs): Baylor College of Medicine Human Genome Sequencing
7  Center, the Broad Institute, the McDonnell Genome Institute at Washington University. The samples were
8  sequenced on the Illumina HiSeq X Ten platform with 150bp paired-end reads. Additionally, the dataset
9  includes 809 whole genomes from the Alzheimer's Disease Neuroimaging Initiative (ADNI) with 756
10 NHW, 28 AA and 25 others. The samples were sequenced on the Illumina HiSeq 2000 platform with 100bp
11 paired-end reads. Whole-genome sequence data on 809 ADNI subjects (cases, mild cognitive impairment,
12 and controls) have been harmonized using the ADSP pipeline for joint analysis. The ADSP Quality Control
13 Work Group performs QC and concordance checks into an overall ADSP VCF file.

14 **COPDGene from the TOPMed Project.** Eligible subjects in COPDGene Study (NCT00608764,
15 www.copdgene.org) were of non-Hispanic white (NHW) or African-American (AA) ancestry, aged 45-80
16 years old, with at least 10 pack-years of smoking and no diagnosed lung disease other than COPD or
17 asthma[62]. IRB approval was obtained at all study centers, and all study participants provided written
18 informed consent. All subjects underwent a baseline survey, including demographics, smoking history, and
19 symptoms; pre- and post-bronchodilator lung function testing; and chest CT scans. Samples from
20 COPDGene were sequenced at the Broad Institute and at the Northwest Genomics Center at the University
21 of Washington. Variants for all TOPMed samples were jointly called by the Informatics Research Center
22 at the University of Michigan. For details on sequencing and variant calling methods,
23 see https://www.nhlbiwgs.org/topmed-whole-genome-sequencing-project-freeze-5b-phases-1-and-2. QC
24 included comparison of annotated and genetic sex and comparison of genotypes from prior SNP array data
25 with genotypes called from sequencing. Samples with questionable identity from either of these checks
26 were excluded from analysis.

27 **Gene annotation of the identified windows**. The windows (single bp or larger) identified as significant at
28 a target FDR threshold are mapped to genes or intergenic regions using the human genome
29 assembly GRCh38.p13 from the Ensembl Release 99[63]. We assign each significant window to its
30 overlapping locus (gene or intergenic region). If the locus is a gene, we report the gene's name; if the locus
31 is intergenic, we report the upstream and downstream genes (enclosed within parentheses and separated by
32 "-"). We also check if the assigned locus has known associations with Alzheimer's disease and lung related
33 traits in the NHGRI-EBI GWAS Catalog[34]. Specifically, we look up associations with the following seven
34 traits for the ADSP: Alzheimer's disease, late-onset Alzheimer's disease, family history of Alzheimer's
35 disease, t-tau measurement, p-tau measurement, amyloid-beta measurement, and beta-amyloid 1-42
36 measurement; and associations with the following 20 traits for the COPDGene: FEV1/FEC ratio, FEV1,
37 FVC, PEF (peak expiratory flow), COPD, response to bronchodilator, asthma, chronic bronchitis, lung
38 carcinoma, lung adenocarcinoma, pulmonary artery enlargement, FEV change measurement, pulmonary
39 function measurement, carbon monoxide exhalation measurement, airway responsiveness measurement,
40 serum IgE measurement, smoking behaviour measurement, smoking status measurement, smoking
41 behaviour, and smoking initiation. These annotations are shown in Supplemental Tables.

42

**Data Availability**

The manuscript used data from existing studies from COPDGene (TopMED, dbGaP phs000951.v4.p4) and the Alzheimer's Disease Sequencing Project (dbGaP phs000572.v8.p4).

**Code Availability**

We have implemented *KnockoffScreen* in a computationally efficient R package that can be applied generally to the analysis of other whole-genome sequencing studies. The package can be accessed at: https://cran.r-project.org/web/packages/KnockoffScreen/index.html.

**References**

1.  RK, C.Y. *et al.* Whole genome sequencing resource identifies 18 new candidate genes for autism spectrum disorder. *Nat Neurosci* **20**, 602-611 (2017).
2.  Taliun, D. *et al.* Sequencing of 53,831 diverse genomes from the NHLBI TOPMed Program. *bioRxiv*, 563866 (2019).
3.  Morrison, A.C. *et al.* Practical Approaches for Whole-Genome Sequence Analysis of Heart- and Blood-Related Traits. *The American Journal of Human Genetics* **100**, 205-215 (2017).
4.  Sazonovs, A. & Barrett, J.C. Rare-Variant Studies to Complement Genome-Wide Association Studies. *Annu Rev Genomics Hum Genet* **19**, 97-112 (2018).
5.  Hormozdiari, F., Kostem, E., Kang, E.Y., Pasaniuc, B. & Eskin, E. Identifying causal variants at loci with multiple signals of association. *Genetics* **198**, 497-508 (2014).
6.  Wang, G., Sarkar, A., Carbonetto, P. & Stephens, M. A simple new approach to variable selection in regression, with application to genetic fine mapping. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **82**, 1273-1300 (2020).
7.  Korthauer, K. *et al.* A practical guide to methods controlling false discoveries in computational biology. *Genome biology* **20**, 118 (2019).
8.  He, X. *et al.* Integrated model of de novo and inherited genetic variants yields greater power to identify risk genes. *PLoS genetics* **9**(2013).
9.  Iossifov, I. *et al.* The contribution of de novo coding mutations to autism spectrum disorder. *Nature* **515**, 216-221 (2014).
10. Consortium, G. Genetic effects on gene expression across human tissues. *Nature* **550**, 204-213 (2017).
11. Liu, Y. *et al.* A statistical framework for mapping risk genes from de novo mutations in whole-genome-sequencing studies. *The American Journal of Human Genetics* **102**, 1031-1047 (2018).
12. Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal statistical society: series B (Methodological)* **57**, 289-300 (1995).
13. Barber, R.F. & Candès, E.J. Controlling the false discovery rate via knockoffs. *The Annals of Statistics* **43**, 2055-2085 (2015).
14. Candes, E., Fan, Y., Janson, L. & Lv, J. Panning for gold:'model-X'knockoffs for high dimensional controlled variable selection. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **80**, 551-577 (2018).
15. Sesia, M., Katsevich, E., Bates, S., Candès, E. & Sabatti, C. Multi-resolution localization of causal variants across the genome. *Nature Communications* **11**, 1093 (2020).
16. Sesia, M., Sabatti, C. & Candès, E.J. Rejoinder: 'Gene hunting with hidden Markov model knockoffs'. *Biometrika* **106**, 35-45 (2019).
17. Romano, Y., Sesia, M. & Candès, E. Deep knockoffs. *Journal of the American Statistical Association*, 1-12 (2019).

18. He, Z., Xu, B., Buxbaum, J. & Ionita-Laza, I. A genome-wide scan statistic framework for whole-genome sequence data analysis. *Nature communications* **10**, 1-11 (2019).

19. Liu, Y. & Xie, J. Cauchy combination test: a powerful test with analytic p-value calculation under arbitrary dependency structures. *Journal of the American Statistical Association*, 1-18 (2019).

20. Hernandez, R.D. *et al.* Ultra-rare variants drive substantial cis-heritability of human gene expression. *bioRxiv*, 219238 (2019).

21. Zhou, W. *et al.* Efficiently controlling for case-control imbalance and sample relatedness in large-scale genetic association studies. *Nat Genet* **50**, 1335-1341 (2018).

22. Chen, Z. *et al.* Threshold for neural tube defect risk by accumulated singleton loss-of-function variants. *Cell research* **28**, 1039-1041 (2018).

23. He, Z., Xu, B., Lee, S. & Ionita-Laza, I. Unified sequence-based association tests allowing for multiple functional annotations and meta-analysis of noncoding variation in metabochip data. *The American Journal of Human Genetics* **101**, 340-352 (2017).

24. Li, B. & Leal, S.M. Methods for detecting associations with rare variants for common diseases: application to analysis of sequence data. *The American Journal of Human Genetics* **83**, 311-321 (2008).

25. Madsen, B.E. & Browning, S.R. A groupwise association test for rare mutations using a weighted sum statistic. *PLoS genetics* **5**(2009).

26. Wu, M.C. *et al.* Rare-variant association testing for sequencing data with the sequence kernel association test. *The American Journal of Human Genetics* **89**, 82-93 (2011).

27. Rentzsch, P., Witten, D., Cooper, G.M., Shendure, J. & Kircher, M. CADD: predicting the deleteriousness of variants throughout the human genome. *Nucleic acids research* **47**, D886-D894 (2019).

28. He, Z., Liu, L., Wang, K. & Ionita-Laza, I. A semi-supervised approach for predicting cell-type specific functional consequences of non-coding variation using MPRAs. *Nature communications* **9**, 1-12 (2018).

29. Liu, Y. *et al.* Acat: A fast and powerful p value combination method for rare-variant analysis in sequencing studies. *The American Journal of Human Genetics* **104**, 410-421 (2019).

30. Gimenez, J.R. & Zou, J. Improving the Stability of the Knockoff Procedure: Multiple Simultaneous Knockoffs and Entropy Maximization. *arXiv preprint arXiv:1810.11378* (2018).

31. Zhou, X. *et al.* Non-coding variability at the APOE locus contributes to the Alzheimer's risk. *Nature communications* **10**, 1-16 (2019).

32. Lee, S., Abecasis, G.R., Boehnke, M. & Lin, X. Rare-variant association analysis: study designs and statistical tests. *The American Journal of Human Genetics* **95**, 5-23 (2014).

33. Sesia, M., Bates, S., Candès, E., Marchini, J. & Sabatti, C. Controlling the false discovery rate in GWAS with population structure. *bioRxiv* (2020).

34. Buniello, A. *et al.* The NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays and summary statistics 2019. *Nucleic acids research* **47**, D1005-D1012 (2019).

35. Jansen, I.E. *et al.* Genome-wide meta-analysis identifies new loci and functional pathways influencing Alzheimer's disease risk. *Nature genetics* **51**, 404-413 (2019).

36. Marioni, R.E. *et al.* GWAS on family history of Alzheimer's disease. *Translational psychiatry* **8**, 1-7 (2018).

37. Dumitriu, A. *et al.* Integrative analyses of proteomics and RNA transcriptomics implicate mitochondrial processes, protein folding pathways and GWAS loci in Parkinson disease. *BMC medical genomics* **9**, 5 (2015).

38. Lee, J.H. *et al.* Fine mapping of 10q and 18q for familial Alzheimer's disease in Caribbean Hispanics. *Molecular psychiatry* **9**, 1042-1051 (2004).

39. McInnes, L.A. *et al.* A complete genome screen for genes predisposing to severe bipolar disorder in two Costa Rican pedigrees. *Proceedings of the National Academy of Sciences* **93**, 13060-13065 (1996).

40. Ho, A. *et al.* Circulating glucuronic acid predicts healthspan and longevity in humans and mice. *Aging (Albany NY)* **11**, 7694 (2019).

41. Xu, Z., Wu, C., Pan, W. & Initiative, A.s.D.N. Imaging-wide association study: Integrating imaging endophenotypes in GWAS. *Neuroimage* **159**, 159-169 (2017).

42. Shi, J. *et al.* Genome-wide association study of recurrent early-onset major depressive disorder. *Molecular psychiatry* **16**, 193-201 (2011).

43. Mez, J. *et al.* Two novel loci, COBL and SLC10A2, for Alzheimer's disease in African Americans. *Alzheimer's & Dementia* **13**, 119-129 (2017).

44. NHLBI Trans-Omics for Precision Medicine. TOPMed Whole Genome Sequencing Project - Freeze 5b, Phases 1 and 2. Vol. 2020 (https://www.nhlbiwgs.org/topmed-whole-genome-sequencing-project-freeze-5b-phases-1-and-2).

45. Janson, L. & Su, W. Familywise error rate control via knockoffs. *Electronic Journal of Statistics* **10**, 960-975 (2016).

46. Schaid, D.J., Chen, W. & Larson, N.B. From genome-wide associations to candidate causal variants by statistical fine-mapping. *Nature Reviews Genetics* **19**, 491-504 (2018).

47. Zhou, W. *et al.* Efficiently controlling for case-control imbalance and sample relatedness in large-scale genetic association studies. *Nature genetics* **50**, 1335-1341 (2018).

48. Liu, D.J. *et al.* Meta-analysis of gene-level tests for rare variant association. *Nature genetics* **46**, 200 (2014).

49. Feng, S., Liu, D., Zhan, X., Wing, M.K. & Abecasis, G.R. RAREMETAL: fast and powerful meta-analysis for rare variants. *Bioinformatics* **30**, 2828-2829 (2014).

50. Lee, S., Teslovich, T.M., Boehnke, M. & Lin, X. General framework for meta-analysis of rare variants in sequencing association studies. *The American Journal of Human Genetics* **93**, 42-53 (2013).

51. Chen, H. *et al.* Efficient Variant Set Mixed Model Association Tests for Continuous and Binary Traits in Large-Scale Whole-Genome Sequencing Studies. *Am J Hum Genet* **104**, 260-274 (2019).

52. Zhou, W. *et al.* Scalable generalized linear mixed model for region-based association tests in large biobanks and cohorts. (Nature Publishing Group, 2020).

53. Zhao, Z. *et al.* UK Biobank whole-exome sequence binary phenome analysis with robust region-based rare-variant test. *The American Journal of Human Genetics* **106**, 3-12 (2020).

54. Gabriel, S.B. *et al.* The structure of haplotype blocks in the human genome. *Science* **296**, 2225-2229 (2002).

55. Anderson, E.C. & Novembre, J. Finding haplotype block boundaries by using the minimum-description-length principle. *The American Journal of Human Genetics* **73**, 336-354 (2003).

56. Wang, L. GEE analysis of clustered binary data with diverging number of covariates. *The Annals of Statistics* **39**, 389-417 (2011).

57. Storey, J.D. The positive false discovery rate: a Bayesian interpretation and the q-value. *The Annals of Statistics* **31**, 2013-2035 (2003).

58. Storey, J.D. & Tibshirani, R. Statistical significance for genomewide studies. *Proceedings of the National Academy of Sciences* **100**, 9440-9445 (2003).

59. Liu, C.-C., Kanekiyo, T., Xu, H. & Bu, G. Apolipoprotein E and Alzheimer disease: risk, mechanisms and therapy. *Nature Reviews Neurology* **9**, 106 (2013).

60. Kukull, W.A. *et al.* Apolipoprotein E in Alzheimer's disease risk and case detection: a case-control study. *Journal of clinical epidemiology* **49**, 1143-1148 (1996).

61. Beecham, G.W. *et al.* The Alzheimer's disease sequencing project: study design and sample selection. *Neurology Genetics* **3**, e194 (2017).

62. Regan, E.A. *et al.* Genetic epidemiology of COPD (COPDGene) study design. *COPD: Journal of Chronic Obstructive Pulmonary Disease* **7**, 32-43 (2011).

63. Yates, A.D. *et al.* Ensembl 2020. *Nucleic acids research* **48**, D682-D688 (2020).

6

## Author Contributions

8  Z.H., L.L. and I.I.-L. developed the concepts for the manuscript and proposed the method. Z.H., L.L., S.M.,
9  H.T., M.G. and I.I.-L designed the analyses and applications and discussed results. Z.H., C.W., Y.L., J.L.
10 and F.L. conducted the analyses. E.K.S., S.G. and M.H.C. helped interpret the results of the TOPMed
11 analyses. Z.H., L.L. and I.I.-L. prepared the manuscript and all authors contributed to editing the paper.

12

## Competing interests

14 The Authors declare no competing interests.

15

16

17