

## **A statistical framework for powerful multi-trait rare variant analysis in large-scale whole-genome sequencing studies**

Xihao Li<sup>1,2</sup>, Han Chen<sup>3,4</sup>, Margaret Sunitha Selvaraj<sup>5,6,7</sup>, Eric Van Buren<sup>8</sup>, Hufeng Zhou<sup>8</sup>, Yuxuan Wang<sup>9</sup>, Ryan Sun<sup>10</sup>, Zachary R. McCaw<sup>1</sup>, Zhi Yu<sup>5,6,7</sup>, Donna K. Arnett<sup>11</sup>, Joshua C. Bis<sup>12</sup>, John Blangero<sup>13</sup>, Eric Boerwinkle<sup>3,14</sup>, Donald W. Bowden<sup>15</sup>, Jennifer A. Brody<sup>12</sup>, Brian E. Cade<sup>6,16,17</sup>, April P. Carson<sup>18</sup>, Jenna C. Carlson<sup>19</sup>, Nathalie Chami<sup>20</sup>, Yii-Der Ida Chen<sup>21</sup>, Joanne E. Curran<sup>13</sup>, Paul S. de Vries<sup>3</sup>, Myriam Fornage<sup>3,22</sup>, Nora Franceschini<sup>23</sup>, Barry I. Freedman<sup>24</sup>, Charles Gu<sup>25</sup>, Nancy L. Heard-Costa<sup>26,27</sup>, Jiang He<sup>28,29</sup>, Lifang Hou<sup>30</sup>, Yi-Jen Hung<sup>31</sup>, Marguerite R. Irvin<sup>32</sup>, Robert C. Kaplan<sup>33,34</sup>, Sharon L.R. Kardia<sup>35</sup>, Tanika Kelly<sup>36</sup>, Iain Konigsberg<sup>37</sup>, Charles Kooperberg<sup>34</sup>, Brian G. Kral<sup>38</sup>, Changwei Li<sup>28,29</sup>, Ruth J.F. Loos<sup>20,39</sup>, Michael C. Mahaney<sup>13</sup>, Lisa W. Martin<sup>40</sup>, Rasika A. Mathias<sup>38</sup>, Ryan L. Minster<sup>19</sup>, Braxton D. Mitchell<sup>41</sup>, May E. Montasser<sup>41</sup>, Alanna C. Morrison<sup>3</sup>, Nicholette D. Palmer<sup>15</sup>, Patricia A. Peyser<sup>35</sup>, Bruce M. Psaty<sup>12,42,43</sup>, Laura M. Raffield<sup>2</sup>, Susan Redline<sup>16,17</sup>, Alexander P. Reiner<sup>34,42</sup>, Stephen S. Rich<sup>44</sup>, Colleen M. Sitlani<sup>12</sup>, Jennifer A. Smith<sup>35</sup>, Kent D. Taylor<sup>21</sup>, Hemant Tiwari<sup>45</sup>, Ramachandran S. Vasan<sup>27,46</sup>, Zhe Wang<sup>20</sup>, Lisa R. Yanek<sup>38</sup>, Bing Yu<sup>3</sup>, NHLBI Trans-Omics for Precision Medicine (TOPMed) Consortium, Kenneth M. Rice<sup>47</sup>, Jerome I. Rotter<sup>21</sup>, Gina M. Peloso<sup>9</sup>, Pradeep Natarajan<sup>5,6,7</sup>, Zilin Li<sup>8,\*</sup>, Zhonghua Liu<sup>48,\*</sup> and Xihong Lin<sup>6,8,49\*</sup>

<sup>1</sup>Department of Biostatistics, University of North Carolina at Chapel Hill, Chapel Hill, NC, USA.

<sup>2</sup>Department of Genetics, University of North Carolina at Chapel Hill, Chapel Hill, NC, USA.

<sup>3</sup>Human Genetics Center, Department of Epidemiology, Human Genetics, and Environmental Sciences, School of Public Health, The University of Texas Health Science Center at Houston, Houston, TX, USA.

<sup>4</sup>Center for Precision Health, School of Biomedical Informatics, The University of Texas Health Science Center at Houston, Houston, TX, USA.

<sup>5</sup>Center for Genomic Medicine and Cardiovascular Research Center, Massachusetts General Hospital, Boston, MA, USA.

<sup>6</sup>Program in Medical and Population Genetics, Broad Institute of Harvard and MIT, Cambridge, MA, USA.

<sup>7</sup>Department of Medicine, Harvard Medical School, Boston, MA, USA.

<sup>8</sup>Department of Biostatistics, Harvard T.H. Chan School of Public Health, Boston, MA, USA.

<sup>9</sup>Department of Biostatistics, Boston University School of Public Health, Boston, MA, USA.

<sup>10</sup>Department of Biostatistics, University of Texas MD Anderson Cancer Center, Houston, TX, USA.

<sup>11</sup>Provost Office, University of South Carolina, Columbia, SC, USA.

<sup>12</sup>Cardiovascular Health Research Unit, Department of Medicine, University of Washington, Seattle, WA, USA.

<sup>13</sup>Department of Human Genetics and South Texas Diabetes and Obesity Institute, School of Medicine, The University of Texas Rio Grande Valley, Brownsville, TX, USA.

<sup>14</sup>Human Genome Sequencing Center, Baylor College of Medicine, Houston, TX, USA.

<sup>15</sup>Department of Biochemistry, Wake Forest University School of Medicine, Winston-Salem, NC, USA.

<sup>16</sup>Division of Sleep and Circadian Disorders, Brigham and Women's Hospital, Boston, MA, USA.

<sup>17</sup>Division of Sleep Medicine, Harvard Medical School, Boston, MA, USA.

<sup>18</sup>Department of Medicine, University of Mississippi Medical Center, Jackson, MS, USA.

<sup>19</sup>Department of Human Genetics and Department of Biostatistics, University of Pittsburgh, Pittsburgh, PA, USA.

<sup>20</sup>The Charles Bronfman Institute for Personalized Medicine, Icahn School of Medicine at Mount Sinai, New York, NY, USA.

<sup>21</sup>The Institute for Translational Genomics and Population Sciences, Department of Pediatrics, The Lundquist Institute for Biomedical Innovation at Harbor-UCLA Medical Center, Torrance, CA, USA.

<sup>22</sup>Brown Foundation Institute of Molecular Medicine, McGovern Medical School, the University of Texas Health Science Center at Houston, Houston, TX, USA.

<sup>23</sup>Department of Epidemiology, University of North Carolina at Chapel Hill, Chapel Hill, NC, USA.

<sup>24</sup>Department of Internal Medicine, Nephrology, Wake Forest University School of Medicine, Winston-Salem, NC, USA.

<sup>25</sup>Division of Biology & Biomedical Sciences, Washington University School of Medicine, St. Louis, MO, USA.

<sup>26</sup>Department of Neurology, Boston University Chobanian & Avedisian School of Medicine, Boston, MA, USA.

<sup>27</sup>Framingham Heart Study, Framingham, MA, USA.

<sup>28</sup>Department of Epidemiology, Tulane University School of Public Health and Tropical Medicine, New Orleans, LA, USA.

<sup>29</sup>Tulane University Translational Science Institute, New Orleans, LA, USA.

<sup>30</sup>Department of Preventive Medicine, Northwestern University, Chicago, IL, USA.

<sup>31</sup>Department of Internal Medicine, Tri-Service General Hospital, National Defense Medical Center, Taipei, Taiwan.

<sup>32</sup>Department of Epidemiology, School of Public Health, University of Alabama at Birmingham, Birmingham, AL, USA.

<sup>33</sup>Department of Epidemiology and Population Health, Albert Einstein College of Medicine, Bronx, NY, USA.

<sup>34</sup>Division of Public Health Sciences, Fred Hutchinson Cancer Center, Seattle, WA, USA.

<sup>35</sup>Department of Epidemiology, School of Public Health, University of Michigan, Ann Arbor, MI, USA.

<sup>36</sup>Department of Medicine, Division of Nephrology, University of Illinois Chicago, Chicago, IL, USA.

<sup>37</sup>Department of Biomedical Informatics, University of Colorado, Aurora, CO, USA.

<sup>38</sup>Department of Medicine, Johns Hopkins University School of Medicine, Baltimore, MD, USA.

<sup>39</sup>Novo Nordisk Foundation Center for Basic Metabolic Research, Faculty of Health and Medical Sciences, University of Copenhagen, Copenhagen, Denmark.

<sup>40</sup>George Washington University School of Medicine and Health Sciences, Washington, DC, USA.

<sup>41</sup>Department of Medicine, University of Maryland School of Medicine, Baltimore, MD, USA.

<sup>42</sup>Departments of Epidemiology, University of Washington, Seattle, WA, USA.

<sup>43</sup>Department of Health Systems and Population Health, University of Washington, Seattle, WA, USA.

<sup>44</sup>Center for Public Health Genomics, University of Virginia, Charlottesville, VA, USA.

<sup>45</sup>Department of Biostatistics, School of Public Health, University of Alabama at Birmingham, Birmingham, AL, USA.

<sup>46</sup>Department of Quantitative and Qualitative Health Sciences, UT Health San Antonio School of Public Health, San Antonio, TX, USA.

<sup>47</sup>Department of Biostatistics, University of Washington, Seattle, WA, USA.

<sup>48</sup>Department of Biostatistics, Mailman School of Public Health, Columbia University, New York, NY, USA.

<sup>49</sup>Department of Statistics, Harvard University, Cambridge, MA, USA.

\*Correspondence should be addressed to Z. Li ([li@hsph.harvard.edu](mailto:li@hsph.harvard.edu)), Z. Liu ([zl2509@cumc.columbia.edu](mailto:zl2509@cumc.columbia.edu)) and X. Lin ([xlin@hsph.harvard.edu](mailto:xlin@hsph.harvard.edu)).

List of consortium members and their affiliations appears at the end of the paper.

1 **Abstract**

2 Large-scale whole-genome sequencing (WGS) studies have improved our  
3 understanding of the contributions of coding and noncoding rare variants to complex  
4 human traits. Leveraging association effect sizes across multiple traits in WGS rare  
5 variant association analysis can improve statistical power over single-trait analysis, and  
6 also detect pleiotropic genes and regions. Existing multi-trait methods have limited  
7 ability to perform rare variant analysis of large-scale WGS data. We propose  
8 MultiSTAAR, a statistical framework and computationally-scalable analytical pipeline for  
9 functionally-informed multi-trait rare variant analysis in large-scale WGS studies.  
10 MultiSTAAR accounts for relatedness, population structure and correlation among  
11 phenotypes by jointly analyzing multiple traits, and further empowers rare variant  
12 association analysis by incorporating multiple functional annotations. We applied  
13 MultiSTAAR to jointly analyze three lipid traits (low-density lipoprotein cholesterol, high-  
14 density lipoprotein cholesterol and triglycerides) in 61,861 multi-ethnic samples from the  
15 Trans-Omics for Precision Medicine (TOPMed) Program. We discovered new  
16 associations with lipid traits missed by single-trait analysis, including rare variants within  
17 an enhancer of *NIPSNAP3A* and an intergenic region on chromosome 1.

18

19

20

21

22

23

24 Advances in next generation sequencing technologies and the decreasing cost of  
25 whole-exome/whole-genome sequencing (WES/WGS) have made it possible to study  
26 the genetic underpinnings of rare variants (i.e. minor allele frequency (MAF) < 1%) in  
27 complex human traits. Large nationwide consortia and biobanks, such as the National  
28 Heart, Lung and Blood Institute (NHLBI)'s Trans-Omics for Precision Medicine  
29 (TOPMed) Program<sup>1</sup>, the National Human Genome Research Institute's Genome  
30 Sequencing Program (GSP) , the National Institute of Health's All of Us Research  
31 Program<sup>2</sup>, and the UK's Biobank WGS Program<sup>3</sup>, are expected to sequence more than  
32 a million of individuals in total, at more than 1 billion genetic variants in both coding and  
33 noncoding regions of the human genome, while also recording thousands of  
34 phenotypes. To mitigate the lack of power of single-variant analyses to identify rare  
35 variant associations<sup>4</sup>, variant set tests have been proposed to analyze the joint effects  
36 of multiple rare variants<sup>5-9</sup>, where most of the work has focused single trait analysis.  
37  
38 Pleiotropy occurs when genetic variants influence multiple traits<sup>10</sup>. There is growing  
39 empirical evidence from genome-wide association studies (GWASs) that many variants  
40 have pleiotropic effects<sup>11,12</sup>. Identifying these effects can provide valuable insights into  
41 the genetic architecture of complex traits<sup>13</sup>. As such, it is of increasing interest to identify  
42 pleiotropic rare variants by jointly analyzing multiple traits in WGS rare variant  
43 association studies (RVASs).  
44  
45 Several existing methods for multi-trait rare variant association analysis, such as  
46 MSKAT<sup>14</sup>, Multi-SKAT<sup>15</sup> and MTAR<sup>16</sup>, have shown that leveraging the cross-phenotype



47 correlation structure can improve the power of multi-trait analyses compared to single-  
48 trait analyses when analyzing pleiotropic genes<sup>14-17</sup>. However, existing methods do not  
49 scale well, and are not feasible when analyzing large-scale WGS studies with hundreds  
50 of millions of rare variants in samples exhibiting relatedness and population structure.  
51 Furthermore, none of the existing multi-trait rare variant analysis methods leverages  
52 functional annotations that predict the biological functionality of variants, resulting in  
53 limited interpretability and power loss. While the STAAR method<sup>18</sup> dynamically  
54 incorporates multiple variant functional annotations to maximize the power of rare  
55 variant association tests, it is designed for single-trait analysis and cannot be directly  
56 applied to multiple traits.

57

58 To overcome these limitations, we propose the Multi-trait variant-Set Test for  
59 Association using Annotation infoRmation (MultiSTAAR), a statistical framework for  
60 multi-trait rare variant analyses of large-scale WGS studies and biobanks. It has several  
61 features. First, by fitting a null Multivariate Linear Mixed Model (MLMM)<sup>19</sup> for multiple  
62 quantitative traits simultaneously, adjusting for ancestry principal components (PCs)<sup>20</sup>  
63 and using a sparse genetic relatedness matrix (GRM)<sup>21,22</sup>, MultiSTAAR scales well but  
64 also accounts for relatedness and population structure, as well as correlations among  
65 the multiple traits. Second, MultiSTAAR enables the incorporation of multiple variant  
66 functional annotations as weights to improve the power of RVASs. Furthermore, we  
67 provide MultiSTAAR via a comprehensive pipeline for large-scale WGS studies, that  
68 facilitates functionally-informed multi-trait analysis of both coding and noncoding rare

69 variants. Third, MultiSTAAR enables conditional multi-trait analysis to assess rare  
70 variant association signals beyond known common and low frequency variants.

71

72 In the current study, we conducted extensive simulation studies to demonstrate the  
73 validity of MultiSTAAR and to assess the power gain of MultiSTAAR by incorporating  
74 multiple relevant variant functional annotations, and its ability in preserving Type I error  
75 rates. We then applied MultiSTAAR to perform WGS RVAS of 61,838 ancestrally  
76 diverse participants from 20 studies from NHLBI's TOPMed consortium by jointly  
77 analyzing three circulating lipid traits: low-density lipoprotein cholesterol (LDL-C), high-  
78 density lipoprotein cholesterol (HDL-C) and triglycerides (TG). We show that  
79 MultiSTAAR is computationally feasible for large-scale WGS multi-trait rare variant  
80 analysis, and in conditional analysis of LDL-C, HDL-C and TG, MultiSTAAR identifies  
81 signals that were missed either by the existing multi-trait rare variant analysis methods  
82 that overlook variant functional annotations, or by single-trait functionally-informed  
83 analysis that ignore correlations between phenotypes.

84

## 85 **Results**

### 86 **Overview of the methods**

87 MultiSTAAR is a statistical framework and an analytic pipeline for jointly analyzing  
88 multiple traits in large-scale WGS rare variant association studies. There are two main  
89 components in the MultiSTAAR framework: (i) fitting null MLMs using ancestry PCs  
90 and sparse GRMs to account for population structure, relatedness and the correlation  
91 between phenotypes, and (ii) testing for associations between each aggregated variant

92 set and multiple traits by dynamically incorporating multiple variant functional  
93 annotations<sup>18</sup> (**Fig. 1a**).

94

95 In WGS RVASs, an important but often underemphasized challenge is selecting  
96 biologically-meaningful and functionally-interpretable analysis units, especially for the  
97 noncoding genome<sup>23,24</sup>. In gene-centric analyses of multiple traits, MultiSTAAR provides  
98 five functional categories (masks) to aggregate coding rare variants of each protein-  
99 coding gene, as well as an additional eight masks of regulatory regions to aggregate  
100 noncoding rare variants. In non-gene-centric analyses of multiple traits, MultiSTAAR  
101 performs agnostic genetic region analyses using sliding windows<sup>18,25</sup> (**Fig. 1b**).

102

103 For each rare variant set analyzed, MultiSTAAR first constructs the multi-trait burden,  
104 SKAT and ACAT-V test statistics (**Methods**). For each type of rare variant test,  
105 MultiSTAAR calculates multiple candidate *P* values using different variant functional  
106 annotations as weights, following the STAAR framework<sup>18</sup>. MultiSTAAR then  
107 aggregates the association strength by combining the *P* values from all annotations  
108 using the ACAT method, that provides robustness to correlation between tests<sup>9</sup>, and  
109 proposes an omnibus test, MultiSTAAR-O, that leverages the advantages of different  
110 type of tests (**Methods**). Furthermore, MultiSTAAR can test multi-trait rare variants'  
111 associations conditional on a set of known associations (**Fig. 1b**).

112

113 **Simulation studies**

114 To evaluate the type I error rates and the power of MultiSTAAR, we performed  
115 simulation studies under several configurations. Following the steps described in Data  
116 Simulation (**Methods**), we generated three quantitative traits with a correlation matrix  
117 similar to the empirical correlation in the three lipid traits<sup>26-28</sup>. We then generated  
118 genotypes by simulating 20,000 sequences for 100 different 1 megabase (Mb) regions,  
119 each of them were generated to mimic the linkage disequilibrium structure of an African  
120 American population by using the calibration coalescent model<sup>29</sup>. Throughout the  
121 simulation studies, we randomly and uniformly selected 5-kilobase (kb) regions from  
122 these 1-Mb regions and considered sample sizes of 10,000 for each replicate. The  
123 simulation studies focused on aggregating uncommon variants with an MAF < 5%.

124

### 125 **Type I error rate evaluations**

126 We performed  $10^8$  simulations to evaluate the type I error rates of the multi-trait burden,  
127 SKAT, ACAT-V and MultiSTAAR-O tests at  $\alpha = 10^{-4}$ ,  $10^{-5}$ , and  $10^{-6}$  (**Supplementary**  
128 **Table 1**). The results show that, for multi-trait rare variant analysis, all four MultiSTAAR  
129 tests controlled the type I error rates at very close to the nominal  $\alpha$  levels.

130

### 131 **Empirical power simulations**

132 We next assessed the power of MultiSTAAR-O for the analysis of multiple phenotypes  
133 under different genetic architectures, while also comparing its power with existing  
134 methods. Specifically, we considered four models, in which variants in the signal region  
135 (variant-phenotype association regions) were associated with (1) one phenotype only,  
136 (2) two positively correlated phenotypes, (3) two negatively correlated phenotypes and

137 (4) all three phenotypes. In addition, we considered different proportions (5%, 15% and  
138 35% on average) of causal variants in the signal region, where causality of variants  
139 depended on different sets of annotations, and the effect size directions of causal  
140 variants were allowed to vary (**Methods**). Power was evaluated as the proportions of  $P$   
141 values less than  $\alpha = 10^{-7}$  based on  $10^4$  simulations. Overall, MultiSTAAR-O  
142 consistently delivered higher power to detect signal regions compared to multi-trait  
143 burden, SKAT and ACAT-V tests, through its incorporation of multiple annotations  
144 (**Extended Data Figs. 2-5, Supplementary Figs. 1-4**). This power advantage was also  
145 robust to the existence of noninformative annotations.

146

#### 147 **Application to the TOPMed lipids WGS data**

148 We applied MultiSTAAR to identify rare variant associations with three quantitative lipid  
149 traits (LDL-C, HDL-C and TG) through a multi-trait analysis using TOPMed Freeze 8  
150 WGS data, comprising 61,838 individuals from 20 multi-ethnic studies (**Supplementary**  
151 **Note**). LDL-C values were adjusted for the usage of lipid-lowering medication<sup>26,30</sup>  
152 (**Methods**), and DNA samples were sequenced at >30x target coverage. Sample- and  
153 variant-level quality control were performed for each participating study<sup>1,26,30</sup>.

154

155 Race/ethnicity was measured using a combination of self-reported race/ethnicity and  
156 study recruitment information<sup>31</sup> (**Supplementary Note**). Of the 61,838 samples, 15,636  
157 (25.3%) were Black or African American, 27,439 (44.4%) were White, 4,461 (7.2%)  
158 were Asian or Asian American, 13,138 (21.2%) were Hispanic/Latino American and  
159 1,164 (1.9%) were Samoans. There were 414 million single-nucleotide variants (SNVs)

160 observed overall, with 6.5 million (1.6%) common variants (MAF > 5%), 5.2 million  
161 (1.2%) low-frequency variants ( $1\% \leq \text{MAF} \leq 5\%$ ) and 402 million (97.2%) rare variants  
162 (MAF < 1%). The study-specific demographics and baseline characteristics are given in  
163 **Supplementary Table 2.**

164

### 165 **Gene-centric multi-trait analysis of coding and noncoding rare variants**

166 We applied MultiSTAAR-O on gene-centric multi-trait analysis of coding and noncoding  
167 rare variants of genes with lipid traits in TOPMed. For coding variants, rare variants  
168 (MAF < 1%) from five coding functional categories (masks) were aggregated,  
169 separately, and analyzed using a joint model for LDL-C, HDL-C and TG, including (1)  
170 putative loss-of-function (stop gain, stop loss and splice) rare variants, (2) missense  
171 rare variants, (3) disruptive missense rare variants, (4) putative loss-of-function and  
172 disruptive missense rare variants and (5) synonymous rare variants of each protein-  
173 coding gene. The putative loss-of-function, missense and synonymous RVs were  
174 defined by GENCODE Variant Effect Predictor (VEP) categories<sup>32</sup>. The disruptive  
175 variants were further defined by MetaSVM<sup>33</sup>, which measures the deleteriousness of  
176 missense mutations. We incorporated 9 annotation principal components (aPCs)<sup>18,26,34</sup>,  
177 CADD<sup>35</sup>, LINSIGHT<sup>36</sup>, FATHMM-XF<sup>37</sup> and MetaSVM<sup>33</sup> (for missense rare variants only)  
178 along with the two MAF-based weights<sup>4</sup> in MultiSTAAR-O (**Supplementary Table 3**).  
179 The overall distribution of MultiSTAAR-O *P* values was well-calibrated for the multi-trait  
180 analysis of coding rare variants (**Extended Data Fig. 1b**). At a Bonferroni-corrected  
181 significance threshold of  $\alpha = 0.05 / (20,000 \times 5) = 5.00 \times 10^{-7}$ , accounting for five  
182 different coding masks across protein-coding genes, MultiSTAAR-O identified 51

183 genome-wide significant associations using unconditional multi-trait analysis (**Extended**  
184 **Data Fig. 1a, Supplementary Table 4**). After conditioning on previously reported  
185 variants associated with LDL-C, HDL-C or TG located within a 1 Mb broader region of  
186 each coding mask in the GWAS Catalog and Million Veteran Program (MVP)<sup>26,38,39</sup>, 34  
187 out of the 51 associations remained significant at the Bonferroni-corrected threshold of  
188  $\alpha = 0.05/51 = 9.80 \times 10^{-4}$  (**Table 1**).

189  
190 For non-coding variants, rare variants from eight noncoding masks were analyzed in a  
191 similar fashion, including (1) promoter rare variants overlaid with CAGE sites<sup>40</sup>, (2)  
192 promoter rare variants overlaid with DHS sites<sup>41</sup>, (3) enhancer rare variants overlaid  
193 with CAGE sites<sup>42,43</sup>, (4) enhancer rare variants overlaid with DHS sites<sup>41,43</sup>, (5)  
194 untranslated region (UTR) rare variants, (6) upstream region rare variants, (7)  
195 downstream region rare variants of each protein-coding gene and (8) rare variants in  
196 ncRNA genes<sup>24</sup>. The promoter rare variants were defined as rare variants in the  $\pm 3$ -  
197 kilobase (kb) window of transcription start sites with the overlap of CAGE sites or DHS  
198 sites. The enhancer rare variants were defined as RVs in GeneHancer-predicted  
199 regions with the overlap of CAGE sites or DHS sites. The UTR, upstream, downstream  
200 and ncRNA rare variants were defined by GENCODE VEP categories<sup>32</sup>. With a well-  
201 calibrated overall distribution of MultiSTAAR-O *P* values (**Extended Data Fig. 1d**) and  
202 at a Bonferroni-corrected significance threshold of  $\alpha = 0.05/(20,000 \times 7) = 3.57 \times$   
203  $10^{-7}$ , accounting for seven different noncoding masks across protein-coding genes,  
204 MultiSTAAR-O identified 76 genome-wide significant associations using unconditional  
205 multi-trait analysis (**Extended Data Fig. 1c, Supplementary Table 5**). After

206 conditioning on known lipids-associated variants<sup>26,38,39</sup>, 6 out of the 76 associations  
207 remained significant at the Bonferroni-corrected threshold of  $\alpha = 0.05/76 = 6.58 \times$   
208  $10^{-4}$  (**Table 2**). These included promoter CAGE and enhancer CAGE rare variants in  
209 *APOA1*, promoter DHS rare variants in *CETP*, enhancer CAGE rare variants in *SPC24*,  
210 and enhancer DHS rare variants in *NIPSNAP3A* and *LIPC*.

211  
212 MultiSTAAR-O further identified 6 genome-wide significant associations using  
213 unconditional multi-trait analysis at  $\alpha = 0.05/20,000 = 2.50 \times 10^{-6}$  accounting for  
214 ncRNA genes (**Extended Data Fig. 1e, Supplementary Table 5**), with 3 rare variant  
215 associations in *RP11-15F12.3*, *RP11-310H4.2* and *MIR4497* remained significant at  
216  $\alpha = 0.05/6 = 8.33 \times 10^{-3}$  after conditioning on known lipids-associated variants<sup>26,38,39</sup>  
217 (**Table 2**).

218  
219 Notably, among the 9 conditionally significant noncoding rare variants associations with  
220 lipid traits, 4 of them were not detected by any of the three single-trait analysis (LDL-C,  
221 HDL-C or TG) using unconditional analysis of STAAR-O, including the associations of  
222 enhancer DHS rare variants in *NIPSNAP3A* and *LIPC* as well as ncRNA rare variants in  
223 *RP11-310H4.2* and *MIR4497* (**Supplementary Table 5**). These results demonstrate  
224 that MultiSTAAR-O can increase power over existing methods, and identify additional  
225 trait-associated signals by leveraging cross-phenotype correlations between multiple  
226 traits.

227

228 **Genetic region multi-trait analysis of rare variants**



229 We next applied MultiSTAAR-O to perform genetic region multi-trait analysis to identify  
230 rare variants associated with lipid traits in TOPMed. Rare variants residing in 2-kilobase  
231 (kb) sliding windows with a 1-kb skip length were aggregated and analyzed using a joint  
232 model for LDL-C, HDL-C and TG. We incorporated 12 quantitative annotations,  
233 including 9 aPCs, CADD, LINSIGHT, FATHMM-XF along with the two MAF weights in  
234 MultiSTAAR-O (**Methods**). The overall distribution of MultiSTAAR-O  $P$  values was well-  
235 calibrated for the multi-trait analysis (**Fig. 2b**). At a Bonferroni-corrected significance  
236 threshold of  $\alpha = 0.05 / (2.65 \times 10^6) = 1.89 \times 10^{-8}$  accounting for 2.65 million 2-kb  
237 sliding windows across the genome, MultiSTAAR-O identified 502 genome-wide  
238 significant associations using unconditional multi-trait analysis (**Fig. 2a, Supplementary**  
239 **Table 6**). By dynamically incorporating multiple functional annotations capturing  
240 different aspects of variant function, MultiSTAAR-O detected more significant sliding  
241 windows and showed consistently smaller  $P$  values for top sliding windows compared  
242 with multi-trait analysis using only MAFs as the weight (**Fig. 2c**). After conditioning on  
243 known lipids-associated variants<sup>26,38,39</sup>, 7 out of the 502 associations remained  
244 significant at the Bonferroni-corrected threshold of  $\alpha = 0.05 / 502 = 9.96 \times 10^{-5}$  (**Table**  
245 **3**), including two sliding windows in *DOCK7* (chromosome 1: 62,651,447 - 62,653,446  
246 bp; chromosome 1: 62,652,447 - 62,654,446 bp) and an intergenic sliding window  
247 (chromosome 1: 145,530,447 - 145,532,446 bp) that were not detected by any of the  
248 three single-trait analysis (LDL-C, HDL-C or TG) using STAAR-O (**Supplementary**  
249 **Table 6**). Notably, all known lipids-associated variants indexed in the previous literature  
250 were at least 1-Mb away from the intergenic sliding window.

251

## 252 **Comparison of MultiSTAAR-O with existing multi-trait rare variant tests**

253 Using TOPMed Freeze 8 WGS data, our gene-centric multi-trait analysis of coding rare  
254 variants identified 34 conditionally significant associations with lipid traits (**Table 1**),  
255 including *NPC1L1* and *SCARB1* missense rare variants that were missed by multi-trait  
256 burden, SKAT and ACAT-V tests (**Supplementary Table 4**). Among the 9 and 7  
257 conditionally significant associations detected in gene-centric multi-trait analysis of  
258 noncoding rare variants and genetic region multi-trait analysis, MultiSTAAR-O identified  
259 1 and 2 associations, respectively, that were missed by multi-trait burden, SKAT and  
260 ACAT-V tests (**Supplementary Tables 5-6**). These associations included enhancer  
261 CAGE rare variants in *SPC24* and two sliding windows in *LDLR* (chromosome 19:  
262 11,104,367 - 11,106,366 bp; chromosome 19: 11,105,367 - 11,107,366 bp).

263

## 264 **Computation cost**

265 The computational cost for MultiSTAAR-O to perform WGS multi-trait rare variant  
266 analysis of  $n = 61,838$  related TOPMed lipids samples was 2 hours using 250 2.10-GHz  
267 computing cores with 12-GB memory for gene-centric coding analysis; or 20 hours  
268 using 250 2.10-GHz computing cores with 24-GB memory for gene-centric noncoding  
269 analysis; 2 hours using 250 2.10-GHz computing cores with 12-GB memory of ncRNA  
270 analysis; and 20 hours using 500 2.10-GHz computing cores with 24-GB memory for  
271 sliding window analysis. Runtime for all analyses scales linearly with the sample size<sup>24</sup>.

272

## 273 **Discussion**

274 In this study, we have introduced MultiSTAAR as a general statistical framework and a  
275 flexible analytical pipeline for performing functionally-informed multi-trait RVAS in large-  
276 scale WGS studies. MultiSTAAR improves power by analyzing multiple traits  
277 simultaneously and dynamically incorporating multiple functional annotations, while  
278 accounting for relatedness and population structure among study samples.

279  
280 By jointly analyzing multiple quantitative traits using a multivariate linear mixed model,  
281 MultiSTAAR explicitly leverages the correlation among multiple phenotypes to enhance  
282 power for detecting additional association signals, outperforming single-trait analyses of  
283 the individual phenotypes. MultiSTAAR also enables conditional multi-trait analysis to  
284 identify putatively novel rare variant associations independent of a set of known  
285 variants. Using TOPMed Freeze 8 WGS data, our gene-centric multi-trait analysis of  
286 noncoding rare variants identified 9 conditionally significant associations with lipid traits  
287 (**Table 2**), including 4 noncoding associations that were missed by single-trait analysis  
288 using STAAR (**Supplementary Table 5**). Our genetic region multi-trait analysis of rare  
289 variants identified 7 conditionally significant 2-kb sliding windows associated with lipid  
290 traits (**Table 3**), including 3 associations that were missed by single-trait analysis using  
291 STAAR (**Supplementary Table 6**).

292  
293 By dynamically incorporating multiple annotations capturing diverse aspects of variant  
294 biological function in the second step, MultiSTAAR further improves power over existing  
295 multi-trait rare variant analysis methods. Our simulation studies demonstrated that  
296 MultiSTAAR-O maintained accurate type I error rates while achieving considerable

297 power gains over multi-trait burden, SKAT and ACAT-V tests that do not incorporate  
298 functional annotation information (**Extended Data Figs. 2-5, Supplementary Figs. 1-**  
299 **4**). Notably, the existing ACAT-V method<sup>9</sup> does not support multi-trait analysis. We  
300 extended it to accommodate multi-trait settings and incorporated the multi-trait ACAT-V  
301 test into the MultiSTAAR framework (**Methods**).

302  
303 Implemented as a flexible analytical pipeline, MultiSTAAR allows for customized input  
304 phenotype selection, variant set definition and user-specified annotation weights to  
305 facilitate functionally-informed multi-trait analyses. In addition to rare variant association  
306 analysis of coding and noncoding regions, MultiSTAAR also provides single-variant  
307 multi-trait analysis for common and low-frequency variants under a given MAF or minor  
308 allele count (MAC) cutoff (e.g.  $MAC \geq 20$ ). Using 61,838 TOPMed lipids samples, it took  
309 8 hours using 250 2.10-GHz computing cores with 12-GB memory for single-variant  
310 multi-trait analysis, which is scalable for large WGS/WES datasets. On the other hand,  
311 MultiSTAAR could be further extended to allow for dynamic windows with data-adaptive  
312 sizes in genetic region analysis<sup>24,44</sup>, to properly leverage synthetic surrogates in the  
313 presence of partially missing phenotypes<sup>45</sup>, and to incorporate summary statistics for  
314 meta-analysis of multiple WGS/WES studies<sup>46</sup>.

315  
316 In summary, MultiSTAAR provides a powerful statistical framework and a  
317 computationally scalable analytical pipeline for large-scale WGS multi-trait analysis with  
318 complex study samples. Compared to single-trait analysis, MultiSTAAR offers a notable  
319 increase in statistical power when analyzing multiple moderately to highly correlated

320 traits, all while maintaining control over type I error rates across various genetic  
321 architectures. As the sample sizes and number of available phenotypes increase in  
322 biobank-scale sequencing studies, our proposed method may contribute to a better  
323 understanding of the genetic architecture of complex traits by elucidating the role of rare  
324 variants with pleiotropic effects.

325

### 326 **Acknowledgments**

327 This work was supported by grants R35-CA197449, U19-CA203654, U01-HG012064,  
328 and U01-HG009088 (X. Lin), NHLBI TOPMed Fellowship (X. Li), R01-HL142711 and  
329 R01-HL127564 (P.N. and G.M.P.), 75N92020D00001, HHSN268201500003I, N01-HC-  
330 95159, 75N92020D00005, N01-HC-95160, 75N92020D00002, N01-HC-95161,  
331 75N92020D00003, N01-HC-95162, 75N92020D00006, N01-HC-95163,  
332 75N92020D00004, N01-HC-95164, 75N92020D00007, N01-HC-95165, N01-HC-95166,  
333 N01-HC-95167, N01-HC-95168, N01-HC-95169, UL1-TR-000040, UL1-TR-001079,  
334 UL1-TR-001420, UL1-TR001881, DK063491, R01-HL071051, R01-HL071205, R01-  
335 HL071250, R01-HL071251, R01-HL071258, R01-HL071259, and UL1-RR033176  
336 (J.I.R.), HHSN268201800001I and U01-HL137162 (K.M.R.), 1R35-HL135818, R01-  
337 HL113338, and HL046389 (S.R.), HL105756 (B.M.P.), HHSN268201600018C,  
338 HHSN268201600001C, HHSN268201600002C, HHSN268201600003C, and  
339 HHSN268201600004C (C.K.), R01-MD012765 and R01-DK117445 (N.F.),  
340 18CDA34110116 from American Heart Association (P.S.d.V.), R01-HL153805, R03-  
341 HL154284 (B.E.C.), HHSN268201700001I, HHSN268201700002I,  
342 HHSN268201700003I, HHSN268201700005I, and HHSN268201700004I (E.B.), U01-

343 HL072524, R01-HL104135-04S1, U01-HL054472, U01-HL054473, U01-HL054495,  
344 U01-HL054509, and R01-HL055673-18S1 (D.K.A.), U01-HL72518, HL087698,  
345 HL49762, HL59684, HL58625, HL071025, HL112064, NR0224103, and M01-  
346 RR000052 (to the Johns Hopkins General Clinical Research Center). This work was  
347 supported by R01 HL92301, R01 HL67348, R01 NS058700, R01 AR48797, R01  
348 DK071891, R01 AG058921, the General Clinical Research Center of the Wake Forest  
349 University School of Medicine (M01 RR07122, F32 HL085989), the American Diabetes  
350 Association, and a pilot grant from the Claude Pepper Older Americans Independence  
351 Center of Wake Forest University Health Sciences (P60 AG10484). The Framingham  
352 Heart Study (FHS) acknowledges the support of contracts NO1-HC-25195,  
353 HHSN268201500001I and 75N92019D00031 from the National Heart, Lung and Blood  
354 Institute and grant supplement R01 HL092577-06S1 for this research. We also  
355 acknowledge the dedication of the FHS study participants without whom this research  
356 would not be possible. R.S.V. is supported in part by the Evans Medical Foundation and  
357 the Jay and Louis Coffman Endowment from the Department of Medicine, Boston  
358 University School of Medicine. The Jackson Heart Study (JHS) is supported and  
359 conducted in collaboration with Jackson State University (HHSN268201800013I),  
360 Tougaloo College (HHSN268201800014I), the Mississippi State Department of Health  
361 (HHSN268201800015I) and the University of Mississippi Medical Center  
362 (HHSN268201800010I, HHSN268201800011I and HHSN268201800012I) contracts  
363 from the National Heart, Lung, and Blood Institute (NHLBI) and the National Institute on  
364 Minority Health and Health Disparities (NIMHD). The authors also wish to thank the  
365 staffs and participants of the JHS. Support for GENOA was provided by the National

366 Heart, Lung and Blood Institute (U01HL054457, U01HL054464, U01HL054481,  
367 R01HL119443, and R01HL087660) of the National Institutes of Health. Collection of the  
368 San Antonio Family Study data was supported in part by National Institutes of Health  
369 (NIH) grants P01 HL045522, MH078143, MH078111 and MH083824; and whole  
370 genome sequencing of SAFS subjects was supported by U01 DK085524 and R01  
371 HL113323. Molecular data for the Trans-Omics in Precision Medicine (TOPMed)  
372 program was supported by the National Heart, Lung and Blood Institute (NHLBI). Core  
373 support including centralized genomic read mapping and genotype calling, along with  
374 variant quality metrics and filtering were provided by the TOPMed Informatics Research  
375 Center (3R01HL-117626-02S1; contract HHSN268201800002I). Core support including  
376 phenotype harmonization, data management, sample-identity QC, and general program  
377 coordination were provided by the TOPMed Data Coordinating Center (R01HL-120393;  
378 U01HL-120393; contract HHSN268201800001I). We gratefully acknowledge the studies  
379 and participants who provided biological samples and data for TOPMed. The full study  
380 specific acknowledgements are detailed in **Supplementary Note**.

381

### 382 **Author contributions**

383 X. Li, H.C., Z. Li, Z. Liu and X. Lin designed the experiments. X. Li, H.C., Z. Li and X.  
384 Lin performed the experiments. X. Li, H.C., M.S.S., E.V.B., Y.W., R.S., Z.R.M., Z.Y.,  
385 D.K.A., J.C.B., J.B., E.B., D.W.B., J.A.B., B.E.C., A.P.C., J.C.C., N.C., Y.D.I.C., J.E.C.,  
386 P.S.d.V., M.F., N.F., B.I.F., C.G., N.L.H.C., J.H., L.H., Y.J.H., M.R.I., R.C.K., S.L.R.K.,  
387 T.K., I.K., C.K., B.G.K., C.L., R.J.F.L., M.C.M., L.W.M., R.A.M., R.L.M., B.D.M., M.E.M.,  
388 A.C.M., N.D.P., P.A.P., B.M.P., L.M.R., S.R., A.P.R., S.S.R., C.M.S., J.A.S., K.D.T.,

389 H.T., R.S.V., Z.W., L.R.Y., B.Y., K.M.R., J.I.R., G.M.P., P.N., Z. Li, Z. Liu and X. Lin  
390 acquired, analyzed or interpreted data. G.M.P., P.N. and the NHLBI TOPMed Lipids  
391 Working Group provided administrative, technical or material support. X. Li, Z. Li, Z. Liu  
392 and X. Lin drafted the manuscript and revised it according to suggestions by the  
393 coauthors. All authors critically reviewed the manuscript, suggested revisions as needed  
394 and approved the final version.

395

### 396 **Competing interests**

397 Z.R.M. is an employee of Insitro. M.E.M. receives research funding from Regeneron  
398 Pharmaceutical Inc., unrelated to this project. B.M.P. serves on the Steering Committee  
399 of the Yale Open Data Access Project funded by Johnson & Johnson. L.M.R. is a  
400 consultant for the TOPMed Administrative Coordinating Center (via Westat). X. Lin is a  
401 consultant of AbbVie Pharmaceuticals and Verily Life Sciences. The remaining authors  
402 declare no competing interests.

403

### 404 **NHLBI Trans-Omics for Precision Medicine (TOPMed) Consortium**

405 Namiko Abe<sup>50</sup>, Gonçalo Abecasis<sup>51</sup>, Francois Aguet<sup>52</sup>, Christine Albert<sup>53</sup>, Laura  
406 Almasy<sup>54</sup>, Alvaro Alonso<sup>55</sup>, Seth Ament<sup>56</sup>, Peter Anderson<sup>57</sup>, Pramod Anugu<sup>58</sup>, Deborah  
407 Applebaum-Bowden<sup>59</sup>, Kristin Ardlie<sup>52</sup>, Dan Arking<sup>60</sup>, Allison Ashley-Koch<sup>61</sup>, Stella  
408 Aslibekyan<sup>62</sup>, Tim Assimes<sup>63</sup>, Paul Auer<sup>64</sup>, Dimitrios Avramopoulos<sup>60</sup>, Najib Ayas<sup>65</sup>,  
409 Adithya Balasubramanian<sup>66</sup>, John Barnard<sup>67</sup>, Kathleen Barnes<sup>68</sup>, R. Graham Barr<sup>69</sup>,  
410 Emily Barron-Casella<sup>60</sup>, Lucas Barwick<sup>70</sup>, Terri Beaty<sup>60</sup>, Gerald Beck<sup>71</sup>, Diane Becker<sup>72</sup>,  
411 Lewis Becker<sup>60</sup>, Rebecca Beer<sup>73</sup>, Amber Beitelshes<sup>56</sup>, Emelia Benjamin<sup>74</sup>, Takis



412 Benos<sup>75</sup>, Marcos Bezerra<sup>76</sup>, Larry Bielak<sup>51</sup>, Thomas Blackwell<sup>51</sup>, Nathan Blue<sup>77</sup>, Russell  
413 Bowler<sup>78</sup>, Ulrich Broeckel<sup>79</sup>, Jai Broome<sup>57</sup>, Deborah Brown<sup>80</sup>, Karen Bunting<sup>50</sup>, Esteban  
414 Burchard<sup>81</sup>, Carlos Bustamante<sup>82</sup>, Erin Buth<sup>83</sup>, Jonathan Cardwell<sup>84</sup>, Vincent Carey<sup>85</sup>,  
415 Julie Carrier<sup>86</sup>, Cara Carty<sup>87</sup>, Richard Casaburi<sup>88</sup>, Juan P Casas Romero<sup>89</sup>, James  
416 Casella<sup>60</sup>, Peter Castaldi<sup>90</sup>, Mark Chaffin<sup>52</sup>, Christy Chang<sup>56</sup>, Yi-Cheng Chang<sup>91</sup>, Daniel  
417 Chasman<sup>92</sup>, Sameer Chavan<sup>84</sup>, Bo-Juen Chen<sup>50</sup>, Wei-Min Chen<sup>93</sup>, Michael Cho<sup>85</sup>,  
418 Seung Hoan Choi<sup>52</sup>, Lee-Ming Chuang<sup>94</sup>, Mina Chung<sup>95</sup>, Ren-Hua Chung<sup>96</sup>, Clary  
419 Clish<sup>97</sup>, Suzy Comhair<sup>98</sup>, Matthew Conomos<sup>83</sup>, Elaine Cornell<sup>99</sup>, Adolfo Correa<sup>100</sup>,  
420 Carolyn Crandall<sup>88</sup>, James Crapo<sup>101</sup>, L. Adrienne Cupples<sup>102</sup>, Jeffrey Curtis<sup>103</sup>, Brian  
421 Custer<sup>104</sup>, Coleen Damcott<sup>56</sup>, Dawood Darbar<sup>105</sup>, Sean David<sup>106</sup>, Colleen Davis<sup>57</sup>,  
422 Michelle Daya<sup>84</sup>, Michael DeBaun<sup>107</sup>, Dawn DeMeo<sup>85</sup>, Ranjan Deka<sup>108</sup>, Scott Devine<sup>56</sup>,  
423 Huyen Dinh<sup>66</sup>, Harsha Doddapaneni<sup>109</sup>, Qing Duan<sup>110</sup>, Shannon Dugan-Perez<sup>111</sup>, Ravi  
424 Duggirala<sup>112</sup>, Jon Peter Durda<sup>113</sup>, Susan K. Dutcher<sup>114</sup>, Charles Eaton<sup>115</sup>, Lynette  
425 Ekunwe<sup>58</sup>, Adel El Boueiz<sup>116</sup>, Patrick Ellinor<sup>117</sup>, Leslie Emery<sup>57</sup>, Serpil Erzurum<sup>118</sup>,  
426 Charles Farber<sup>93</sup>, Jesse Farek<sup>66</sup>, Tasha Fingerlin<sup>119</sup>, Matthew Flickinger<sup>51</sup>, Chris  
427 Frazar<sup>57</sup>, Mao Fu<sup>56</sup>, Stephanie M. Fullerton<sup>57</sup>, Lucinda Fulton<sup>120</sup>, Stacey Gabriel<sup>52</sup>,  
428 Weiniu Gan<sup>73</sup>, Shanshan Gao<sup>84</sup>, Yan Gao<sup>58</sup>, Margery Gass<sup>121</sup>, Heather Geiger<sup>122</sup>, Bruce  
429 Gelb<sup>123</sup>, Mark Geraci<sup>124</sup>, Soren Germer<sup>50</sup>, Robert Gerszten<sup>125</sup>, Auyon Ghosh<sup>85</sup>, Richard  
430 Gibbs<sup>66</sup>, Chris Gignoux<sup>63</sup>, Mark Gladwin<sup>126</sup>, David Glahn<sup>127</sup>, Stephanie Gogarten<sup>57</sup>, Da-  
431 Wei Gong<sup>56</sup>, Harald Goring<sup>128</sup>, Sharon Graw<sup>129</sup>, Kathryn J. Gray<sup>130</sup>, Daniel Grine<sup>84</sup>, Colin  
432 Gross<sup>51</sup>, Yue Guan<sup>56</sup>, Xiuqing Guo<sup>131</sup>, Namrata Gupta<sup>132</sup>, Jeff Haessler<sup>121</sup>, Michael  
433 Hall<sup>133</sup>, Yi Han<sup>66</sup>, Patrick Hanly<sup>134</sup>, Daniel Harris<sup>135</sup>, Nicola L. Hawley<sup>136</sup>, Ben Heavner<sup>83</sup>,  
434 Susan Heckbert<sup>137</sup>, Ryan Hernandez<sup>81</sup>, David Herrington<sup>138</sup>, Craig Hersh<sup>139</sup>, Bertha

435 Hidalgo<sup>62</sup>, James Hixson<sup>140</sup>, Brian Hobbs<sup>141</sup>, John Hokanson<sup>84</sup>, Elliott Hong<sup>56</sup>, Karin  
436 Hoth<sup>142</sup>, Chao (Agnes) Hsiung<sup>143</sup>, Jianhong Hu<sup>66</sup>, Haley Huston<sup>144</sup>, Chii Min Hwu<sup>145</sup>,  
437 Rebecca Jackson<sup>146</sup>, Deepti Jain<sup>57</sup>, Cashell Jaquish<sup>147</sup>, Jill Johnsen<sup>148</sup>, Andrew  
438 Johnson<sup>73</sup>, Craig Johnson<sup>57</sup>, Rich Johnston<sup>55</sup>, Kimberly Jones<sup>60</sup>, Hyun Min Kang<sup>149</sup>,  
439 Shannon Kelly<sup>150</sup>, Eimear Kenny<sup>123</sup>, Michael Kessler<sup>56</sup>, Alyna Khan<sup>57</sup>, Ziad Khan<sup>66</sup>,  
440 Wonji Kim<sup>151</sup>, John Kimoff<sup>152</sup>, Greg Kinney<sup>153</sup>, Barbara Konkle<sup>154</sup>, Holly Kramer<sup>155</sup>,  
441 Christoph Lange<sup>156</sup>, Ethan Lange<sup>84</sup>, Leslie Lange<sup>157</sup>, Cathy Laurie<sup>57</sup>, Cecelia Laurie<sup>57</sup>,  
442 Meryl LeBoff<sup>85</sup>, Jonathon LeFaive<sup>51</sup>, Jiwon Lee<sup>85</sup>, Sandra Lee<sup>66</sup>, Wen-Jane Lee<sup>145</sup>,  
443 David Levine<sup>57</sup>, Daniel Levy<sup>73</sup>, Joshua Lewis<sup>56</sup>, Xiaohui Li<sup>131</sup>, Yun Li<sup>110</sup>, Henry Lin<sup>131</sup>,  
444 Honghuang Lin<sup>158</sup>, Simin Liu<sup>159</sup>, Yongmei Liu<sup>160</sup>, Yu Liu<sup>161</sup>, Steven Lubitz<sup>117</sup>, Kathryn  
445 Lunetta<sup>162</sup>, James Luo<sup>73</sup>, Ulysses Magalang<sup>163</sup>, Barry Make<sup>60</sup>, Ani Manichaikul<sup>93</sup>, Alisa  
446 Manning<sup>164</sup>, JoAnn Manson<sup>85</sup>, Melissa Marton<sup>122</sup>, Susan Mathai<sup>84</sup>, Susanne May<sup>83</sup>,  
447 Patrick McArdle<sup>56</sup>, Merry-Lynn McDonald<sup>165</sup>, Sean McFarland<sup>151</sup>, Stephen McGarvey<sup>166</sup>,  
448 Daniel McGoldrick<sup>167</sup>, Caitlin McHugh<sup>83</sup>, Becky McNeil<sup>168</sup>, Hao Mei<sup>58</sup>, James Meigs<sup>169</sup>,  
449 Vipin Menon<sup>66</sup>, Luisa Mestroni<sup>129</sup>, Ginger Metcalf<sup>66</sup>, Deborah A Meyers<sup>170</sup>, Emmanuel  
450 Mignot<sup>171</sup>, Julie Mikulla<sup>73</sup>, Nancy Min<sup>58</sup>, Mollie Minear<sup>172</sup>, Matt Moll<sup>90</sup>, Zeineen Momin<sup>66</sup>,  
451 Courtney Montgomery<sup>173</sup>, Donna Muzny<sup>66</sup>, Josyf C Mychaleckyj<sup>93</sup>, Girish Nadkarni<sup>123</sup>,  
452 Rakhi Naik<sup>60</sup>, Take Naseri<sup>174</sup>, Sergei Nekhai<sup>175</sup>, Sarah C. Nelson<sup>83</sup>, Bonnie Neltner<sup>84</sup>,  
453 Caitlin Nessner<sup>66</sup>, Deborah Nickerson<sup>176</sup>, Osuji Nkechinyere<sup>66</sup>, Kari North<sup>110</sup>, Jeff  
454 O'Connell<sup>177</sup>, Tim O'Connor<sup>56</sup>, Heather Ochs-Balcom<sup>178</sup>, Geoffrey Okwuonu<sup>66</sup>, Allan  
455 Pack<sup>179</sup>, David T. Paik<sup>180</sup>, James Pankow<sup>181</sup>, George Papanicolaou<sup>73</sup>, Cora Parker<sup>182</sup>,  
456 Juan Manuel Peralta<sup>112</sup>, Marco Perez<sup>63</sup>, James Perry<sup>56</sup>, Ulrike Peters<sup>183</sup>, Lawrence S  
457 Phillips<sup>55</sup>, Jacob Pleiness<sup>51</sup>, Toni Pollin<sup>56</sup>, Wendy Post<sup>184</sup>, Julia Powers Becker<sup>185</sup>,

458 Meher Preethi Boorgula<sup>84</sup>, Michael Preuss<sup>123</sup>, Pankaj Qasba<sup>73</sup>, Dandi Qiao<sup>85</sup>, Zhaohui  
459 Qin<sup>55</sup>, Nicholas Rafaels<sup>186</sup>, Mahitha Rajendran<sup>66</sup>, D.C. Rao<sup>120</sup>, Laura Rasmussen-  
460 Torvik<sup>187</sup>, Aakrosh Ratan<sup>93</sup>, Robert Reed<sup>56</sup>, Catherine Reeves<sup>188</sup>, Elizabeth Regan<sup>189</sup>,  
461 Muagututi'a Sefuiva Reupena<sup>190</sup>, Rebecca Robillard<sup>191</sup>, Nicolas Robine<sup>122</sup>, Dan  
462 Roden<sup>192</sup>, Carolina Roselli<sup>52</sup>, Ingo Ruczinski<sup>60</sup>, Alexi Runnels<sup>122</sup>, Pamela Russell<sup>84</sup>,  
463 Sarah Ruuska<sup>193</sup>, Kathleen Ryan<sup>56</sup>, Ester Cerdeira Sabino<sup>194</sup>, Danish Saleheen<sup>195</sup>,  
464 Shabnam Salimi<sup>196</sup>, Sejal Salvi<sup>66</sup>, Steven Salzberg<sup>60</sup>, Kevin Sandow<sup>197</sup>, Vijay G.  
465 Sankaran<sup>198</sup>, Jireh Santibanez<sup>66</sup>, Karen Schwander<sup>120</sup>, David Schwartz<sup>84</sup>, Frank  
466 Scurba<sup>126</sup>, Christine Seidman<sup>199</sup>, Jonathan Seidman<sup>200</sup>, Vivien Sheehan<sup>201</sup>, Stephanie  
467 L. Sherman<sup>202</sup>, Amol Shetty<sup>56</sup>, Aniket Shetty<sup>84</sup>, Wayne Hui-Heng Sheu<sup>145</sup>, M. Benjamin  
468 Shoemaker<sup>203</sup>, Brian Silver<sup>204</sup>, Edwin Silverman<sup>85</sup>, Robert Skomro<sup>205</sup>, Albert Vernon  
469 Smith<sup>206</sup>, Josh Smith<sup>57</sup>, Nicholas Smith<sup>207</sup>, Tanja Smith<sup>50</sup>, Sylvia Smoller<sup>208</sup>, Beverly  
470 Snively<sup>209</sup>, Michael Snyder<sup>210</sup>, Tamar Sofer<sup>125</sup>, Nona Sotoodehnia<sup>57</sup>, Adrienne M. Stilp<sup>57</sup>,  
471 Garrett Storm<sup>211</sup>, Elizabeth Streeten<sup>56</sup>, Jessica Lasky Su<sup>212</sup>, Yun Ju Sung<sup>120</sup>, Jody  
472 Sylvia<sup>85</sup>, Adam Szpiro<sup>57</sup>, Frédéric Sériès<sup>213</sup>, Daniel Taliun<sup>51</sup>, Hua Tang<sup>210</sup>, Margaret  
473 Taub<sup>60</sup>, Matthew Taylor<sup>129</sup>, Simeon Taylor<sup>56</sup>, Marilyn Telen<sup>61</sup>, Timothy A. Thornton<sup>57</sup>,  
474 Machiko Threlkeld<sup>214</sup>, Lesley Tinker<sup>215</sup>, David Tirschwell<sup>57</sup>, Sarah Tishkoff<sup>216</sup>, Catherine  
475 Tong<sup>217</sup>, Russell Tracy<sup>218</sup>, Michael Tsai<sup>181</sup>, Dhananjay Vaidya<sup>60</sup>, David Van Den Berg<sup>219</sup>,  
476 Peter VandeHaar<sup>51</sup>, Scott Vrieze<sup>181</sup>, Tarik Walker<sup>84</sup>, Robert Wallace<sup>142</sup>, Avram Walts<sup>84</sup>,  
477 Fei Fei Wang<sup>57</sup>, Heming Wang<sup>220</sup>, Jiongming Wang<sup>221</sup>, Karol Watson<sup>88</sup>, Jennifer Watt<sup>66</sup>,  
478 Daniel E. Weeks<sup>222</sup>, Joshua Weinstock<sup>149</sup>, Bruce Weir<sup>57</sup>, Scott T Weiss<sup>223</sup>, Lu-Chen  
479 Weng<sup>117</sup>, Jennifer Wessel<sup>224</sup>, Cristen Willer<sup>103</sup>, Kayleen Williams<sup>83</sup>, L. Keoki Williams<sup>225</sup>,  
480 Scott Williams<sup>226</sup>, Carla Wilson<sup>85</sup>, James Wilson<sup>227</sup>, Lara Winterkorn<sup>122</sup>, Quenna

481 Wong<sup>57</sup>, Baojun Wu<sup>228</sup>, Joseph Wu<sup>180</sup>, Huichun Xu<sup>56</sup>, Ivana Yang<sup>84</sup>, Ketian Yu<sup>51</sup>,  
482 Seyede Maryam Zekavat<sup>52</sup>, Yingze Zhang<sup>229</sup>, Snow Xueyan Zhao<sup>101</sup>, Wei Zhao<sup>230</sup>,  
483 Xiaofeng Zhu<sup>231</sup>, Elad Ziv<sup>232</sup>, Michael Zody<sup>50</sup>, Sebastian Zoellner<sup>51</sup>, Mariza de  
484 Andrade<sup>233</sup>, Lisa de las Fuentes<sup>234</sup>  
485  
486 50 - New York Genome Center, New York, New York, 10013, US; 51 - University of  
487 Michigan, Ann Arbor, Michigan, 48109, US; 52 - Broad Institute, Cambridge,  
488 Massachusetts, 2142, US; 53 - Cedars Sinai, Boston, Massachusetts, 2114, US; 54 -  
489 Children's Hospital of Philadelphia, University of Pennsylvania, Philadelphia,  
490 Pennsylvania, 19104, US; 55 - Emory University, Atlanta, Georgia, 30322, US; 56 -  
491 University of Maryland, Baltimore, Maryland, 21201, US; 57 - University of Washington,  
492 Seattle, Washington, 98195, US; 58 - University of Mississippi, Jackson, Mississippi,  
493 38677, US; 59 - National Institutes of Health, Bethesda, Maryland, 20892, US; 60 -  
494 Johns Hopkins University, Baltimore, Maryland, 21218, US; 61 - Duke University,  
495 Durham, North Carolina, 27708, US; 62 - University of Alabama, Birmingham, Alabama,  
496 35487, US; 63 - Stanford University, Stanford, California, 94305, US; 64 - Medical  
497 College of Wisconsin, Milwaukee, Wisconsin, 53211, US; 65 - Providence Health Care,  
498 Medicine, Vancouver, CA; 66 - Baylor College of Medicine Human Genome Sequencing  
499 Center, Houston, Texas, 77030, US; 67 - Cleveland Clinic, Cleveland, Ohio, 44195, US;  
500 68 - Tempus, University of Colorado Anschutz Medical Campus, Aurora, Colorado,  
501 80045, US; 69 - Columbia University, New York, New York, 10032, US; 70 - The  
502 Emmes Corporation, LTRC, Rockville, Maryland, 20850, US; 71 - Cleveland Clinic,  
503 Quantitative Health Sciences, Cleveland, Ohio, 44195, US; 72 - Johns Hopkins

504 University, Medicine, Baltimore, Maryland, 21218, US; 73 - National Heart, Lung, and  
505 Blood Institute, National Institutes of Health, Bethesda, Maryland, 20892, US; 74 -  
506 Boston University, Massachusetts General Hospital, Boston University School of  
507 Medicine, Boston, Massachusetts, 2114, US; 75 - University of Florida, Epidemiology,  
508 Gainesville, Florida, 32610, US; 76 - Fundação de Hematologia e Hemoterapia de  
509 Pernambuco - Hemope, Recife, 52011-000, BR; 77 - University of Utah, Obstetrics and  
510 Gynecology, Salt Lake City, Utah, 84132, US; 78 - National Jewish Health, National  
511 Jewish Health, Denver, Colorado, 80206, US; 79 - Medical College of Wisconsin,  
512 Pediatrics, Milwaukee, Wisconsin, 53226, US; 80 - University of Texas Health at  
513 Houston, Pediatrics, Houston, Texas, 77030, US; 81 - University of California, San  
514 Francisco, San Francisco, California, 94143, US; 82 - Stanford University, Biomedical  
515 Data Science, Stanford, California, 94305, US; 83 - University of Washington,  
516 Biostatistics, Seattle, Washington, 98195, US; 84 - University of Colorado at Denver,  
517 Denver, Colorado, 80204, US; 85 - Brigham & Women's Hospital, Boston,  
518 Massachusetts, 2115, US; 86 - University of Montreal, US; 87 - Washington State  
519 University, Pullman, Washington, 99164, US; 88 - University of California, Los Angeles,  
520 Los Angeles, California, 90095, US; 89 - Brigham & Women's Hospital, US; 90 -  
521 Brigham & Women's Hospital, Medicine, Boston, Massachusetts, 2115, US; 91 -  
522 National Taiwan University, Taipei, 10617, TW; 92 - Brigham & Women's Hospital,  
523 Division of Preventive Medicine, Boston, Massachusetts, 2215, US; 93 - University of  
524 Virginia, Charlottesville, Virginia, 22903, US; 94 - National Taiwan University, National  
525 Taiwan University Hospital, Taipei, 10617, TW; 95 - Cleveland Clinic, Cleveland Clinic,  
526 Cleveland, Ohio, 44195, US; 96 - National Health Research Institute Taiwan, Miaoli

527 County, 350, TW; 97 - Broad Institute, Metabolomics Platform, Cambridge,  
528 Massachusetts, 2142, US; 98 - Cleveland Clinic, Immunity and Immunology, Cleveland,  
529 Ohio, 44195, US; 99 - University of Vermont, Burlington, Vermont, 5405, US; 100 -  
530 University of Mississippi, Population Health Science, Jackson, Mississippi, 39216, US;  
531 101 - National Jewish Health, Denver, Colorado, 80206, US; 102 - Boston University,  
532 Biostatistics, Boston, Massachusetts, 2115, US; 103 - University of Michigan, Internal  
533 Medicine, Ann Arbor, Michigan, 48109, US; 104 - Vitalant Research Institute, San  
534 Francisco, California, 94118, US; 105 - University of Illinois at Chicago, Chicago, Illinois,  
535 60607, US; 106 - University of Chicago, Chicago, Illinois, 60637, US; 107 - Vanderbilt  
536 University, Nashville, Tennessee, 37235, US; 108 - University of Cincinnati, Cincinnati,  
537 Ohio, 45220, US; 109 - Baylor College of Medicine Human Genome Sequencing  
538 Center, Houston, Texas, 77030; 110 - University of North Carolina, Chapel Hill, North  
539 Carolina, 27599, US; 111 - Baylor College of Medicine Human Genome Sequencing  
540 Center, BCM, Houston, Texas, 77030, US; 112 - University of Texas Rio Grande Valley  
541 School of Medicine, Edinburg, Texas, 78539, US; 113 - University of Vermont,  
542 Pathology and Laboratory Medicine, Burlington, Vermont, 5405, US; 114 - Washington  
543 University in St Louis, Genetics, St Louis, Missouri, 63110, US; 115 - Brown University,  
544 Providence, Rhode Island, 2912, US; 116 - Harvard University, Channing Division of  
545 Network Medicine, Cambridge, Massachusetts, 2138, US; 117 - Massachusetts General  
546 Hospital, Boston, Massachusetts, 2114, US; 118 - Cleveland Clinic, Lerner Research  
547 Institute, Cleveland, Ohio, 44195, US; 119 - National Jewish Health, Center for Genes,  
548 Environment and Health, Denver, Colorado, 80206, US; 120 - Washington University in  
549 St Louis, St Louis, Missouri, 63130, US; 121 - Fred Hutchinson Cancer Research

550 Center, Seattle, Washington, 98109, US; 122 - New York Genome Center, New York  
551 City, New York, 10013, US; 123 - Icahn School of Medicine at Mount Sinai, New York,  
552 New York, 10029, US; 124 - University of Pittsburgh, Pittsburgh, Pennsylvania, US; 125  
553 - Beth Israel Deaconess Medical Center, Boston, Massachusetts, 2215, US; 126 -  
554 University of Pittsburgh, Pittsburgh, Pennsylvania, 15260, US; 127 - Boston Children's  
555 Hospital, Harvard Medical School, Department of Psychiatry, Boston, Massachusetts,  
556 2115, US; 128 - University of Texas Rio Grande Valley School of Medicine, San  
557 Antonio, Texas, 78229, US; 129 - University of Colorado Anschutz Medical Campus,  
558 Aurora, Colorado, 80045, US; 130 - Mass General Brigham, Obstetrics and  
559 Gynecology, Boston, Massachusetts, 2115, US; 131 - Lundquist Institute, Torrance,  
560 California, 90502, US; 132 - Broad Institute, Broad Institute, Cambridge,  
561 Massachusetts, 2142, US; 133 - University of Mississippi, Cardiology, Jackson,  
562 Mississippi, 39216, US; 134 - University of Calgary, Medicine, Calgary, CA; 135 -  
563 University of Maryland, Genetics, Philadelphia, Pennsylvania, 19104, US; 136 - Yale  
564 University, Department of Chronic Disease Epidemiology, New Haven, Connecticut,  
565 6520, US; 137 - University of Washington, Epidemiology, Seattle, Washington, 98195-  
566 9458, US; 138 - Wake Forest Baptist Health, Winston-Salem, North Carolina, 27157,  
567 US; 139 - Brigham & Women's Hospital, Channing Division of Network Medicine,  
568 Boston, Massachusetts, 2115, US; 140 - University of Texas Health at Houston,  
569 Houston, Texas, 77225, US; 141 - Regeneron Genetics Center, Boston,  
570 Massachusetts, 2115, US; 142 - University of Iowa, Iowa City, Iowa, 52242, US; 143 -  
571 National Health Research Institute Taiwan, Institute of Population Health Sciences,  
572 NHRI, Miaoli County, 350, TW; 144 - Blood Works Northwest, Seattle, Washington,

573 98104, US; 145 - Taichung Veterans General Hospital Taiwan, Taichung City, 407, TW;  
574 146 - Oklahoma State University Medical Center, Internal Medicine, Division of  
575 Endocrinology, Diabetes and Metabolism, Columbus, Ohio, 43210, US; 147 - National  
576 Heart, Lung, and Blood Institute, National Institutes of Health, NHLBI, Bethesda,  
577 Maryland, 20892, US; 148 - University of Washington, Medicine, Seattle, Washington,  
578 98109, US; 149 - University of Michigan, Biostatistics, Ann Arbor, Michigan, 48109, US;  
579 150 - University of California, San Francisco, San Francisco, California, 94118, US; 151  
580 - Harvard University, Cambridge, Massachusetts, 2138, US; 152 - McGill University,  
581 Montréal, QC H3A 0G4, CA; 153 - University of Colorado at Denver, Epidemiology,  
582 Aurora, Colorado, 80045, US; 154 - Blood Works Northwest, Medicine, Seattle,  
583 Washington, 98104, US; 155 - Loyola University, Public Health Sciences, Maywood,  
584 Illinois, 60153, US; 156 - Harvard School of Public Health, Biostats, Boston,  
585 Massachusetts, 2115, US; 157 - University of Colorado at Denver, Medicine, Aurora,  
586 Colorado, 80048, US; 158 - Boston University, University of Massachusetts Chan  
587 Medical School, Worcester, Massachusetts, 1655, US; 159 - Brown University,  
588 Epidemiology and Medicine, Providence, Rhode Island, 2912, US; 160 - Duke  
589 University, Cardiology, Durham, North Carolina, 27708, US; 161 - Stanford University,  
590 Cardiovascular Institute, Stanford, California, 94305, US; 162 - Boston University,  
591 Boston, Massachusetts, 2215, US; 163 - The Ohio State University, Division of  
592 Pulmonary, Critical Care and Sleep Medicine, Columbus, Ohio, 43210, US; 164 - Broad  
593 Institute, Harvard University, Massachusetts General Hospital; 165 - University of  
594 Alabama, University of Alabama at Birmingham, Birmingham, Alabama, 35487, US; 166  
595 - Brown University, Epidemiology, Providence, Rhode Island, 2912, US; 167 - University



596 of Washington, Genome Sciences, Seattle, Washington, 98195, US; 168 - RTI  
597 International, US; 169 - Massachusetts General Hospital, Medicine, Boston,  
598 Massachusetts, 2114, US; 170 - University of Arizona, Tucson, Arizona, 85721, US; 171  
599 - Stanford University, Center For Sleep Sciences and Medicine, Palo Alto, California,  
600 94304, US; 172 - National Institute of Child Health and Human Development, National  
601 Institutes of Health, Bethesda, Maryland, 20892, US; 173 - Oklahoma Medical Research  
602 Foundation, Genes and Human Disease, Oklahoma City, Oklahoma, 73104, US; 174 -  
603 Ministry of Health, Government of Samoa, Apia, WS; 175 - Howard University,  
604 Washington, District of Columbia, 20059, US; 176 - University of Washington,  
605 Department of Genome Sciences, Seattle, Washington, 98195, US; 177 - University of  
606 Maryland, Balitmore, Maryland, 21201, US; 178 - University at Buffalo, Buffalo, New  
607 York, 14260, US; 179 - University of Pennsylvania, Division of Sleep  
608 Medicine/Department of Medicine, Philadelphia, Pennsylvania, 19104-3403, US; 180 -  
609 Stanford University, Stanford Cardiovascular Institute, Stanford, California, 94305, US;  
610 181 - University of Minnesota, Minneapolis, Minnesota, 55455, US; 182 - RTI  
611 International, Biostatistics and Epidemiology Division, Research Triangle Park, North  
612 Carolina, 27709-2194, US; 183 - Fred Hutchinson Cancer Research Center, Fred Hutch  
613 and UW, Seattle, Washington, 98109, US; 184 - Johns Hopkins University,  
614 Cardiology/Medicine, Baltimore, Maryland, 21218, US; 185 - University of Colorado at  
615 Denver, Medicine, Denver, Colorado, 80204, US; 186 - University of Colorado at  
616 Denver, CCPM, Denver, Colorado, 80045, US; 187 - Northwestern University, Chicago,  
617 Illinois, 60208, US; 188 - New York Genome Center, New York Genome Center, New  
618 York City, New York, 10013, US; 189 - National Jewish Health, Medicine, Denver,

619 Colorado, 80206, US; 190 - Lutia I Puava Ae Mapu I Fagalele, Apia, WS; 191 -  
620 University of Ottawa, Sleep Research Unit, University of Ottawa Institute for Mental  
621 Health Research, Ottawa, ON K1Z 7K4, CA; 192 - Vanderbilt University, Medicine,  
622 Pharmacology, Biomedical Informatics, Nashville, Tennessee, 37235, US; 193 -  
623 University of Washington, Seattle, Washington, 98104, US; 194 - Universidade de Sao  
624 Paulo, Faculdade de Medicina, Sao Paulo, 1310000, BR; 195 - Columbia University,  
625 New York, New York, 10027, US; 196 - University of Maryland, Pathology, Seattle,  
626 Washington, 98195, US; 197 - Lundquist Institute, TGPS, Torrance, California, 90502,  
627 US; 198 - Harvard University, Division of Hematology/Oncology, Boston,  
628 Massachusetts, 2115, US; 199 - Harvard Medical School, Genetics, Boston,  
629 Massachusetts, 2115, US; 200 - Harvard Medical School, Boston, Massachusetts,  
630 2115, US; 201 - Emory University, Pediatrics, Atlanta, Georgia, 30307, US; 202 - Emory  
631 University, Human Genetics, Atlanta, Georgia, 30322, US; 203 - Vanderbilt University,  
632 Medicine/Cardiology, Nashville, Tennessee, 37235, US; 204 - UMass Memorial Medical  
633 Center, Worcester, Massachusetts, 1655, US; 205 - University of Saskatchewan,  
634 Saskatoon, SK S7N 5C9, CA; 206 - University of Michigan; 207 - University of  
635 Washington, Epidemiology, Seattle, Washington, 98195, US; 208 - Albert Einstein  
636 College of Medicine, New York, New York, 10461, US; 209 - Wake Forest Baptist  
637 Health, Biostatistical Sciences, Winston-Salem, North Carolina, 27157, US; 210 -  
638 Stanford University, Genetics, Stanford, California, 94305, US; 211 - University of  
639 Colorado at Denver, Genomic Cardiology, Aurora, Colorado, 80045, US; 212 - Brigham  
640 & Women's Hospital, Channing Department of Medicine, Boston, Massachusetts, 2115,  
641 US; 213 - Université Laval, Quebec City, G1V 0A6, CA; 214 - University of Washington,

642 University of Washington, Department of Genome Sciences, Seattle, Washington,  
643 98195, US; 215 - Fred Hutchinson Cancer Research Center, Cancer Prevention  
644 Division of Public Health Sciences, Seattle, Washington, 98109, US; 216 - University of  
645 Pennsylvania, Genetics, Philadelphia, Pennsylvania, 19104, US; 217 - University of  
646 Washington, Department of Biostatistics, Seattle, Washington, 98195, US; 218 -  
647 University of Vermont, Pathology & Laboratory Medicine, Burlington, Vermont, 5405,  
648 US; 219 - University of Southern California, USC Methylation Characterization Center,  
649 University of Southern California, California, 90033, US; 220 - Brigham & Women's  
650 Hospital, Mass General Brigham, Boston, Massachusetts, 2115, US; 221 - University of  
651 Michigan, US; 222 - University of Pittsburgh, Department of Human Genetics,  
652 Pittsburgh, Pennsylvania, 15260, US; 223 - Brigham & Women's Hospital, Channing  
653 Division of Network Medicine, Department of Medicine, Boston, Massachusetts, 2115,  
654 US; 224 - Indiana University, Epidemiology, Indianapolis, Indiana, 46202, US; 225 -  
655 Henry Ford Health System, Detroit, Michigan, 48202, US; 226 - Case Western Reserve  
656 University; 227 - Beth Israel Deaconess Medical Center, Cardiology, Cambridge,  
657 Massachusetts, 2139, US; 228 - Henry Ford Health System, Department of Medicine,  
658 Detroit, Michigan, 48202, US; 229 - University of Pittsburgh, Medicine, Pittsburgh,  
659 Pennsylvania, 15260, US; 230 - University of Michigan, Department of Epidemiology,  
660 Ann Arbor, Michigan, 48109, US; 231 - Case Western Reserve University, Department  
661 of Population and Quantitative Health Sciences, Cleveland, Ohio, 44106, US; 232 -  
662 University of California, San Francisco, Medicine, San Francisco, California, 94143, US;  
663 233 - Mayo Clinic, Health Quantitative Sciences Research, Rochester, Minnesota,

664 55905, US; 234 - Washington University in St Louis, Department of Medicine,

665 Cardiovascular Division, St. Louis, Missouri, 63110, US

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687 **TABLES**

688 **Table 1 | TOPMed Gene-centric coding multi-trait analysis results of both**  
 689 **unconditional analysis and analysis conditional on known lipids-associated**  
 690 **variants.** A total of 61,838 samples from the TOPMed Program were considered in the  
 691 analysis. Results for the conditionally significant genes (unconditional MultiSTAAR-O  
 692  $P < 5.00 \times 10^{-7}$ ; conditional MultiSTAAR-O  $P < 9.80 \times 10^{-4}$ ) are presented in the  
 693 table. MultiSTAAR-O is a two-sided test. Chr. no., chromosome number; Category,  
 694 functional category; No. of SNVs, number of rare variants (MAF < 1%) of the particular  
 695 coding functional category in the gene; MultiSTAAR-O, MultiSTAAR-O  $P$  value; Variants  
 696 (adjusted), adjusted variants in the conditional analysis.

697

Gene	Chr. no.	Category	No. of SNVs	MultiSTAAR-O (Unconditional)	MultiSTAAR-O (Conditional)	Variants (adjusted)
<i>PCSK9</i>	1	Putative loss-of-function	14	1.14E-115	2.66E-08	rs12117661, rs2495491, rs11591147, rs67608943, rs72646508, rs693668, rs28362261, rs28362263, rs141502002, rs505151, rs28362286
<i>APOB</i>	2	Putative loss-of-function	29	8.04E-28	5.76E-27	rs12478327, rs72654432, rs1042034, rs676210, rs533617, rs17240441, rs34722314, rs563290, rs10692845
<i>ABCA1</i>	9	Putative loss-of-function	28	2.04E-21	5.41E-21	rs2150867, rs33918808, rs112853430, rs4149307, rs9282541, rs1883025, rs1800978
<i>LDLR</i>	19	Putative loss-of-function	19	8.81E-21	7.16E-21	rs140753491, rs138294113, rs17242353, rs17242843, rs10422256, rs72658860, rs11669576, rs2738447, rs72658867, rs2738464, rs6511728, rs3760782, rs59168178, rs2278426, rs112942459
<i>PCSK9</i>	1	Missense	271	8.94E-71	1.29E-10	rs12117661, rs2495491, rs11591147, rs67608943, rs72646508, rs693668, rs28362261, rs28362263, rs141502002, rs505151, rs28362286
<i>APOB</i>	2	Missense	1407	5.57E-08	4.31E-08	rs12478327, rs72654432, rs1042034, rs676210, rs533617, rs17240441, rs34722314, rs563290, rs10692845
<i>ABCG5</i>	2	Missense	242	5.75E-08	9.81E-08	rs114780578, rs11887534, rs4245791
<i>NPC1L1</i>	7	Missense	477	3.10E-08	1.60E-07	rs217381
<i>LPL</i>	8	Missense	149	9.57E-19	7.14E-04	rs6996383, rs268, rs328, rs3289, rs13702, rs15285, rs78810414, rs28550053, rs12676079, rs55682243
<i>ABCA1</i>	9	Missense	597	3.63E-46	1.75E-33	rs2150867, rs33918808, rs112853430, rs4149307, rs9282541, rs1883025, rs1800978
<i>SCARB1</i>	12	Missense	192	6.77E-15	3.55E-15	rs6488913, rs4765127, rs1716407, rs825456, rs1672875, rs10846744, rs10773112, rs187471874, rs10773119, rs1973688, rs1601935, rs2043082, rs10468017, rs1532085, rs436965, rs35980001, rs1800588, rs2070895, rs113298164, rs35571500, rs247617, rs17231506, rs34498052, rs34119551, rs34065661, rs1597000001*, rs7499892, rs5883, rs289719, rs11860407, rs189866004,
<i>LIPC</i>	15	Missense	246	2.54E-20	6.66E-15	
<i>CETP</i>	16	Missense	168	8.84E-14	2.09E-04	

						rs5880
<i>LCAT</i>	16	Missense	107	9.18E-14	3.06E-17	rs111315946, rs150660813, rs4986970, rs35673026, rs1109166, rs548291389, rs140753491, rs138294113, rs17242353, rs17242843, rs10422256, rs72658860, rs11669576, rs2738447, rs72658867, rs2738464, rs6511728, rs3760782, rs59168178, rs2278426, rs112942459
<i>LDLR</i>	19	Missense	342	7.92E-58	2.12E-57	rs3761077, rs150641967, rs187429064, rs2074304
<i>TM6SF2</i>	19	Missense	120	7.06E-08	6.16E-07	rs12117661, rs2495491, rs11591147, rs67608943, rs72646508, rs693668, rs28362261, rs28362263, rs141502002, rs505151, rs28362286
<i>PCSK9</i>	1	Putative loss-of-function and disruptive missense	71	1.14E-107	8.22E-17	
<i>APOB</i>	2	Putative loss-of-function and disruptive missense	75	9.96E-12	9.86E-12	rs12478327, rs72654432, rs1042034, rs676210, rs533617, rs17240441, rs34722314, rs563290, rs10692845
<i>NPC1L1</i>	7	Putative loss-of-function and disruptive missense	303	1.79E-09	8.29E-09	rs217381
<i>ABCA1</i>	9	Putative loss-of-function and disruptive missense	357	7.85E-33	2.66E-33	rs2150867, rs33918808, rs112853430, rs4149307, rs9282541, rs1883025, rs1800978
<i>APOC3</i>	11	Putative loss-of-function and disruptive missense	15	2.86E-126	3.01E-06	rs509728, rs61905072, rs66505542, rs7102314, rs964184, rs75198898, rs142958146, rs2075291, rs3135506, rs651821, rs45611741, rs662799, rs10750097, rs9804646, rs978880643, rs2070669, rs76353203, rs138326449, rs147210663, rs140621530, rs525028, rs141469619, rs188287950, rs202207736
<i>SCARB1</i>	12	Putative loss-of-function and disruptive missense	60	3.49E-17	2.14E-17	rs6488913, rs4765127, rs1716407, rs825456, rs1672875, rs10846744, rs10773112, rs187471874, rs10773119
<i>LIPC</i>	15	Putative loss-of-function and disruptive missense	130	1.01E-19	1.49E-17	rs1973688, rs1601935, rs2043082, rs10468017, rs1532085, rs436965, rs35980001, rs1800588, rs2070895, rs113298164
<i>LCAT</i>	16	Putative loss-of-function and disruptive missense	88	2.38E-16	5.07E-17	rs111315946, rs150660813, rs4986970, rs35673026, rs1109166, rs548291389
<i>LDLR</i>	19	Putative loss-of-function and disruptive missense	221	6.97E-72	1.57E-71	rs140753491, rs138294113, rs17242353, rs17242843, rs10422256, rs72658860, rs11669576, rs2738447, rs72658867, rs2738464, rs6511728, rs3760782, rs59168178, rs2278426, rs112942459
<i>PCSK9</i>	1	Disruptive missense	57	7.03E-19	1.33E-12	rs12117661, rs2495491, rs11591147, rs67608943, rs72646508, rs693668, rs28362261, rs28362263, rs141502002, rs505151, rs28362286
<i>APOB</i>	2	Disruptive missense	46	5.78E-09	4.48E-09	rs12478327, rs72654432, rs1042034, rs676210, rs533617, rs17240441, rs34722314, rs563290, rs10692845
<i>NPC1L1</i>	7	Disruptive missense	276	3.34E-09	1.57E-08	rs217381
<i>ABCA1</i>	9	Disruptive missense	329	1.17E-22	1.59E-23	rs2150867, rs33918808, rs112853430, rs4149307, rs9282541, rs1883025, rs1800978
<i>APOC3</i>	11	Disruptive missense	6	2.38E-29	3.93E-04	rs509728, rs61905072, rs66505542, rs7102314, rs964184, rs75198898, rs142958146, rs2075291, rs3135506,

						rs651821, rs45611741, rs662799, rs10750097, rs9804646, rs978880643, rs2070669, rs76353203, rs138326449, rs147210663, rs140621530, rs525028, rs141469619, rs188287950, rs202207736
<i>SCARB1</i>	12	Disruptive missense	51	4.44E-16	2.86E-16	rs6488913, rs4765127, rs1716407, rs825456, rs1672875, rs10846744, rs10773112, rs187471874, rs10773119, rs1973688, rs1601935, rs2043082, rs10468017, rs1532085, rs436965, rs35980001, rs1800588, rs2070895, rs113298164, rs111315946, rs150660813, rs4986970, rs35673026, rs1109166, rs548291389, rs140753491, rs138294113, rs17242353, rs17242843, rs10422256, rs72658860, rs11669576, rs2738447, rs72658867, rs2738464, rs6511728, rs3760782, rs59168178, rs2278426, rs112942459
<i>LIPC</i>	15	Disruptive missense	112	2.19E-18	2.65E-16	
<i>LCAT</i>	16	Disruptive missense	84	2.85E-14	6.44E-15	
<i>LDLR</i>	19	Disruptive missense	203	2.22E-59	5.13E-59	

698 \* Samoan-specific missense variant.

699

700

701

702

703

704

705 **Table 2 | TOPMed Gene-centric noncoding multi-trait analysis results of both**  
 706 **unconditional analysis and analysis conditional on known lipids-associated**  
 707 **variants.** A total of 61,838 samples from the TOPMed Program were considered in the  
 708 analysis. Results for the conditionally significant genes (unconditional MultiSTAAR-O  
 709  $P < 3.57 \times 10^{-7}$  and conditional MultiSTAAR-O  $P < 6.58 \times 10^{-4}$  for 7 different  
 710 noncoding masks across protein-coding genes; unconditional MultiSTAAR-O  $P <$   
 711  $2.50 \times 10^{-6}$  and conditional MultiSTAAR-O  $P < 8.33 \times 10^{-3}$  for ncRNA genes) are  
 712 presented in the table. MultiSTAAR-O is a two-sided test. Chr. no., chromosome  
 713 number; Category, functional category; No. of SNVs, number of rare variants (MAF <  
 714 1%) of the particular noncoding functional category in the gene; MultiSTAAR-O,  
 715 MultiSTAAR-O  $P$  value; Variants (adjusted), adjusted variants in the conditional  
 716 analysis; n/a, no variant adjusted in the conditional analysis.  
 717

Gene	Chr. no.	Category	No. of SNVs	MultiSTAAR-O (Unconditional)	MultiSTAAR-O (Conditional)	Variants (adjusted)
<i>APOA1</i>	11	Promoter (CAGE)	230	2.33E-07	9.45E-07	rs509728, rs61905072, rs66505542, rs7102314, rs964184, rs75198898, rs142958146, rs2075291, rs3135506, rs651821, rs45611741, rs662799, rs10750097, rs9804646, rs978880643, rs2070669, rs76353203, rs138326449, rs147210663, rs140621530, rs525028, rs141469619, rs188287950, rs202207736
<i>CETP</i>	16	Promoter (DHS)	411	1.21E-12	5.75E-04	rs35571500, rs247617, rs17231506, rs34498052, rs34119551, rs34065661, rs1597000001*, rs7499892, rs5883, rs289719, rs11860407, rs189866004, rs5880
<i>APOA1</i>	11	Enhancer (CAGE)	642	1.88E-24	6.23E-04	rs509728, rs61905072, rs66505542, rs7102314, rs964184, rs75198898, rs142958146, rs2075291, rs3135506, rs651821, rs45611741, rs662799, rs10750097, rs9804646, rs978880643, rs2070669, rs76353203, rs138326449, rs147210663, rs140621530, rs525028, rs141469619, rs188287950, rs202207736
<i>SPC24</i>	19	Enhancer (CAGE)	366	1.33E-08	4.88E-04	rs140753491, rs138294113, rs17242353, rs17242843, rs10422256, rs72658860, rs11669576, rs2738447, rs72658867, rs2738464, rs6511728, rs3760782, rs59168178, rs2278426, rs112942459
<i>NIPSNA P3A</i>	9	Enhancer (DHS)	767	2.63E-08	8.46E-06	rs2150867, rs33918808, rs112853430, rs4149307, rs9282541, rs1883025, rs1800978
<i>LIPC</i>	15	Enhancer (DHS)	3714	4.26E-08	1.25E-04	rs1973688, rs1601935, rs2043082, rs10468017, rs1532085, rs436965, rs35980001, rs1800588, rs2070895, rs113298164
<i>RP11-310H4.2</i>	7	ncRNA	154	1.69E-06	1.69E-06	n/a



<i>MIR4497</i>	12	ncRNA	23	1.37E-06	1.42E-06	rs5800864
<i>RP11-15F12.3</i>	18	ncRNA	64	7.53E-11	7.50E-03	rs77960347, rs117623631, rs9958734, rs7229562, rs8086351, rs10048323, rs8084172

---

718 \* Samoan-specific missense variant.

719

720

721

722

723

724

725

726

727

728

729 **Table 3 | TOPMed Genetic region (2-kb sliding window) multi-trait analysis results**  
 730 **of both unconditional analysis and analysis conditional on known lipid-**  
 731 **associated variants.** A total of 61,838 samples from the TOPMed Program were  
 732 considered in the analysis. Results for the conditionally significant sliding windows  
 733 (unconditional MultiSTAAR-O  $P < 1.89 \times 10^{-8}$  and conditional MultiSTAAR-O  $P <$   
 734  $9.96 \times 10^{-5}$ ) are presented in the table. MultiSTAAR-O is a two-sided test. Chr. no.,  
 735 chromosome number; Start location, start location of the 2-kb sliding window; End  
 736 location, end location of the 2-kb sliding window; No. of SNVs, number of rare variants  
 737 (MAF < 1%) in the 2-kb sliding window; MultiSTAAR-O, MultiSTAAR-O  $P$  value;  
 738 Variants (adjusted), adjusted variants in the conditional analysis; n/a, no variant  
 739 adjusted in the conditional analysis. Physical positions of each window are on build  
 740 hg38.

Chr. no.	Start location	End location	Gene	No. of SNVs	MultiSTAAR-O (Unconditional)	MultiSTAAR-O (Conditional)	Variants (adjusted)
1	55,051,447	55,053,446	<i>PCSK9</i>	327	7.11E-11	6.60E-08	rs12117661, rs2495491, rs11591147, rs67608943, rs72646508, rs693668, rs28362261, rs28362263, rs141502002, rs505151, rs28362286
1	55,052,447	55,054,446	<i>PCSK9</i>	320	9.37E-09	9.07E-06	rs12117661, rs2495491, rs11591147, rs67608943, rs72646508, rs693668, rs28362261, rs28362263, rs141502002, rs505151, rs28362286
1	62,651,447	62,653,446	<i>DOCK7</i>	277	5.08E-09	7.56E-10	rs67461605
1	62,652,447	62,654,446	<i>DOCK7</i>	257	4.87E-09	7.24E-10	rs67461605
1	145,530,447	145,532,446	<i>intergenic</i>	233	5.12E-09	5.12E-09	n/a
19	11,104,367	11,106,366	<i>LDLR</i>	336	1.15E-12	8.33E-13	rs140753491, rs138294113, rs17242353, rs17242843, rs10422256, rs72658860, rs11669576, rs2738447, rs72658867, rs2738464, rs6511728, rs3760782, rs59168178, rs2278426, rs112942459
19	11,105,367	11,107,366	<i>LDLR</i>	338	5.97E-14	5.55E-15	rs140753491, rs138294113, rs17242353, rs17242843, rs10422256, rs72658860, rs11669576, rs2738447, rs72658867, rs2738464, rs6511728, rs3760782, rs59168178, rs2278426, rs112942459

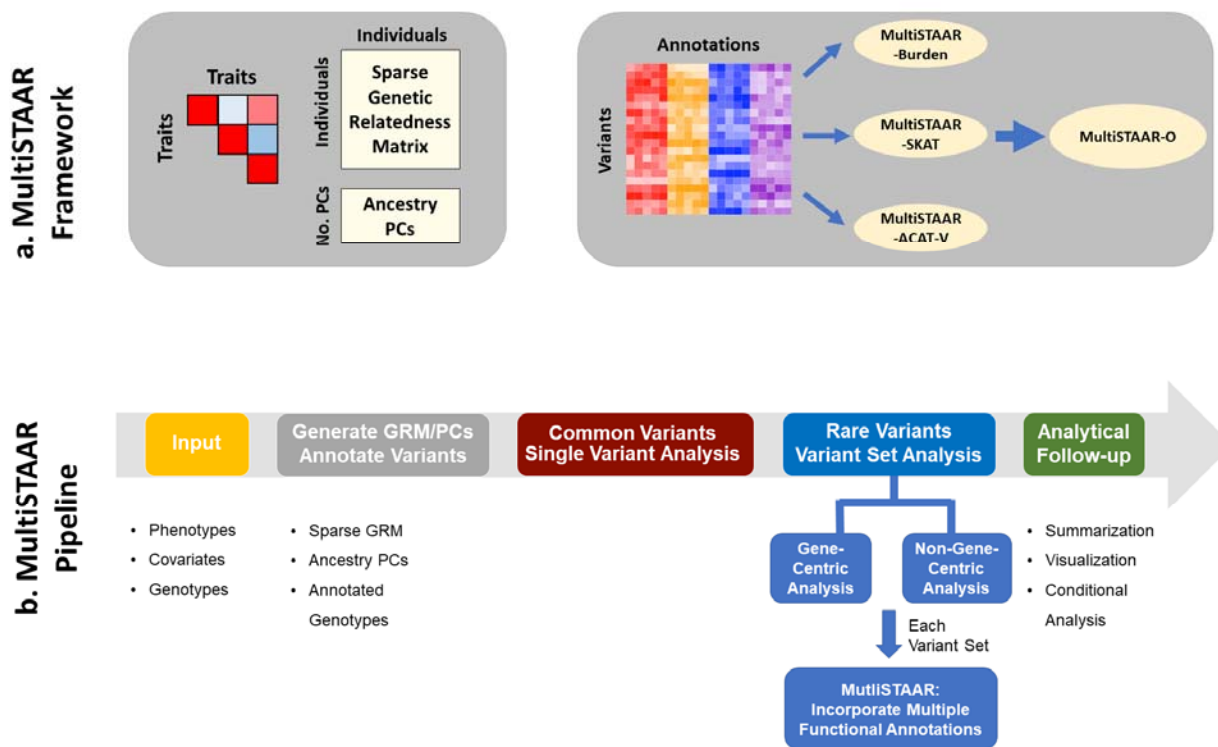
741 \* Samoan-specific missense variant.

742

743

744 **FIGURES**

745 **Fig. 1 | MultiSTAAR framework and pipeline. a**, MultiSTAAR framework. (i) Fit null  
 746 Multivariate Linear Mixed Models (MLMMs) using sparse GRM and ancestry PCs to  
 747 account for population structure, relatedness and the correlation between phenotypes.  
 748 (ii) Test for associations between each variant set and multiple traits by dynamically  
 749 incorporating multiple variant functional annotations. **b**, MultiSTAAR pipeline. (i) Prepare  
 750 the input data of MultiSTAAR, including genotypes, multiple phenotypes and covariates.  
 751 (ii) Calculate sparse GRM, ancestry PCs and annotate all variants in the genome. (iii)  
 752 Perform single variant analysis for common variants. (iv) Define the rare variant analysis  
 753 units, including gene-centric analysis of five coding functional categories and eight  
 754 noncoding functional categories and non-gene-centric analysis of sliding windows. (v)  
 755 Provide result summarization and perform analytical follow-up via conditional analysis.  
 756

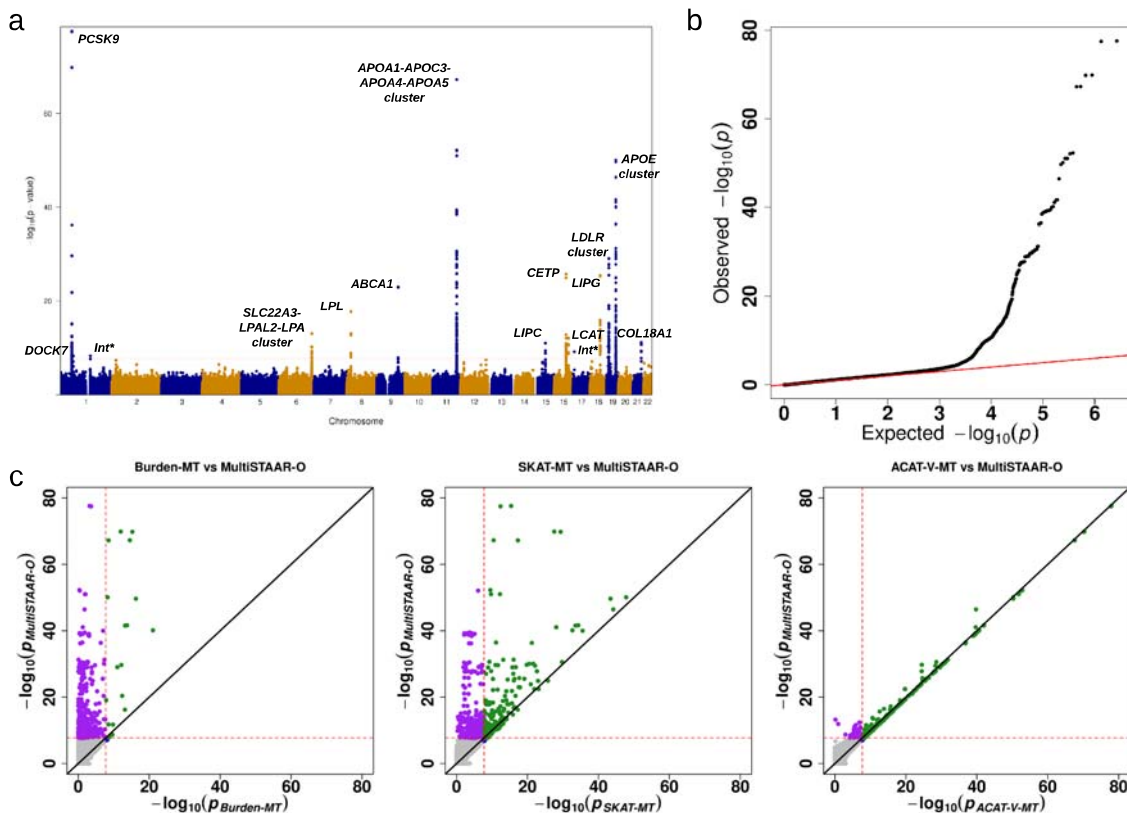


757

758

759

760 **Fig. 2 | TOPMed Genetic region (2-kb sliding window) unconditional multi-trait**  
 761 **analysis results of low-density lipoprotein cholesterol (LDL-C), high-density**  
 762 **lipoprotein cholesterol (HDL-C) and triglycerides (TG) using TOPMed data. a,**  
 763 **Manhattan plot showing the associations of 2.65 million 2-kb sliding windows versus**  
 764  **$-\log_{10}(P)$  of MultiSTAAR-O. The horizontal line indicates a genome-wide  $P$  value**  
 765 **threshold of  $1.89 \times 10^{-8}$  ( $n = 61,838$ ). b, Quantile-quantile plot of 2-kb sliding window**  
 766 **MultiSTAAR-O  $P$  values ( $n = 61,838$ ). c, Scatterplot of  $P$  values for the 2-kb sliding**  
 767 **windows comparing MultiSTAAR-O with Burden-MT, SKAT-MT and ACAT-V-MT tests**  
 768 **(MT is short for Multi-Trait). Each dot represents a sliding window with x-axis label being**  
 769 **the  $-\log_{10}(P)$  of the conventional multi-trait test and y-axis label being the  $-\log_{10}(P)$  of**  
 770 **MultiSTAAR-O ( $n = 61,838$ ). Burden-MT, SKAT-MT, ACAT-V-MT and MultiSTAAR-O**  
 771 **are two-sided tests. Int\*, intergenic sliding window.**



772

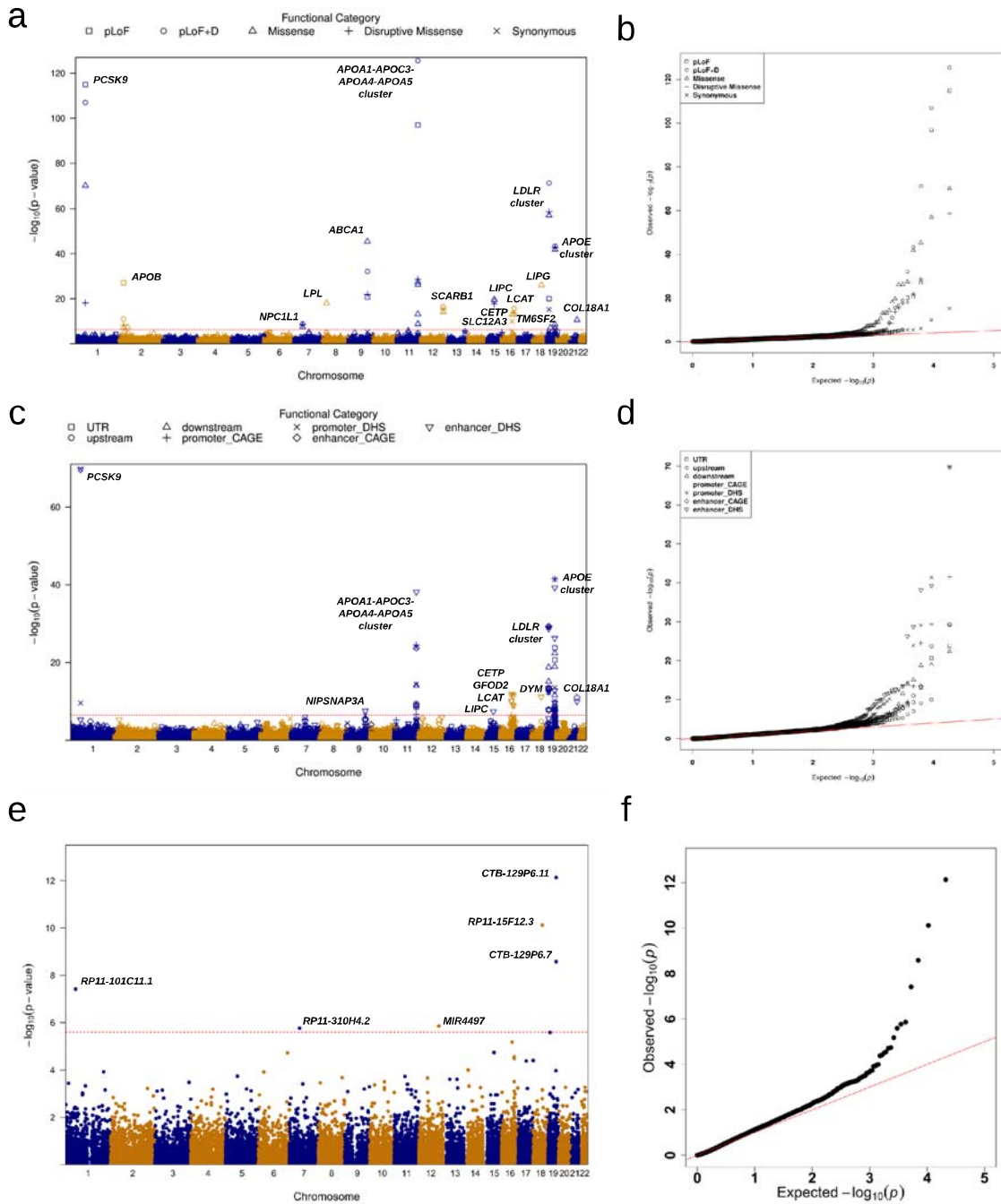
773

774

775 **EXTENDED DATA FIGURES**

776 **Extended Data Fig. 1 | Manhattan plots and Q-Q plots for unconditional**  
777 **gene-centric coding, noncoding and ncRNA analysis of low-density lipoprotein**  
778 **cholesterol (LDL-C), high-density lipoprotein cholesterol (HDL-C) and triglycerids**  
779 **(TG) using TOPMed data (n = 61,838). a,** Manhattan plots for unconditional gene-  
780 centric coding analysis of protein-coding gene. The horizontal line indicates a genome-  
781 wide MultiSTAAR-O  $P$  value threshold of  $5.00 \times 10^{-7}$ . The significant threshold is  
782 defined by multiple comparisons using the Bonferroni correction ( $0.05 / (20,000 \times 5) =$   
783  $5.00 \times 10^{-7}$ ). Different symbols represent the MultiSTAAR-O  $P$  value of the protein-  
784 coding gene using different functional categories (putative loss-of-function, putative  
785 loss-of-function and disruptive missense, missense, disruptive missense, synonymous).  
786 **b,** Quantile-quantile plots for unconditional gene-centric coding analysis of protein-  
787 coding gene. Different symbols represent the MultiSTAAR-O  $P$ -value of the gene using  
788 different functional categories. **c,** Manhattan plots for unconditional gene-centric  
789 noncoding analysis of protein-coding gene. The horizontal line indicates a genome-wide  
790 MultiSTAAR-O  $P$  value threshold of  $3.57 \times 10^{-7}$ . The significant threshold is defined by  
791 multiple comparisons using the Bonferroni correction ( $0.05 / (20,000 \times 7) = 3.57 \times 10^{-7}$ ).  
792 Different symbols represent the MultiSTAAR-O  $P$  value of the protein-coding gene using  
793 different functional categories (upstream, downstream, UTR, promoter\_CAGE,  
794 promoter\_DHS, enhancer\_CAGE, enhancer\_DHS). Promoter\_CAGE and  
795 promoter\_DHS are the promoters with overlap of Cap Analysis of Gene Expression  
796 (CAGE) sites and DNase hypersensitivity (DHS) sites for a given gene, respectively.  
797 Enhancer\_CAGE and enhancer\_DHS are the enhancers in GeneHancer predicted  
798 regions with the overlap of CAGE sites and DHS sites for a given gene, respectively. **d,**  
799 Quantile-quantile plots for unconditional gene-centric noncoding analysis of protein-  
800 coding gene. Different symbols represent the MultiSTAAR-O  $P$ -value of the gene using  
801 different functional categories. **e,** Manhattan plots for unconditional gene-centric  
802 noncoding analysis of ncRNA gene. The horizontal line indicates a genome-wide  
803 MultiSTAAR-O  $P$  value threshold of  $2.50 \times 10^{-6}$ . The significant threshold is defined by  
804 multiple comparisons using the Bonferroni correction ( $0.05 / 20,000 = 2.50 \times 10^{-6}$ ). **f,**  
805 Quantile-quantile plots for unconditional gene-centric noncoding analysis of ncRNA

806 gene. In panels, **a**, **c** and **e**, the chromosome number are indicated by the colors of  
 807 dots. In all panels, MultiSTAAR-O is a two-sided test.

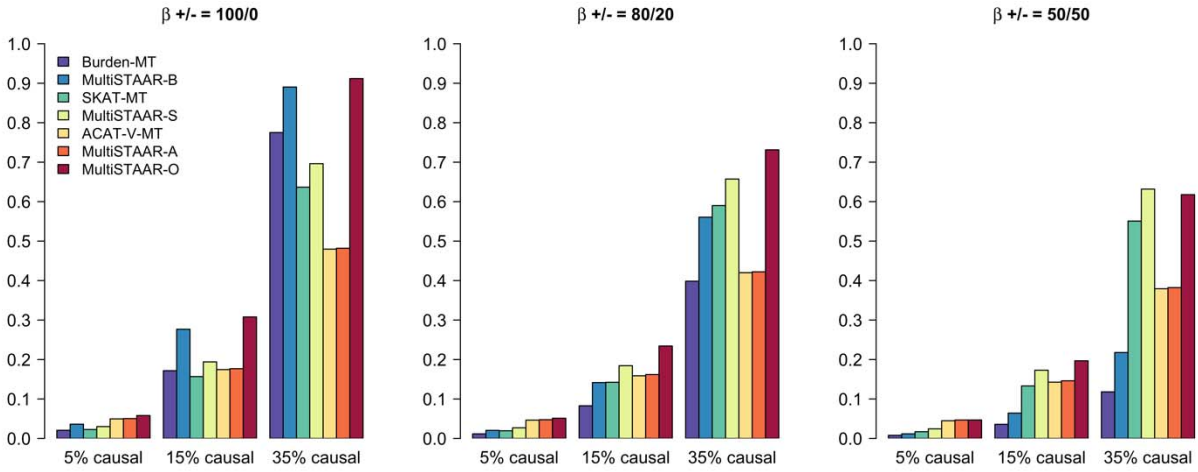


808

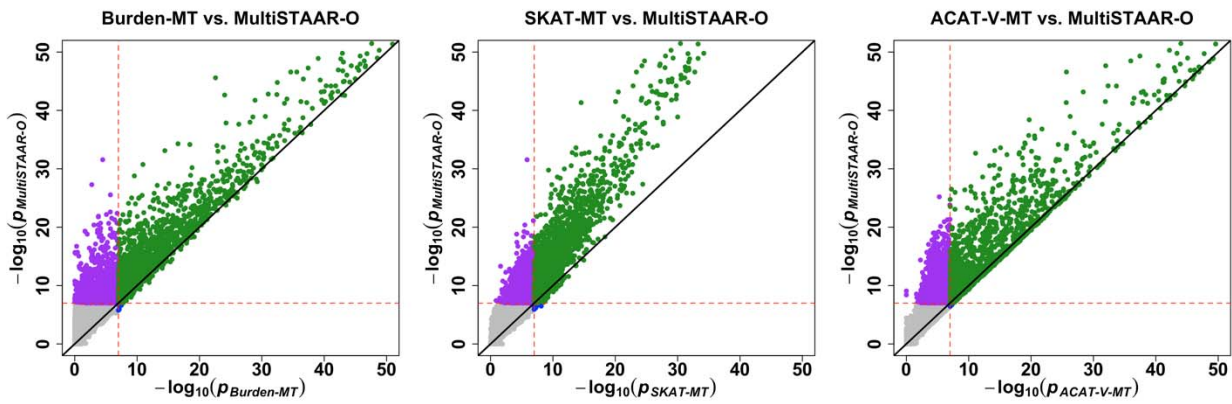
809

810

811 **Extended Data Fig. 2 | Power comparisons of Burden-MT, SKAT-MT, ACAT-V-MT**  
812 **(MT is short for Multi-Trait) and MultiSTAAR methods when variants in the signal**  
813 **region are associated with one phenotype.** Multi-trait Burden, SKAT and ACAT-V  
814 tests implemented in MultiSTAAR are denoted by Burden-MT, SKAT-MT and ACAT-V-  
815 MT. MultiSTAAR methods incorporating ten functional annotations are denoted by  
816 MultiSTAAR-B, MultiSTAAR-S, MultiSTAAR-A and MultiSTAAR-O. In each simulation  
817 replicate, a 5-kb region was randomly selected as the signal region. Within each signal  
818 region, variants were randomly generated to be causal based on the multivariate logistic  
819 model and on average there were 5%, 15% or 35% causal variants in the signal region.  
820 The effect sizes of causal variants were  $\beta_j = c_0 |\log_{10} MAF_j|$ , where  $c_0$  was set to be  
821 0.13. The barplot of power in the top panel consider settings in which the effect sizes for  
822 the causal variants are 100% positive (0% negative), 80% positive (20% negative), and  
823 50% positive (50% negative). The scatterplot of  $P$  values in the bottom panel compare  
824 MultiSTAAR-O to Burden-MT, SKAT-MT and ACAT-V-MT when 15% of variants in the  
825 signal region are causal variants with all positive effect sizes. Power was estimated as  
826 the proportion of the  $P$  values less than  $\alpha = 10^{-7}$  based on  $10^4$  replicates. Burden-MT,  
827 SKAT-MT, ACAT-V-MT, MultiSTAAR-B, MultiSTAAR-S, MultiSTAAR-A and  
828 MultiSTAAR-O are two-sided tests. Total sample size considered was 10,000.



829



830

831

832

833

834

835

836

837

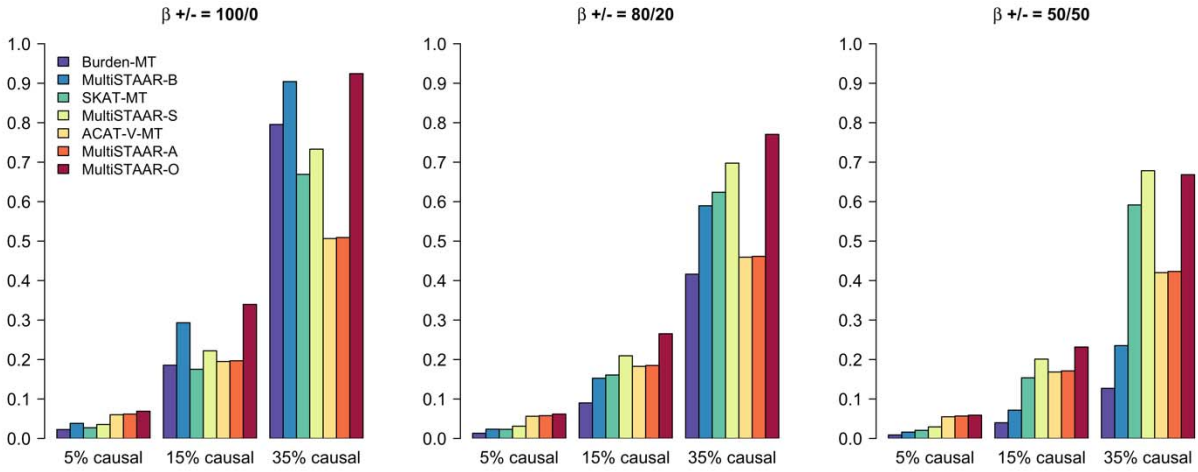
838

839

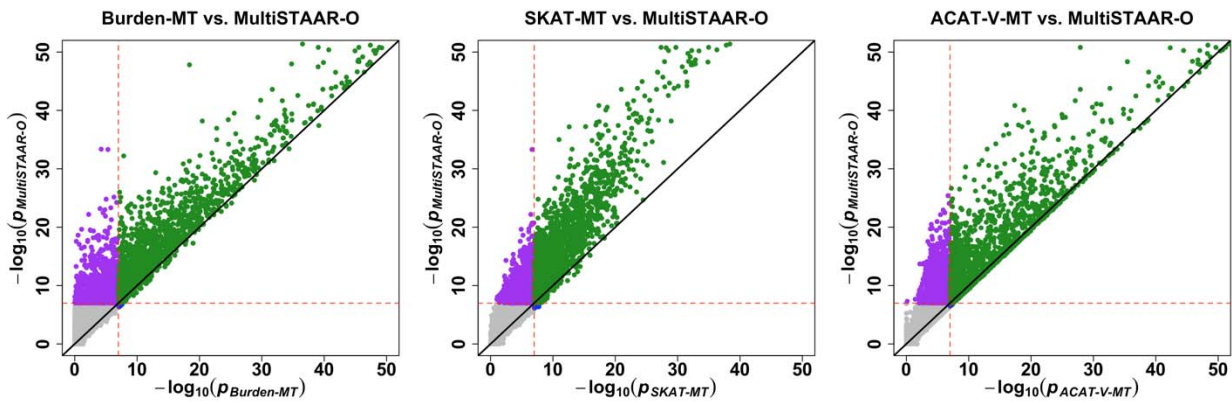
840



841 **Extended Data Fig. 3 | Power comparisons of Burden-MT, SKAT-MT, ACAT-V-MT**  
842 **(MT is short for Multi-Trait) and MultiSTAAR methods when variants in the signal**  
843 **region are associated with two positively correlated phenotypes.** In each  
844 simulation replicate, a 5-kb region was randomly selected as the signal region. Within  
845 each signal region, variants were randomly generated to be causal based on the  
846 multivariate logistic model and on average there were 5%, 15% or 35% causal variants  
847 in the signal region. The effect sizes of causal variants were  $\beta_j = c_0 |\log_{10} MAF_j|$ , where  
848  $c_0$  was set to be 0.1. The barplot of power in the top panel consider settings in which the  
849 effect sizes for the causal variants are 100% positive (0% negative), 80% positive (20%  
850 negative), and 50% positive (50% negative). The scatterplot of  $P$  values in the bottom  
851 panel compare MultiSTAAR-O to Burden-MT, SKAT-MT and ACAT-V-MT when 15% of  
852 variants in the signal region are causal variants with all positive effect sizes. Power was  
853 estimated as the proportion of the  $P$  values less than  $\alpha = 10^{-7}$  based on  $10^4$  replicates.  
854 Burden-MT, SKAT-MT, ACAT-V-MT, MultiSTAAR-B, MultiSTAAR-S, MultiSTAAR-A and  
855 MultiSTAAR-O are two-sided tests. Total sample size considered was 10,000.  
856  
857



858



859

860

861

862

863

864

865

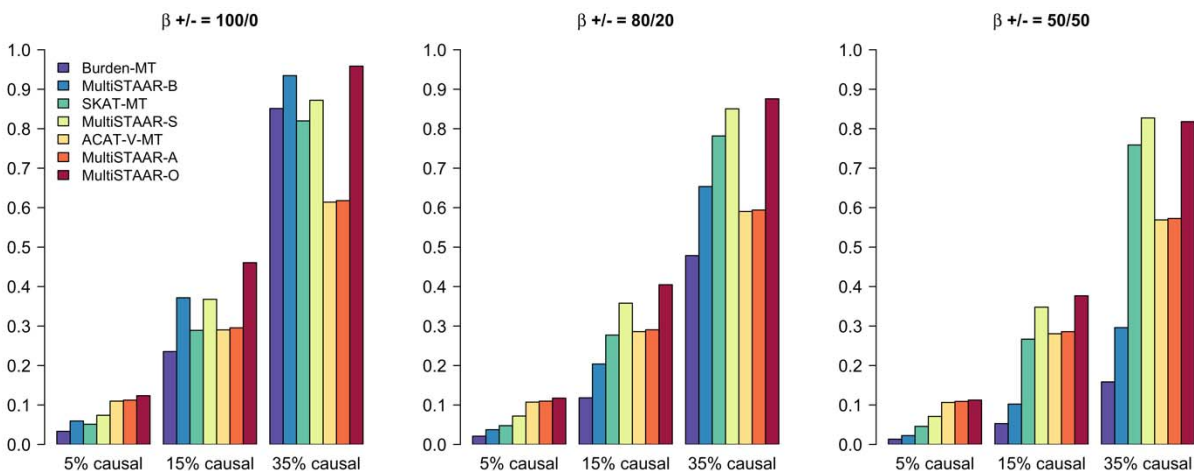
866

867

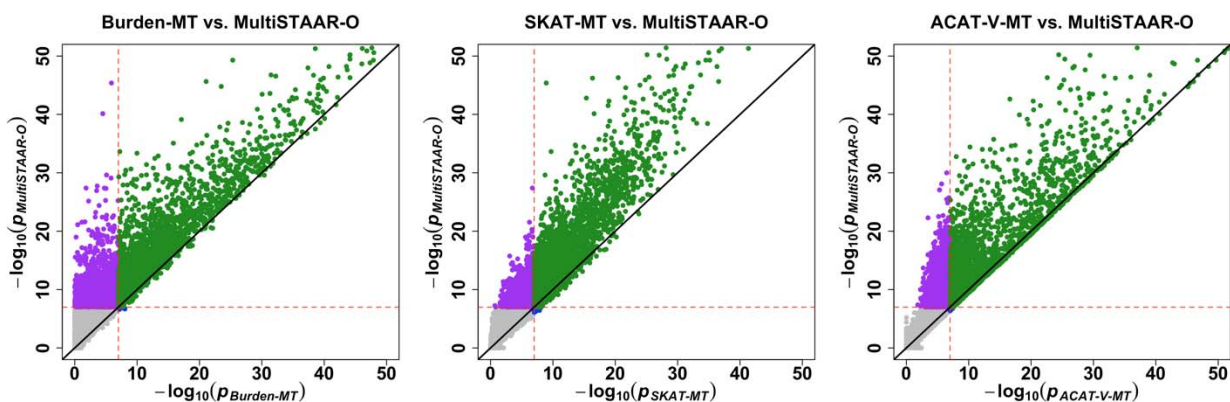
868

869

870 **Extended Data Fig. 4 | Power comparisons of Burden-MT, SKAT-MT, ACAT-V-MT**  
871 **(MT is short for Multi-Trait) and MultiSTAAR methods when variants in the signal**  
872 **region are associated with two negatively correlated phenotypes.** In each  
873 simulation replicate, a 5-kb region was randomly selected as the signal region. Within  
874 each signal region, variants were randomly generated to be causal based on the  
875 multivariate logistic model and on average there were 5%, 15% or 35% causal variants  
876 in the signal region. The effect sizes of causal variants were  $\beta_j = c_0 |\log_{10} MAF_j|$ , where  
877  $c_0$  was set to be 0.1. The barplot of power in the top panel consider settings in which the  
878 effect sizes for the causal variants are 100% positive (0% negative), 80% positive (20%  
879 negative), and 50% positive (50% negative). The scatterplot of  $P$  values in the bottom  
880 panel compare MultiSTAAR-O to Burden-MT, SKAT-MT and ACAT-V-MT when 15% of  
881 variants in the signal region are causal variants with all positive effect sizes. Power was  
882 estimated as the proportion of the  $P$  values less than  $\alpha = 10^{-7}$  based on  $10^4$  replicates.  
883 Burden-MT, SKAT-MT, ACAT-V-MT, MultiSTAAR-B, MultiSTAAR-S, MultiSTAAR-A and  
884 MultiSTAAR-O are two-sided tests. Total sample size considered was 10,000.  
885  
886



887



888

889

890

891

892

893

894

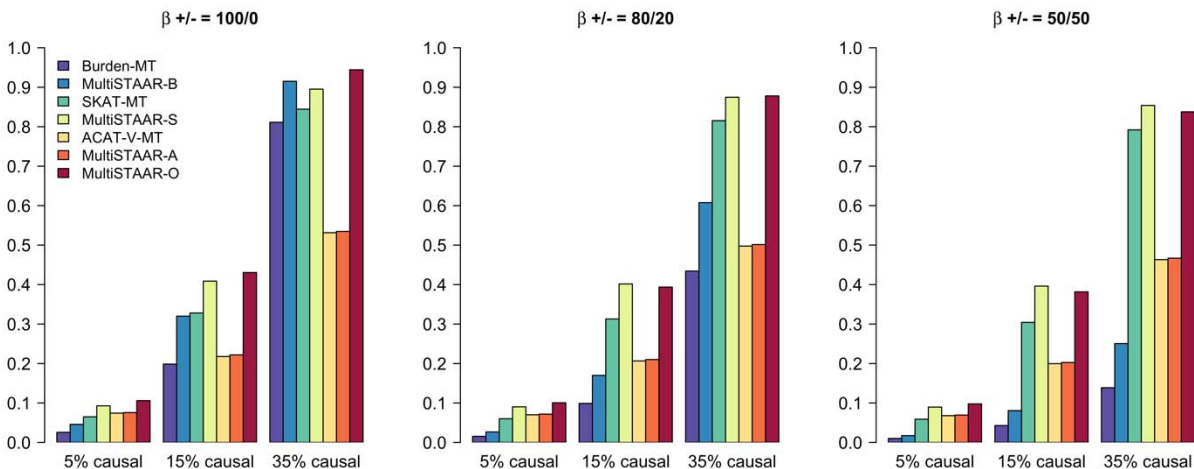
895

896

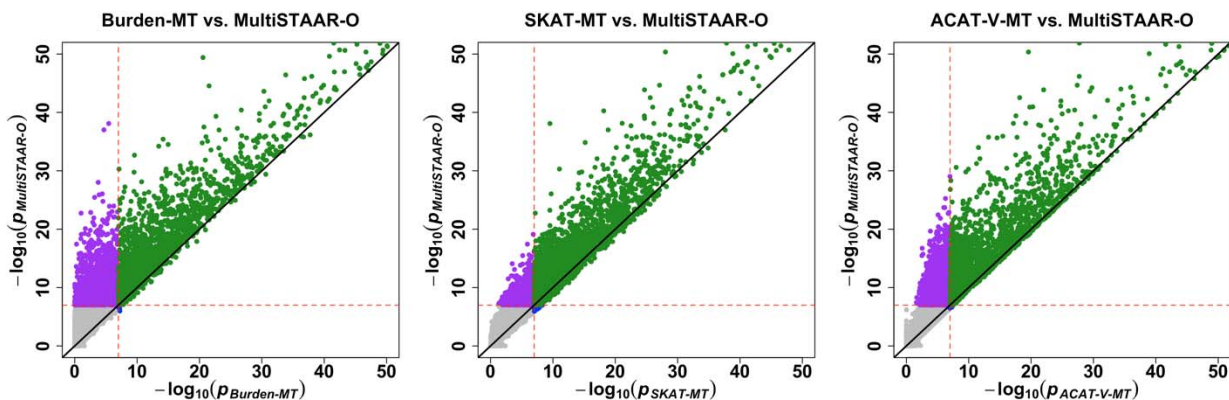
897

898

899 **Extended Data Fig. 5 | Power comparisons of Burden-MT, SKAT-MT, ACAT-V-MT**  
900 **(MT is short for Multi-Trait) and MultiSTAAR methods when variants in the signal**  
901 **region are associated with three phenotypes.** In each simulation replicate, a 5-kb  
902 region was randomly selected as the signal region. Within each signal region, variants  
903 were randomly generated to be causal based on the multivariate logistic model and on  
904 average there were 5%, 15% or 35% causal variants in the signal region. The effect  
905 sizes of causal variants were  $\beta_j = c_0 |\log_{10} MAF_j|$ , where  $c_0$  was set to be 0.07. The  
906 barplot of power in the top panel consider settings in which the effect sizes for the  
907 causal variants are 100% positive (0% negative), 80% positive (20% negative), and  
908 50% positive (50% negative). The scatterplot of  $P$  values in the bottom panel compare  
909 MultiSTAAR-O to Burden-MT, SKAT-MT and ACAT-V-MT when 15% of variants in the  
910 signal region are causal variants with all positive effect sizes. Power was estimated as  
911 the proportion of the  $P$  values less than  $\alpha = 10^{-7}$  based on  $10^4$  replicates. Burden-MT,  
912 SKAT-MT, ACAT-V-MT, MultiSTAAR-B, MultiSTAAR-S, MultiSTAAR-A and  
913 MultiSTAAR-O are two-sided tests. Total sample size considered was 10,000.  
914  
915



916



917

918

919

920

921

922

923

924

925

926

927

## 928 References

- 929 1. Taliun, D. *et al.* Sequencing of 53,831 diverse genomes from the NHLBI  
930 TOPMed Program. *Nature* **590**, 290-299 (2021).
- 931 2. The “All of Us” Research Program. *New England Journal of Medicine* **381**, 668-  
932 676 (2019).
- 933 3. Halldorsson, B.V. *et al.* The sequences of 150,119 genomes in the UK Biobank.  
934 *Nature* **607**, 732-740 (2022).
- 935 4. Lee, S., Abecasis, Gonçalo R., Boehnke, M. & Lin, X. Rare-Variant Association  
936 Analysis: Study Designs and Statistical Tests. *The American Journal of Human*  
937 *Genetics* **95**, 5-23 (2014).
- 938 5. Li, B. & Leal, S.M. Methods for Detecting Associations with Rare Variants for  
939 Common Diseases: Application to Analysis of Sequence Data. *The American*  
940 *Journal of Human Genetics* **83**, 311-321 (2008).
- 941 6. Madsen, B.E. & Browning, S.R. A Groupwise Association Test for Rare  
942 Mutations Using a Weighted Sum Statistic. *PLOS Genetics* **5**, e1000384 (2009).
- 943 7. Morris, A.P. & Zeggini, E. An evaluation of statistical approaches to rare variant  
944 analysis in genetic association studies. *Genetic Epidemiology* **34**, 188-193  
945 (2010).
- 946 8. Wu, Michael C. *et al.* Rare-Variant Association Testing for Sequencing Data with  
947 the Sequence Kernel Association Test. *The American Journal of Human*  
948 *Genetics* **89**, 82-93 (2011).
- 949 9. Liu, Y. *et al.* ACAT: A Fast and Powerful p Value Combination Method for Rare-  
950 Variant Analysis in Sequencing Studies. *The American Journal of Human*  
951 *Genetics* **104**, 410-421 (2019).
- 952 10. Solovieff, N., Cotsapas, C., Lee, P.H., Purcell, S.M. & Smoller, J.W. Pleiotropy in  
953 complex traits: challenges and strategies. *Nature Reviews Genetics* **14**, 483-495  
954 (2013).
- 955 11. Sivakumaran, S. *et al.* Abundant Pleiotropy in Human Complex Diseases and  
956 Traits. *The American Journal of Human Genetics* **89**, 607-618 (2011).
- 957 12. Abdellaoui, A., Yengo, L., Verweij, K.J.H. & Visscher, P.M. 15 years of GWAS  
958 discovery: Realizing the promise. *The American Journal of Human Genetics*  
959 (2023).
- 960 13. Watanabe, K. *et al.* A global overview of pleiotropy and genetic architecture in  
961 complex traits. *Nature Genetics* **51**, 1339-1348 (2019).
- 962 14. Wu, B. & Pankow, J.S. Sequence Kernel Association Test of Multiple Continuous  
963 Phenotypes. *Genetic Epidemiology* **40**, 91-100 (2016).
- 964 15. Dutta, D., Scott, L., Boehnke, M. & Lee, S. Multi-SKAT: General framework to  
965 test for rare-variant association with multiple phenotypes. *Genetic Epidemiology*  
966 **43**, 4-23 (2019).
- 967 16. Luo, L. *et al.* Multi-trait analysis of rare-variant association summary statistics  
968 using MTAR. *Nature Communications* **11**, 2850 (2020).
- 969 17. Broadaway, K.A. *et al.* A Statistical Approach for Testing Cross-Phenotype  
970 Effects of Rare Variants. *The American Journal of Human Genetics* **98**, 525-540  
971 (2016).

- 972 18. Li, X. *et al.* Dynamic incorporation of multiple in silico functional annotations  
973 empowers rare variant association analysis of large whole-genome sequencing  
974 studies at scale. *Nature Genetics* **52**, 969-983 (2020).
- 975 19. Sammel, M., Lin, X. & Ryan, L. Multivariate linear mixed models for multiple  
976 outcomes. *Statistics in Medicine* **18**, 2479-2492 (1999).
- 977 20. Conomos, M.P., Miller, M.B. & Thornton, T.A. Robust Inference of Population  
978 Structure for Ancestry Prediction and Correction of Stratification in the Presence  
979 of Relatedness. *Genetic Epidemiology* **39**, 276-293 (2015).
- 980 21. Conomos, Matthew P., Reiner, Alexander P., Weir, Bruce S. & Thornton,  
981 Timothy A. Model-free Estimation of Recent Genetic Relatedness. *The American*  
982 *Journal of Human Genetics* **98**, 127-148 (2016).
- 983 22. Gogarten, S.M. *et al.* Genetic association testing using the GENESIS  
984 R/Bioconductor package. *Bioinformatics* **35**, 5346-5348 (2019).
- 985 23. Lee, P.H. *et al.* Principles and methods of in-silico prioritization of non-coding  
986 regulatory variants. *Human Genetics* **137**, 15-30 (2018).
- 987 24. Li, Z. *et al.* A framework for detecting noncoding rare-variant associations of  
988 large-scale whole-genome sequencing studies. *Nature Methods* **19**, 1599-1611  
989 (2022).
- 990 25. Morrison, A.C. *et al.* Practical approaches for whole-genome sequence analysis  
991 of heart-and blood-related traits. *The American Journal of Human Genetics* **100**,  
992 205-215 (2017).
- 993 26. Selvaraj, M.S. *et al.* Whole genome sequence analysis of blood lipid levels in  
994 >66,000 individuals. *Nature Communications* **13**, 5995 (2022).
- 995 27. Liu, Z. & Lin, X. Multiple phenotype association tests using summary statistics in  
996 genome-wide association studies. *Biometrics* **74**, 165-175 (2018).
- 997 28. Teslovich, T.M. *et al.* Biological, clinical and population relevance of 95 loci for  
998 blood lipids. *Nature* **466**, 707 (2010).
- 999 29. Schaffner, S.F. *et al.* Calibrating a coalescent simulation of human genome  
1000 sequence variation. *Genome research* **15**, 1576-1583 (2005).
- 1001 30. Natarajan, P. *et al.* Deep-coverage whole genome sequences and blood lipids  
1002 among 16,324 individuals. *Nature Communications* **9**, 3391 (2018).
- 1003 31. Stilp, A.M. *et al.* A System for Phenotype Harmonization in the National Heart,  
1004 Lung, and Blood Institute Trans-Omics for Precision Medicine (TOPMed)  
1005 Program. *American Journal of Epidemiology* (2021).
- 1006 32. Frankish, A. *et al.* GENCODE reference annotation for the human and mouse  
1007 genomes. *Nucleic acids research* **47**, D766-D773 (2019).
- 1008 33. Dong, C. *et al.* Comparison and integration of deleteriousness prediction  
1009 methods for nonsynonymous SNVs in whole exome sequencing studies. *Human*  
1010 *Molecular Genetics* **24**, 2125-2137 (2014).
- 1011 34. Li, Z. *et al.* A framework for detecting noncoding rare variant associations of  
1012 large-scale whole-genome sequencing studies. *bioRxiv*, 2021.11.05.467531  
1013 (2021).
- 1014 35. Kircher, M. *et al.* A general framework for estimating the relative pathogenicity of  
1015 human genetic variants. *Nature Genetics* **46**, 310-315 (2014).



- 1016 36. Huang, Y.-F., Gulko, B. & Siepel, A. Fast, scalable prediction of deleterious  
1017 noncoding variants from functional and population genomic data. *Nature*  
1018 *Genetics* **49**, 618-624 (2017).
- 1019 37. Rogers, M.F. *et al.* FATHMM-XF: accurate prediction of pathogenic point  
1020 mutations via extended features. *Bioinformatics* **34**, 511-513 (2017).
- 1021 38. Buniello, A. *et al.* The NHGRI-EBI GWAS Catalog of published genome-wide  
1022 association studies, targeted arrays and summary statistics 2019. *Nucleic Acids*  
1023 *Research* **47**, D1005-D1012 (2019).
- 1024 39. Klarin, D. *et al.* Genetics of blood lipids among ~300,000 multi-ethnic participants  
1025 of the Million Veteran Program. *Nature Genetics* **50**, 1514-1523 (2018).
- 1026 40. Forrest, A.R. *et al.* A promoter-level mammalian expression atlas. *Nature* **507**,  
1027 462 (2014).
- 1028 41. Abascal, F. *et al.* Expanded encyclopaedias of DNA elements in the human and  
1029 mouse genomes. *Nature* **583**, 699-710 (2020).
- 1030 42. Andersson, R. *et al.* An atlas of active enhancers across human cell types and  
1031 tissues. *Nature* **507**, 455-461 (2014).
- 1032 43. Fishilevich, S. *et al.* GeneHancer: genome-wide integration of enhancers and  
1033 target genes in GeneCards. *Database* **2017**(2017).
- 1034 44. Li, Z. *et al.* Dynamic Scan Procedure for Detecting Rare-Variant Association  
1035 Regions in Whole-Genome Sequencing Studies. *The American Journal of*  
1036 *Human Genetics* **104**, 802-814 (2019).
- 1037 45. McCaw, Z.R., Gao, J., Lin, X. & Gronsbell, J. Leveraging a machine learning  
1038 derived surrogate phenotype to improve power for genome-wide association  
1039 studies of partially missing phenotypes in population biobanks. *bioRxiv*,  
1040 2022.12.12.520180 (2022).
- 1041 46. Li, X. *et al.* Powerful, scalable and resource-efficient meta-analysis of rare variant  
1042 associations in large whole genome sequencing studies. *Nature Genetics* **55**,  
1043 154-164 (2023).
- 1044 47. Chen, H. *et al.* Control for Population Structure and Relatedness for Binary Traits  
1045 in Genetic Association Studies via Logistic Mixed Models. *The American Journal*  
1046 *of Human Genetics* **98**, 653-666 (2016).
- 1047 48. Chen, H. *et al.* Efficient Variant Set Mixed Model Association Tests for  
1048 Continuous and Binary Traits in Large-Scale Whole-Genome Sequencing  
1049 Studies. *The American Journal of Human Genetics* **104**, 260-274 (2019).
- 1050 49. Liu, Y. & Xie, J. Cauchy combination test: a powerful test with analytic p-value  
1051 calculation under arbitrary dependency structures. *Journal of the American*  
1052 *Statistical Association* **115**, 393-402 (2020).

1055

1056

1057

1058 **Methods**

1059 **Ethics statement**

1060 This study relied on analyses of genetic data from TOPMed cohorts. The study has  
1061 been approved by the TOPMed Publications Committee, TOPMed Lipids Working  
1062 Group and all the participating cohorts, including Old Order Amish (phs000956.v1.p1),  
1063 Atherosclerosis Risk in Communities Study (phs001211), Mt Sinai BioMe Biobank  
1064 (phs001644), Coronary Artery Risk Development in Young Adults (phs001612),  
1065 Cleveland Family Study (phs000954), Cardiovascular Health Study (phs001368),  
1066 Diabetes Heart Study (phs001412), Framingham Heart Study (phs000974), Genetic  
1067 Study of Atherosclerosis Risk (phs001218), Genetic Epidemiology Network of  
1068 Arteriopathy (phs001345), Genetic Epidemiology Network of Salt Sensitivity  
1069 (phs001217), Genetics of Lipid Lowering Drugs and Diet Network (phs001359),  
1070 Hispanic Community Health Study - Study of Latinos (phs001395), Hypertension  
1071 Genetic Epidemiology Network and Genetic Epidemiology Network  
1072 of Arteriopathy (phs001293), Jackson Heart Study (phs000964), Multi-Ethnic Study of  
1073 Atherosclerosis (phs001416), San Antonio Family Heart Study (phs001215), Genome-  
1074 wide Association Study of Adiposity in Samoans (phs000972), Taiwan Study of  
1075 Hypertension using Rare Variants (phs001387), and Women's Health Initiative  
1076 (phs001237), where the accession numbers are provided in parenthesis. The use of  
1077 human genetics data from TOPMed cohorts was approved by the Harvard T.H. Chan  
1078 School of Public Health IRB (IRB13-0353).

1079

1080 **Notation and model**

1081 Suppose there are  $n$  subjects with a total of  $M$  variants sequenced across the whole  
 1082 genome. For the  $i$ -th subject, let  $\mathbf{Y}_i = (y_{i1}, y_{i2}, \dots, y_{iK})^T$  denote a vector of  $K$  quantitative  
 1083 phenotypes;  $\mathbf{X}_i = (x_{i1}, x_{i2}, \dots, x_{iq})^T$  denotes  $q$  covariates, such as age, gender and  
 1084 ancestral principal components;  $\mathbf{G}_i = (G_{i1}, G_{i2}, \dots, G_{ip})^T$  denotes the genotype matrix of  
 1085 the  $p$  genetic variants in a variant set. Since these  $K$  phenotypes may be defined on  
 1086 different measurement scales, we assume that each phenotype has been rescaled to  
 1087 have zero mean and unit variance.

1088

1089 When the data consist of unrelated samples, we consider the following Multivariate  
 1090 Linear Model (MLM)

$$\mathbf{Y}_i = \begin{bmatrix} y_{i1} \\ y_{i2} \\ \vdots \\ y_{iK} \end{bmatrix} = \begin{bmatrix} \alpha_{0,1} + \mathbf{X}_i^T \boldsymbol{\alpha}_1 + \mathbf{G}_i^T \boldsymbol{\beta}_1 \\ \alpha_{0,2} + \mathbf{X}_i^T \boldsymbol{\alpha}_2 + \mathbf{G}_i^T \boldsymbol{\beta}_2 \\ \vdots \\ \alpha_{0,K} + \mathbf{X}_i^T \boldsymbol{\alpha}_K + \mathbf{G}_i^T \boldsymbol{\beta}_K \end{bmatrix} + \begin{bmatrix} \varepsilon_{i1} \\ \varepsilon_{i2} \\ \vdots \\ \varepsilon_{iK} \end{bmatrix}, \#(1)$$

1091 where  $\alpha_{0,k}$  is an intercept,  $\boldsymbol{\alpha}_k = (\alpha_{1,k}, \alpha_{2,k}, \dots, \alpha_{q,k})^T$  and  $\boldsymbol{\beta}_k = (\beta_{1,k}, \beta_{2,k}, \dots, \beta_{p,k})^T$  are  
 1092 column vectors of regression coefficients for covariates  $\mathbf{X}_i$  and genotype  $\mathbf{G}_i$  in  
 1093 phenotype  $k$ , respectively. The error terms  $\boldsymbol{\varepsilon}_i = (\varepsilon_{i1}, \varepsilon_{i2}, \dots, \varepsilon_{iK})^T$  are independent and  
 1094 identically distributed and follow a multivariate normal distribution with mean a vector of  
 1095 zeros and variance-covariance matrix  $\boldsymbol{\Sigma}_{K \times K}$ , assumed identical for all subjects. For all  $n$   
 1096 subjects, using matrix notation we can write model (1) as

$$\mathbf{Y}_{n \times K} = \mathbf{1}_n \boldsymbol{\alpha}_0^T + \mathbf{X}_{n \times q} \boldsymbol{\alpha}_{q \times K} + \mathbf{G}_{n \times p} \boldsymbol{\beta}_{p \times K} + \boldsymbol{\varepsilon}_{n \times K}, \#(2)$$

1097 where  $\mathbf{1}_n$  is a column vector of 1's with length  $n$ ,  $\boldsymbol{\alpha}_0 = (\alpha_{0,1}, \alpha_{0,2}, \dots, \alpha_{0,K})^T$  is a column  
 1098 vector of regression intercepts, the  $k$ -th columns of  $\boldsymbol{\alpha}_{q \times K}$  and  $\boldsymbol{\beta}_{p \times K}$  are  $\boldsymbol{\alpha}_k$  and  $\boldsymbol{\beta}_k$ ,  
 1099 respectively, and  $\boldsymbol{\varepsilon}_{n \times K} = (\boldsymbol{\varepsilon}_1, \boldsymbol{\varepsilon}_2, \dots, \boldsymbol{\varepsilon}_n)^T \sim \text{MatrixNormal}_{n,K}(\mathbf{0}_{n \times K}, \mathbf{I}_{n \times n}, \boldsymbol{\Sigma}_{K \times K})$  follows a

1100 matrix normal distribution. We calculate the scaled residual for each subject on each  
 1101 phenotype, defined as  $\hat{\mathbf{e}}_{n \times K} = (\mathbf{Y}_{n \times K} - \hat{\boldsymbol{\mu}}_{n \times K}) \hat{\boldsymbol{\Sigma}}_{K \times K}^{-1}$ , where  $\hat{\boldsymbol{\mu}}_{n \times K}$  (a matrix of fitted  
 1102 values) and  $\hat{\boldsymbol{\Sigma}}_{K \times K}$  are estimated under the null MLM  $\mathbf{Y}_{n \times K} = \mathbf{1}_n \boldsymbol{\alpha}_0^T + \mathbf{X}_{n \times q} \boldsymbol{\alpha}_{q \times K} + \boldsymbol{\varepsilon}_{n \times K}$ ,  
 1103 where no variant has any effect on any outcome.

1104

1105 When the data consist of related samples, we consider the following Multivariate Linear  
 1106 Mixed Model (MLMM)<sup>19,47,48</sup>

$$\mathbf{Y}_i = \begin{bmatrix} y_{i1} \\ y_{i2} \\ \vdots \\ y_{iK} \end{bmatrix} = \begin{bmatrix} \alpha_{0,1} + \mathbf{X}_i^T \boldsymbol{\alpha}_1 + \mathbf{G}_i^T \boldsymbol{\beta}_1 \\ \alpha_{0,2} + \mathbf{X}_i^T \boldsymbol{\alpha}_2 + \mathbf{G}_i^T \boldsymbol{\beta}_2 \\ \vdots \\ \alpha_{0,K} + \mathbf{X}_i^T \boldsymbol{\alpha}_K + \mathbf{G}_i^T \boldsymbol{\beta}_K \end{bmatrix} + \begin{bmatrix} b_{i1} \\ b_{i2} \\ \vdots \\ b_{iK} \end{bmatrix} + \begin{bmatrix} \varepsilon_{i1} \\ \varepsilon_{i2} \\ \vdots \\ \varepsilon_{iK} \end{bmatrix}, \quad \#(3)$$

1107 where the random effects  $b_{ik}$  account for relatedness and remaining population  
 1108 structure unaccounted by ancestral PCs<sup>20</sup>. We assume that  $\mathbf{b}_{n \times K} = (b_{ik})_{n \times K} \sim$   
 1109  $\text{MatrixNormal}_{n,K}(\mathbf{0}_{n \times K}, \boldsymbol{\Phi}_{n \times n}, \boldsymbol{\Theta}_{K \times K})$  with a variance component matrix  $\boldsymbol{\Theta}_{K \times K}$  and a  
 1110 sparse genetic relatedness matrix  $\boldsymbol{\Phi}_{n \times n}$ <sup>21,22</sup>. For all  $n$  subjects, using matrix notation we  
 1111 can rewrite equation (3) as

$$\mathbf{Y}_{n \times K} = \mathbf{1}_n \boldsymbol{\alpha}_0^T + \mathbf{X}_{n \times q} \boldsymbol{\alpha}_{q \times K} + \mathbf{G}_{n \times p} \boldsymbol{\beta}_{p \times K} + \mathbf{b}_{n \times K} + \boldsymbol{\varepsilon}_{n \times K}. \#(4)$$

1112 We calculate the scaled residual for each subject on each phenotype, defined as  
 1113  $\hat{\mathbf{e}}_{n \times K} = (\mathbf{Y}_{n \times K} - \hat{\boldsymbol{\mu}}_{n \times K}) \hat{\boldsymbol{\Sigma}}_{K \times K}^{-1}$ , where  $\hat{\boldsymbol{\mu}}_{n \times K}$  and  $\hat{\boldsymbol{\Sigma}}_{K \times K}$  are estimated under the null MLMM  
 1114  $\mathbf{Y}_{n \times K} = \mathbf{1}_n \boldsymbol{\alpha}_0^T + \mathbf{X}_{n \times q} \boldsymbol{\alpha}_{q \times K} + \mathbf{b}_{n \times K} + \boldsymbol{\varepsilon}_{n \times K}$ . Under both MLM and MLMM, our goal is to  
 1115 test for an association between a set of  $p$  genetic variants and  $K$  quantitative  
 1116 phenotypes, adjusting for covariates and relatedness. This corresponds to testing  
 1117  $H_0: \boldsymbol{\beta}_1 = \boldsymbol{\beta}_2 = \dots \boldsymbol{\beta}_K = \mathbf{0}$ .

1118

## 1119 **Multi-trait rare variant association tests using MultiSTAAR**

1120 Single-trait score-based aggregation methods<sup>5-9</sup> can be extended to allow for jointly  
1121 testing the association between rare variants in a variant set and multiple quantitative  
1122 phenotypes. For a given variant set, let  $\mathbf{S}_{p \times K} = (\mathbf{S}_{jk})_{p \times K} = (\mathbf{G}_{n \times p})^T \hat{\mathbf{e}}_{n \times K}$  denote the  
1123 matrix of score statistics where  $S_{jk}$  is the score statistic for the  $j$ -th variant on the  $k$ -th  
1124 phenotype. For multi-trait burden test using MultiSTAAR (Burden-MT), we consider test  
1125 statistic

$$Q_{Burden-MT} = \left( \sum_{j=1}^p w_j \mathbf{S}_{j\cdot} \right) \hat{\mathbf{V}}^{-1} \left( \sum_{j=1}^p w_j \mathbf{S}_{j\cdot} \right)^T,$$

1126 where  $w_j$  is the weight defined as a function of the MAF for the  $j$ -th variant<sup>4,18</sup>,  $\mathbf{S}_{j\cdot} =$   
1127  $(S_{j1}, S_{j2}, \dots, S_{jK})$  is the  $j$ -th row of  $\mathbf{S}$  and  $\hat{\mathbf{V}}$  is the estimated variance-covariance matrix of  
1128  $\sum_{j=1}^p w_j \mathbf{S}_{j\cdot} = \mathbf{w}^T \mathbf{S}$ .  $Q_{Burden-MT}$  asymptotically follows a standard chi-square distribution  
1129 with  $K$  degrees of freedom under the null hypothesis, and its  $P$  value can be obtained  
1130 analytically while accounting for LD between variants and correlation between  
1131 phenotypes.

1132

1133 For multi-trait SKAT using MultiSTAAR (SKAT-MT), we consider the statistic

$$Q_{SKAT-MT} = \sum_{k=1}^K \sum_{j=1}^p w_j^2 S_{jk}^2.$$

1134  $Q_{SKAT-MT}$  asymptotically follows a mixture of chi-square distributions under the null  
1135 hypothesis, and its  $P$  value can be obtained analytically while accounting for LD  
1136 between variants and correlation between phenotypes<sup>14,15</sup>.

1137

1138 For multi-trait ACAT-V using MultiSTAAR (ACAT-V-MT), we propose test statistic

$$Q_{ACAT-V-MT} = \overline{w^2 \text{MAF}(1 - \text{MAF})} \tan((0.5 - p_0)\pi) + \sum_{j=1}^{p'} w_j^2 \text{MAF}_j(1 - \text{MAF}_j) \tan((0.5 - p_j)\pi),$$

1139 where  $p'$  is the number of variants with a minor allele count (MAC) greater than 10 and

1140  $p_j$  is the multi-trait association  $P$  value of individual variant  $j$  for those variants with a

1141  $\text{MAC} > 10$ , whose test statistic is given by the  $K$  degrees of freedom multivariate score

1142 test

$$Q_j = \mathbf{S}_j \cdot \widehat{\mathbf{V}}_{\mathbf{S}_j}^{-1} \mathbf{S}_j^T,$$

1143 where  $\widehat{\mathbf{V}}_{\mathbf{S}_j}$  is the estimated variance-covariance matrix of  $\mathbf{S}_j$ ;  $p_0$  is the multi-trait burden

1144 test  $P$  value of extremely rare variants with an  $\text{MAC} \leq 10$  as described above and

1145  $\overline{w^2 \text{MAF}(1 - \text{MAF})}$  is the average of the weights  $w_j^2 \text{MAF}_j(1 - \text{MAF}_j)$  among the

1146 extremely rare variants with an  $\text{MAC} \leq 10$ .  $Q_{ACAT-V-MT}$  is approximated well by a scaled

1147 Cauchy distribution under the null hypothesis, and its  $P$  value can be obtained

1148 analytically while accounting for LD between variants and correlation between

1149 phenotypes<sup>9,49</sup>. Note that when  $K = 1$ , the multi-trait burden, SKAT, and ACAT-V tests

1150 reduce to the original single-trait burden, SKAT and ACAT-V tests.

1151

1152 Suppose we have a collection of  $L$  annotations, let  $A_{jl}$  denote the  $l$ -th annotation for the

1153  $j$ th variant in the variant set. We define the functionally-informed multi-trait burden,

1154 SKAT and ACAT-V test statistics weighted by the  $l$ -th annotation as follows

$$Q_{Burden-MT,l,(a_1,a_2)} = \left( \sum_{j=1}^p \hat{\pi}_{jl} w_{j,(a_1,a_2)} \mathbf{S}_j \right) \hat{\mathbf{V}}_{l,(a_1,a_2)}^{-1} \left( \sum_{j=1}^p \hat{\pi}_{jl} w_{j,(a_1,a_2)} \mathbf{S}_j \right)^T,$$

$$Q_{SKAT-MT,l,(a_1,a_2)} = \sum_{k=1}^K \sum_{j=1}^p \hat{\pi}_{jl} w_{j,(a_1,a_2)}^2 S_{jk}^2,$$

$$\begin{aligned} Q_{ACAT-V-MT,l,(a_1,a_2)} &= \overline{\hat{\pi}_{\cdot l} w_{(a_1,a_2)}^2} \text{MAF}(1 - \text{MAF}) \tan\left((0.5 - p_{0,l})\pi\right) \\ &+ \sum_{j=1}^{M'} \hat{\pi}_{jl} w_{j,(a_1,a_2)}^2 \text{MAF}_j(1 - \text{MAF}_j) \tan\left((0.5 - p_j)\pi\right), \end{aligned}$$

1155 where  $\hat{\pi}_{jl} = \frac{\text{rank}(A_{jl})}{M}$ ,  $w_{j,(a_1,a_2)} = \text{Beta}(\text{MAF}_j; a_1, a_2)$  with  $(a_1, a_2) \in \mathcal{A} = \{(1,25), (1,1)\}$ ,

1156  $\hat{\mathbf{V}}_{l,(a_1,a_2)}$  is the estimated variance-covariance matrix of  $\sum_{j=1}^p \hat{\pi}_{jl} w_{j,(a_1,a_2)} \mathbf{S}_j$  and

1157  $\overline{\hat{\pi}_{\cdot l} w_{(a_1,a_2)}^2} \text{MAF}(1 - \text{MAF})$  is the average of the weights  $\hat{\pi}_{jl} w_{j,(a_1,a_2)}^2 \text{MAF}_j(1 -$

1158  $\text{MAF}_j)$  among the extremely rare variants with  $\text{MAC} \leq 10$ . Finally, we define the

1159 omnibus MultiSTAAR-O test statistic as

$$\begin{aligned} T_{MultiSTAAR-O} &= \frac{1}{3|\mathcal{A}|} \sum_{(a_1,a_2) \in \mathcal{A}} [T_{MultiSTAAR-B(a_1,a_2)} + T_{MultiSTAAR-S(a_1,a_2)} \\ &+ T_{MultiSTAAR-A(a_1,a_2)}] \\ &= \frac{1}{3|\mathcal{A}|} \sum_{(a_1,a_2) \in \mathcal{A}} \sum_{l=0}^L \left[ \frac{\tan\{(0.5 - p_{Burden-MT,l,(a_1,a_2)})\pi\}}{L+1} \right. \\ &+ \left. \frac{\tan\{(0.5 - p_{SKAT-MT,l,(a_1,a_2)})\pi\}}{L+1} + \frac{\tan\{(0.5 - p_{ACAT-V-MT,l,(a_1,a_2)})\pi\}}{L+1} \right], \end{aligned}$$

1160 and the  $P$  value of  $T_{MultiSTAAR-O}$  can be calculated by

$$p_{MultiSTAAR-O} = \frac{1}{2} - \frac{\{\arctan(T_{MultiSTAAR-O})\}}{\pi}.$$

1161

## 1162 **Data simulation**

### 1163 *Type I error rate simulations*

1164 We performed simulation studies to evaluate how accurately MultiSTAAR controls the  
1165 type I error rate. We generated three quantitative traits from a multivariate linear model,  
1166 conditional on two covariates

$$Y_i = \begin{bmatrix} Y_{i1} \\ Y_{i2} \\ Y_{i3} \end{bmatrix} = \begin{bmatrix} 0.5X_{i1} + 0.5X_{i2} \\ 0.5X_{i1} + 0.5X_{i2} \\ 0.5X_{i1} + 0.5X_{i2} \end{bmatrix} + \begin{bmatrix} \varepsilon_{i1} \\ \varepsilon_{i2} \\ \varepsilon_{i3} \end{bmatrix},$$

1167 where  $X_{i1} \sim N(0,1)$ ,  $X_{i2} \sim \text{Bernoulli}(0.5)$  and

$$\begin{bmatrix} \varepsilon_{i1} \\ \varepsilon_{i2} \\ \varepsilon_{i3} \end{bmatrix} \sim MVN \left( \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1.0 & -0.1 & 0.2 \\ -0.1 & 1.0 & -0.4 \\ 0.2 & -0.4 & 1.0 \end{bmatrix} \right).$$

1168

1169 The correlation matrix of error terms  $\varepsilon_i = (\varepsilon_{i1}, \varepsilon_{i2}, \varepsilon_{i3})^T$  was chosen to mimic the  
1170 correlations between three lipid traits LDL-C, HDL-C and TG, estimated from the  
1171 TOPMed data<sup>26</sup>. We considered a sample size of 10,000 and generated genotypes by  
1172 simulating 20,000 sequences for 100 different regions each spanning 1 Mb. The data  
1173 generation used the calibration coalescent model (COSI)<sup>29</sup> with parameters set to mimic  
1174 the LD structure of African Americans. In each simulation replicate, 10 annotations were  
1175 generated as  $A_1, \dots, A_{10}$  all independently and identically distributed as  $N(0,1)$  for each  
1176 variant, and we randomly selected 5-kb regions from these 1-Mb regions for type I error  
1177 rate simulations. We applied MultiSTAAR-B, MultiSTAAR-S, MultiSTAAR-A and  
1178 MultiSTAAR-O by incorporating MAFs and the 10 annotations together with Burden-MT,



1179 SKAT-MT and ACAT-V-MT tests. We repeated the procedure with  $10^8$  replicates to  
1180 examine the type I error rate at levels  $\alpha = 10^{-4}, 10^{-5}$ , and  $10^{-6}$ .

1181

### 1182 *Empirical power simulations*

1183 Next, we carried out simulation studies under a variety of configurations to assess the  
1184 the power of MultiSTAAR-O, and how its incorporation of multiple functional annotations  
1185 affects power compared to the multi-trait burden, SKAT, and ACAT-V tests implemented  
1186 in MultiSTAAR. In each simulation replicate, we randomly selected 5-kb regions from a  
1187 1-Mb region for power evaluations. For each selected 5-kb region, we generated three  
1188 quantitative traits from a multivariate linear model

$$\mathbf{Y}_i = \begin{bmatrix} Y_{i1} \\ Y_{i2} \\ Y_{i3} \end{bmatrix} = \begin{bmatrix} 0.5X_{i1} + 0.5X_{i2} + \mathbf{G}_i^T \boldsymbol{\beta}_1 \\ 0.5X_{i1} + 0.5X_{i2} + \mathbf{G}_i^T \boldsymbol{\beta}_2 \\ 0.5X_{i1} + 0.5X_{i2} + \mathbf{G}_i^T \boldsymbol{\beta}_3 \end{bmatrix} + \begin{bmatrix} \varepsilon_{i1} \\ \varepsilon_{i2} \\ \varepsilon_{i3} \end{bmatrix},$$

1189 where  $X_{1i}, X_{2i}, \varepsilon_i$  were defined as in the type I error rate simulations,

1190  $\mathbf{G}_i = (G_{i1}, G_{i2}, \dots, G_{ip})^T$  and  $\boldsymbol{\beta}_k = (\beta_{1,k}, \beta_{2,k}, \dots, \beta_{p,k})^T$  were the genotypes and effect sizes  
1191 of the  $p$  genetic variants in the signal region.

1192

1193 The genetic effect of variant  $j$  on phenotype  $k$  was defined as  $\beta_{j,k} = c_j d_k \gamma_j$  to allow for  
1194 heterogeneous effect sizes among variants and phenotypes. Specifically, we generated  
1195 the causal variant indicator  $c_j$  according to a logistic model

$$\text{logit } P(c_j = 1) = \delta_0 + \delta_{l_1} A_{j,l_1} + \delta_{l_2} A_{j,l_2} + \delta_{l_3} A_{j,l_3} + \delta_{l_4} A_{j,l_4} + \delta_{l_5} A_{j,l_5},$$

1196 where  $\{l_1, \dots, l_5\} \subset \{1, \dots, 10\}$  were randomly sampled for each region. For different

1197 regions, causality of variants depended on different sets of annotations. We set

1198  $\delta_{l_i} = \log(5)$  for all annotations and varied the proportions of causal variants in the signal

1199 region by setting  $\delta_0 = \text{logit}(0.0015)$ ,  $\text{logit}(0.015)$  and  $\text{logit}(0.18)$  which corresponds to  
1200 averaging 5%, 15% and 35% causal variants in the signal region, respectively. We  
1201 considered four scenarios of phenotypic indicator  $d_k$  that reflect different underlying  
1202 genetic architectures across phenotypes:  $(d_1, d_2, d_3) = (1, 0, 0)$ ,  $(1, 0, 1)$ ,  $(1, 1, 0)$  and  
1203  $(1, 1, 1)$ . These correspond to causal variants in the signal region being associated with  
1204 (1) one phenotype only, (2) two positively correlated phenotypes, (3) two negatively  
1205 correlated phenotypes and (4) all three phenotypes. We modeled the absolute effect  
1206 sizes of causal variants using  $|\gamma_j| = c_0 |\log_{10} \text{MAF}_j|$ , such that it was a decreasing  
1207 function of MAF.  $c_0$  was set to be 0.13, 0.1, 0.1 and 0.07, respectively, to ensure a  
1208 decent power of tests under each scenario. We additionally varied the proportions of  
1209 causal variant effect size directions (signs of  $r_j$ ) by randomly generating 100%, 80%,  
1210 and 50% variants on average to have positive effects. We applied MultiSTAAR-B,  
1211 MultiSTAAR-S, MultiSTAAR-A, and MultiSTAAR-O using MAFs and all 10 annotations  
1212 together with Burden-MT, SKAT-MT and ACAT-V-MT tests. We repeated the procedure  
1213 with  $10^4$  replicates to examine the power at level  $\alpha = 10^{-7}$ . The sample size was 10,000  
1214 across all scenarios.

1215

## 1216 **Lipid Traits**

1217 Conventionally measured plasma lipids, including LDL-C, HDL-C, and triglycerides,  
1218 were included for analysis. LDL-C was either calculated by the Friedewald equation  
1219 when triglycerides were <400 mg/dl or directly measured. Given the average effect of  
1220 statins, when statins were present, LDL-C was adjusted by dividing by 0.7. Triglycerides

1221 were natural log transformed for analysis. Phenotypes were harmonized by each cohort  
1222 and deposited into the dbGaP TOPMed Exchange Area.

1223

### 1224 **Multi-trait analysis of lipid levels in the TOPMed WGS data**

1225 The TOPMed WGS data consist of multi-ethnic related samples<sup>1</sup>. Race/ethnicity was  
1226 defined using a combination of self-reported race/ethnicity from participant  
1227 questionnaires and study recruitment information (**Supplementary Note**)<sup>31</sup>. In this  
1228 study, we applied MultiSTAAR to perform multi-trait rare variant analysis of three  
1229 quantitative lipid traits (LDL-C, HDL-C and TG) using 20 study cohorts from the  
1230 TOPMed Freeze 8 WGS data. LDL-C was adjusted for the presence of medications as  
1231 before<sup>30</sup>. For each study, we first fit a linear regression model adjusting for age, age<sup>2</sup>,  
1232 sex for each race/ethnicity-specific group. In addition, for Old Order Amish (OOA), we  
1233 also adjusted for *APOB* p.R3527Q in LDL-C and TC analyses and adjusted for *APOC3*  
1234 p.R19Ter in TG and HDL-C analyses<sup>30</sup>.

1235

1236 We performed rank-based inverse normal transformation of the residuals of LDL-C,  
1237 HDL-C and TG within each race/ethnicity-specific group. We then fit a multivariate linear  
1238 mixed model for the rank normalized residuals, adjusting for 11 ancestral principal  
1239 components, ethnicity group indicators, and a variance component for empirically  
1240 derived sparse kinship matrix to account for population structure, relatedness and  
1241 correlation between phenotypes.

1242

1243 We next applied MultiSTAAR-O to perform multi-trait variant set analyses for rare  
1244 variants (MAF < 1%) by scanning the genome, including gene-centric analysis of each  
1245 protein-coding gene using five coding variant functional categories (putative loss-of-  
1246 function rare variants, missense rare variants, disruptive missense rare variants,  
1247 putative loss-of-function and disruptive missense rare variants and synonymous rare  
1248 variants); seven noncoding variant functional categories (promoter rare variants overlaid  
1249 with CAGE sites, promoter rare variants overlaid with DHS sites, enhancer rare variants  
1250 overlaid with CAGE sites, enhancer rare variants overlaid with DHS sites, UTR rare  
1251 variants, upstream region rare variants, downstream region rare variants) and rare  
1252 variants in ncRNA genes; and genetic region analysis using 2-kb sliding windows  
1253 across the genome with a 1-kb skip length. The WGS multi-trait rare variant analysis  
1254 was performed using the R packages MultiSTAAR (version 0.9.7,  
1255 <https://github.com/xihaoli/MultiSTAAR>) and STAARpipeline (version 0.9.7,  
1256 <https://github.com/xihaoli/STAARpipeline>). The WGS rare variant single-trait analysis of  
1257 LDL-C, HDL-C and TG was performed using the R package STAARpipeline (version  
1258 0.9.7, <https://github.com/xihaoli/STAARpipeline>). Both multi-trait and single-trait  
1259 analyses results were summarized and visualized using the R package  
1260 STAARpipelineSummary (version 0.9.7,  
1261 <https://github.com/xihaoli/STAARpipelineSummary>).

1262

## 1263 **Genome build**

1264 All genome coordinates are given in NCBI GRCh38/UCSC hg38.

1265

1266 **Statistics and reproducibility**

1267 Sample size was not predetermined. The multi-trait analysis consists of 20 study  
1268 cohorts of TOPMed Freeze 8 and had 61,838 samples with lipid traits. We did not use  
1269 any study design that required randomization or blinding.

1270

1271 **Data availability**

1272 This paper used the TOPMed Freeze 8 WGS data and lipids phenotype data. Genotype  
1273 and phenotype data are both available in database of Genotypes and Phenotypes. The  
1274 TOPMed WGS data were from the following twenty study cohorts (accession numbers  
1275 provided in parentheses): Old Order Amish (phs000956.v1.p1), Atherosclerosis Risk in  
1276 Communities Study (phs001211), Mt Sinai BioMe Biobank (phs001644), Coronary  
1277 Artery Risk Development in Young Adults (phs001612), Cleveland Family Study  
1278 (phs000954), Cardiovascular Health Study (phs001368), Diabetes Heart Study  
1279 (phs001412), Framingham Heart Study (phs000974), Genetic Study of Atherosclerosis  
1280 Risk (phs001218), Genetic Epidemiology Network of Arteriopathy (phs001345), Genetic  
1281 Epidemiology Network of Salt Sensitivity (phs001217), Genetics of Lipid Lowering  
1282 Drugs and Diet Network (phs001359), Hispanic Community Health Study - Study of  
1283 Latinos (phs001395), Hypertension Genetic Epidemiology Network and Genetic  
1284 Epidemiology Network of Arteriopathy (phs001293), Jackson Heart Study (phs000964),  
1285 Multi-Ethnic Study of Atherosclerosis (phs001416), San Antonio Family Heart Study  
1286 (phs001215), Genome-wide Association Study of Adiposity in Samoans (phs000972),  
1287 Taiwan Study of Hypertension using Rare Variants (phs001387), and Women's Health

1288 Initiative (phs001237). The sample sizes, ancestry and phenotype summary statistics of  
1289 these cohorts are given in **Supplementary Table 2**.

1290

1291 The functional annotation data are publicly available and were downloaded from the  
1292 following links: GRCh38 CADD v1.4 (<https://cadd.gs.washington.edu/download>);  
1293 ANNOVAR dbNSFP v3.3a ([https://annovar.openbioinformatics.org/en/latest/user-](https://annovar.openbioinformatics.org/en/latest/user-guide/download)  
1294 [guide/download](https://annovar.openbioinformatics.org/en/latest/user-guide/download)); LINSIGHT (<https://github.com/CshISiepelLab/LINSIGHT>); FATHMM-  
1295 XF (<http://fathmm.biocompute.org.uk/fathmm-xf>); FANTOM5 CAGE  
1296 (<https://fantom.gsc.riken.jp/5/data>); GeneCards (<https://www.genecards.org>; v4.7 for  
1297 [hg38](https://www.genecards.org)); and Umap/Bismap (<https://bismap.hoffmanlab.org>; 'before March 2020' version).  
1298 In addition, recombination rate and nucleotide diversity were obtained from Gazal et  
1299 al<sup>50</sup>. The whole-genome individual functional annotation data was assembled from a  
1300 variety of sources and the computed annotation principal components are available at  
1301 the Functional Annotation of Variant-Online Resource (FAVOR) site  
1302 (<https://favor.genohub.org>)<sup>51</sup> and the FAVOR database  
1303 (<https://doi.org/10.7910/DVN/1VGTJI>)<sup>52</sup>.

1304

### 1305 **Code availability**

1306 MultiSTAAR is implemented as an open source R package available at  
1307 <https://github.com/xihaoli/MultiSTAAR> and  
1308 <https://content.sph.harvard.edu/xlin/software.html>. Data analysis was performed in R  
1309 (4.1.0). STAAR v0.9.7 and MultiSTAAR v0.9.7 were used in simulation and real data  
1310 analysis and implemented as open-source R packages available at

1311 <https://github.com/xihaoli/STAAR> and <https://github.com/xihaoli/MultiSTAAR>. The  
1312 assembled functional annotation data were downloaded from FAVOR using Wget  
1313 (<https://www.gnu.org/software/wget/wget.html>).

1314

## 1315 **References**

- 1316 50. Gazal, S. *et al.* Linkage disequilibrium–dependent architecture of human complex  
1317 traits shows action of negative selection. *Nature Genetics* **49**, 1421-1427 (2017).  
1318 51. Zhou, H. *et al.* FAVOR: functional annotation of variants online resource and  
1319 annotator for variation across the human genome. *Nucleic Acids Research* **51**,  
1320 D1300-D1311 (2023).  
1321 52. Zhou, H., Arapoglou, T., Li, X., Li, Z. & Lin, X. FAVOR Essential Database. V1  
1322 Edition (Harvard Dataverse, 2022).