

1 **Expanding the pool of public controls for GWAS via a method for combining genotypes**  
2 **from arrays and sequencing**

3 Ravi Mathur<sup>\*1</sup>, Fang Fang<sup>\*1</sup>, Nathan Gaddis<sup>1</sup>, Dana B. Hancock<sup>1</sup>, Michael H. Cho<sup>2,3</sup>, John E.  
4 Hokanson<sup>4</sup>, Laura J. Bierut<sup>5</sup>, Sharon M. Lutz<sup>6</sup>, Kendra Young<sup>4</sup>, Albert V. Smith<sup>7,8</sup>, NHLBI Trans-  
5 Omics for Precision Medicine (TOPMed) Consortium<sup>^</sup>, Edwin K. Silverman<sup>2,3</sup>, Grier P. Page<sup>1,9</sup>,  
6 Eric O. Johnson<sup>+1,9</sup>

7 1. GenOmics, Bioinformatics, and Translational Research Center, RTI International,  
8 Research Triangle Park, North Carolina, USA

9 2. Channing Division of Network Medicine, Brigham and Women's Hospital, Boston,  
10 Massachusetts, USA

11 3. Division of Pulmonary and Critical Care Medicine, Brigham and Women's Hospital,  
12 Boston, Massachusetts, USA

13 4. Department of Epidemiology, Colorado School of Public Health, University of Colorado  
14 Denver, Denver, Colorado, USA

15 5. Department of Psychiatry, Washington University, St. Louis, MO, USA

16 6. PRecisiOn Medicine Translational Research (PROMoTeR) Center, Department of  
17 Population Medicine, Harvard Medical School and Harvard Pilgrim Health Care, Boston,  
18 Massachusetts, USA

19 7. Center for Statistical Genetics, University of Michigan, Ann Arbor, Michigan, USA

20 8. Department of Biostatistics, University of Michigan, Ann Arbor, Michigan, USA

21 9. Fellow Program, RTI International, Research Triangle Park, North Carolina, USA

22

23 \* Authors contributed equally to this manuscript

24 ^ A complete list of NHLBI Trans-Omics for Precision Medicine (TOPMed) Consortium members  
25 is shown in Supplementary Table 1.

26 + Corresponding Author

27 **Abstract**

28           Genome-wide association studies (GWAS) have made impactful discoveries for complex diseases,  
29 often by amassing very large sample sizes. Yet, GWAS of many diseases remain underpowered,  
30 especially for non-European ancestries. One cost-effective approach to increase sample size is to  
31 combine existing case-only cohorts with public controls, but this approach is limited by the need for a  
32 large overlap in variants across genotyping arrays and the scarcity of non-European controls. We  
33 developed and validated a protocol, Genotyping Array-WGS Merge (GAWMerge), for combining  
34 genotypes from arrays and whole genome sequencing, ensuring complete variant overlap, and allowing  
35 for diverse samples like Trans-Omics for Precision Medicine to be used. Our protocol involves phasing,  
36 imputation, and filtering. We illustrated its ability to control type I error and recover known disease-  
37 associated signals across technologies, independent datasets, and ancestries in smoking-related cohorts.  
38 GAWMerge enables genetic studies to leverage existing cohorts to validly increase sample size and  
39 enhance discovery.

40

41

42           Genome-wide association studies (GWAS) offer a powerful tool for identifying genetic  
43 variants for complex diseases, especially when large sample sizes are amassed. For diseases  
44 with limited sample sizes or for which case-only cohorts are available, public controls, who are  
45 not assessed for the disease, can be used without bias to cost effectively improve statistical  
46 power and novel locus discovery, if the disease prevalence is low in the general population.<sup>1-5</sup>  
47 Combining cases and controls in this way is feasible even with samples genotyped on different  
48 array-based technologies<sup>6-9</sup>. A significant limitation of combining disease study cases with public  
49 controls is that unbiased results are only achieved using the intersecting set of variants  
50 genotyped across all arrays and cohorts being combined.<sup>9</sup> This limitation effectively prevents  
51 combining cohorts where the number of shared genotyped variants is too small to form the  
52 basis for imputation or to provide whole genome coverage. An in-depth comparison of the  
53 Illumina HumanHap, Illumina OmniExpress, and Affymetrix 6.0 arrays found over 2,000,000  
54 single nucleotide polymorphisms (SNPs) in union but only 75,000 variants that intersect across  
55 all arrays<sup>10</sup>. Additionally, reliance on array-based technology prevents use of expanding whole  
56 genome sequencing (WGS) resources with high representation of non-European ancestry  
57 groups, like the Trans-Omics for Precision Medicine (TOPMed) program, for public controls.  
58 Being able to combine case and public control genotypes from array- and/or sequencing-based  
59 platforms opens up the increasing set of WGS resources for new GWAS. As of January 2021,  
60 there are at least 217 case-only studies containing >136,000 samples across many genotyping  
61 platforms in the database of Genetics and Phenotypes (dbGaP) (query = 'case set[Study  
62 Design]'). There are >227,000 public controls with WGS data in resources such as TOPMed  
63 (>155,000 samples)<sup>11</sup>, UK BioBank (>50,000 samples)<sup>12</sup>, Gabriella Miller Kids First Pediatric  
64 Consortium (>21,000 samples), and GenomeAsia100K Project (>1,700 individuals)<sup>13</sup>, which are  
65 eligible to be combined with these case-only datasets for GWAS.

66           The NHLBI-supported TOPMed program<sup>11</sup> with its collection of >155,000 human  
67 subjects with WGS data affords an unparalleled opportunity to leverage public controls and

68 greatly expand GWAS sample sizes. With such a large sample size and one of the most  
69 genetically diverse datasets (40% European, 31% African, 16% Hispanic, 9% Asian, and 4%  
70 Others) available, TOPMed has the potential to overcome the aforementioned challenges of  
71 applying public controls, as the WGS data should overlap all variants measured on arrays, and  
72 the representation of non-European populations will enhance the availability of diverse public  
73 controls.

74 While incorporating public controls to maximize the utility of genetic discovery is  
75 desirable, there is no established approach to validly combine array- and sequencing-based  
76 genotype data. Each of these technologies has its own strengths, weaknesses, and different  
77 inter- and intra-technology measurement properties that complicate combining data across  
78 technologies. Here, we developed a protocol, Genotype Array-WGS Merge (GAWMerge), to  
79 combine genotypes from array and WGS to conduct GWAS analyses. We illustrate our  
80 protocol's validity and its utility using TOPMed WGS samples as public controls combined with  
81 case-only array-genotyped cohorts.

82

83

84

## 85 **Results**

86 **Protocol to Integrate Array and WGS data.** GAWMerge is a protocol that we developed to  
87 integrate array and WGS genotyping technologies that minimizes false positives while  
88 discovering true association signals. Details of the protocol development process are provided  
89 in the **Methods** section. The final protocol consists of eight major steps (Figure 1): (1) select  
90 control dataset(s) with WGS genotype data; (2) extract the SNPs from the WGS data of the  
91 control samples that match those for the array-genotyped case samples; (3) independently  
92 subject the case and control samples to the same quality control (QC) procedure (further details  
93 in the **Methods**); (4) phase the case and control samples with the same software (further details  
94 in the **Methods**); (5) merge the phased case and control data and impute to the desired  
95 reference genome (e.g., 1000Genome, TOPMed reference panel); (6) filter out genotyped SNPs  
96 with low quality (empirical  $ER^2 < 0.9$ )<sup>14</sup> and re-impute; (7) test SNP associations with phenotype  
97 of interest in case and control samples combined; and (8) filter association results for minor  
98 allele frequency (MAF), imputation quality ( $R^2$ ), and difference in imputation quality.

99 For selection of controls in **step 1**, it is crucial to choose samples with an ancestral  
100 composition consistent with the case samples, as population stratification is a strong  
101 confounding factor for GWAS analysis. Additional demographic (e.g., age, sex) and clinical  
102 variables (e.g., smoking status) should be considered based on the datasets being combined.

103 Our previous work<sup>9</sup> suggested potential bias in association testing when using  
104 genotypes imputed from the full sets of SNPs from different genotyping arrays. Starting from the  
105 intersection of genotyped SNP sets avoids such bias (**step 2**). We employed the same strategy  
106 for merging array and WGS genotypes, but because of the full genome coverage of WGS, the  
107 entire set of array SNPs were used. The array and WGS data were then independently QC'd  
108 using the same QC steps (**step 3**). This then was followed by phasing, merging, and imputation  
109 (**steps 4-5**). To further reduce potential bias between the array-genotyped and WGS-derived

110 SNPs, a second round of imputation is performed after removing genotyped SNPs with low  
111 empirical  $R^2$  ( $ER^2 < 0.9$ , **step 6**, Supplementary Figure 1). Finally, following association testing  
112 (**step 7**), filtering based on MAF ( $> 0.01$ ), imputation quality ( $R^2 > 0.8$ ), and imputation quality  
113 difference between cases (i.e., array data) and controls (i.e., WGS data) is **step 8** ( $|R^2_{array} -$   
114  $R^2_{WGS}| < 0.1$ , Supplementary Figure 2) which minimizes technical variation in the combined  
115 case/control data. More details regarding the development of the protocol can be found in the  
116 “Protocol Development” section of **Methods**.

117

118 **Protocol Evaluation Design.** To evaluate the performance of GAWMerge, we used three  
119 smoking-related datasets: Collaborative Genetic Study of Nicotine Dependence (COGEND)<sup>15,16</sup>,  
120 Genetic Epidemiology of COPD (COPDGene) study<sup>17</sup>, and Evaluation of COPD Longitudinally  
121 to Identify Predictive Surrogate End-points (ECLIPSE)<sup>18</sup>. As indicated in Table 1, the three  
122 datasets have different array platforms, providing the opportunity to assess the performance of  
123 the protocol in different settings. In both COPDGene and ECLIPSE, the COPD diagnosis  
124 followed the Global Initiative for Chronic Obstructive Lung Disease (GOLD) severity  
125 classifications, and COPD cases were defined as GOLD Grade 2–4 COPD (moderate, severe,  
126 and very severe COPD)<sup>19</sup>. The study design to evaluate GAWMerge across (a) genotyping  
127 technology (ensuring no technology driven false positives), (b) type-I error (ensuring minimal  
128 false positive associations), and (c) recovery of known GWAS hits (demonstrating capture of  
129 true positives) is presented in Figure 2.

130

131 **Reproducibility across genotyping technologies.** COPDGene has both array and WGS  
132 genotype data on the same samples available through TOPMed. Genotypes derived from array  
133 and whole genome sequencing data for the same samples should be consistent but are often  
134 not.<sup>20,21</sup> To evaluate the consistency of genotyping, we performed a technical comparison of

135 array and WGS data using the same set of samples from COPDGene (n=3,235 with African-  
136 American ancestry). The array data were phased independently and integrated with the WGS  
137 phased data available in TOPMed, followed by imputation and association testing using  
138 genotyping platform as the outcome. If the array- and WGS-derived genotypes for the same set  
139 of samples were equivalent, one would expect to observe no significant associations, but in fact  
140 we observed many false positives (Supplementary Figure 3).

141 We suspected that the false positives we observed derived from the phasing step since  
142 phasing of the array and WGS genotypes was based on different sets of variants. In addition,  
143 the TOPMed phased WGS data were derived from the samples of all studies<sup>22</sup>, which is  
144 different from the sample set we used, the COPDGene cohort, for phasing the array data. We  
145 repeated the technical comparison, using the same set of QC-validated variants and samples  
146 (Figure 2a) as the basis for separate phasing of the array and WGS data, followed by the  
147 subsequent steps in GAWMerge (Figure 1). The array data were specified as the case group for  
148 association testing, and the WGS data were specified as the control group, for European  
149 ancestry (EA) and African ancestry (AA) separately. The results (Supplementary Figure 3)  
150 confirmed that phasing based on a common set of variants and samples followed by the  
151 additional steps of GAWMerge eliminated false positives and made array and WGS data  
152 comparable for conducting GWAS.

153

154 **Controlling type I error in case-only vs. public control GWAS.** We assessed type-I  
155 error in a comprehensive analysis involving three smoking-related datasets and their meta-  
156 analysis, as shown in Figure 2b. To fully leverage the large sample size of the COPDGene  
157 dataset, we evenly divided the EA samples into two subsets: EA1 and EA2. COPDGene EA1  
158 included all participants diagnosed with COPD (N=2,736) and randomly sampled participants  
159 with no COPD (N=515). The resulting ratio of individuals with COPD in COPDGene EA1 (84%)



160 was close to the ratio in ECLIPSE EA (87%). Three GWAS were conducted to assess type-I  
161 error, as follows: (1) array data from COPDGene EA1 (N=3,251) vs. WGS from ECLIPSE EA  
162 (N=1,461); (2) array data from COGENE EA (N=1,961) vs. WGS data from COPDGene EA2  
163 (with no COPD, N=3,251); and (3) array data from COGENE AA (N=712) vs. WGS from  
164 COPDGene AA (N=1,710). All association models include ten principal components as  
165 covariates to account for population substructure. COPDGene, COGENE, and ECLIPSE are all  
166 smoking cohorts and ratios of COPD were consistent across array and WGS datasets, thus we  
167 expected no genome-wide significant association signals (controlled type 1 error). Applying  
168 GAWMerge to these data we observed no false positive signals in each separate GWAS  
169 analysis (Supplementary Figure 5) and in their meta-analysis (Figure 3) results.

170

171 **Recovery of known COPD loci in case-only vs. public control GWAS.** The last  
172 evaluation step was to recover known GWAS hits for COPD.<sup>19,23</sup> As shown in Figure 2c, we  
173 conducted three GWAS for COPD, as follows: (1) COPD cases from COPDGene EA with WGS  
174 data (N=2,736) vs. controls from COGENE EA with array data (N=1,961); (2) COPD cases from  
175 ECLIPSE EA with array data (N=1,764) vs. controls from COPDGene EA with WGS data  
176 (N=2,475); and (3) COPD cases from COPDGene AA with WGS data (N=813) vs. controls from  
177 COGENE AA with array data (N=712). Because COPD is highly comorbid with smoking history,  
178 only smokers (current and former) were used as controls to compare with COPD cases across  
179 these GWAS analyses. All association models include ten principal components as covariates  
180 to account for population substructure. Results for each GWAS analysis are presented in  
181 Supplementary Figure 6. Meta-analysis of the 3 analyses successfully recovered 5 out of 7 loci  
182 reported as COPD-associated (Figure 4 and Table 2) at genome-wide significance ( $P < 5 \times$   
183  $10^{-8}$ , Supplementary Table 2). The direction of association for all recovered SNPs was the  
184 same as previously reported<sup>24</sup>. The two SNPs that did not exceed the genome-wide significance

185 threshold were nominally associated at  $P < 0.05$  in our analysis. These two SNPs were missing  
186 in Analysis 1 (COPD cases with WGS data from COPDGene EA Vs. smoking controls with  
187 array data from COGEND EA) due to the filters applied with the protocol; the reduced power  
188 caused by their missingness likely explain the lower significance level observed.

## 189 Discussion

190 In summary, we present GAWMerge, a protocol for integrating array and WGS genotype  
191 data to conduct GWAS with a case-only and public control design. This protocol overcomes  
192 previous obstacles to using public controls<sup>9</sup>. The ability to use WGS data for public controls 1)  
193 ensures complete overlap with variants on any array used for genotyping of cases, and 2)  
194 provides a much larger pool of public controls to draw from, especially for non-Europeans, from  
195 ancestrally diverse resources like TOPMed. In our proof-of-concept study, we applied  
196 GAWMerge to WGS data from TOPMed (specifically, COPDGene and ECLIPSE cohorts) as  
197 public controls for array-genotyped case datasets. We first showed that the two genotyping  
198 technologies are compatible by comparing array- and WGS-derived genotypes for the same  
199 samples from COPDGene and demonstrating a lack of false positives. We then showed that  
200 GAWMerge controls type I error, as evidenced by the expected lack of genome-wide significant  
201 findings in a GWAS meta-analysis comparing smoker cases vs. smoker controls from  
202 independent datasets. Lastly, GAWMerge recovered known COPD-associated findings from  
203 Hobbs et al.<sup>24</sup> including *CHRNA3* on chromosome 15, *FAM13A* on chromosome 4, *CYP2A6* on  
204 chromosome 19, *TGFB2* on chromosome 1, and *HHIP* on chromosome 4. The key aspects of  
205 the protocol that provide these unbiased findings are 1) phasing the array and WGS data  
206 independently using only the intersection of variants across technologies and 2) including the  
207 empirical  $R^2$  and  $R^2$  difference filters to remove poorly imputed and differently imputed variants.

208 The development of GAWMerge was done with TOPMed WGS and array genotyped-  
209 data, although it can be applied using any case-only array-genotyped data with other WGS data  
210 resources (e.g., UK BioBank<sup>12</sup>, Gabrielle Miller Kids First and/or GenomeAsia 100K<sup>13</sup> data). To  
211 incorporate new data, it will be important to identify the phenotypic data which will be used to  
212 combine controls with available cases. For example, we selected controls based on the smoking  
213 status of the cohorts to minimize bias due to smoking. Additional phenotypic and clinical data,

214 such as sex and age distributions, should be considered when selecting the most appropriate  
215 controls for combining with available cases. In this study we combined cases and controls with  
216 the same ancestry to minimize bias. Further work is needed to evaluate GAWMerge for trans-  
217 ancestry and mega analysis GWAS<sup>25</sup>. GAWMerge was developed with imputation using the  
218 thousand genomes reference population, although method can be applied using other reference  
219 populations, such as the TOPMed reference population on the Michigan Imputation Server<sup>14</sup>.  
220 Since TOPMed samples are used as controls in GAWMerge, there will be sample overlap  
221 between the input data and the TOPMed reference population, which may cause bias and must  
222 be applied cautiously. Further work is needed to evaluate the bias of such an imputation  
223 strategy.

224 GAWMerge has some limitations. First, careful consideration of not only ancestry, sex,  
225 and age distributions, but other systematic differences between a given case-only cohort and  
226 public controls, like smoking status, is essential to unbiased use of public controls and  
227 application of GAWMerge. All association analysis conducted included ten principal components  
228 as covariates to account for population substructure, although applying GWAS in as  
229 homogenous population as possible is desirable. This requirement places some limits on the  
230 public controls that can be used for any given case-only cohort. Second, the additional QC  
231 steps might mask some real trait-associated variants. In the attempt to recover the known  
232 genetic variants associated with COPD, there were two loci (*RIN3* and *MMP3/12*) not reaching  
233 the genome-wide significance in the meta-analysis (Table 2). The three SNPs were filtered out  
234 in the first GWAS, comparing COPD cases in COPDGene EA with WGS data and smoking  
235 controls in COGEND EA with array data, due to high  $R^2$  difference between the WGS and array  
236 data. Thus, GAWMerge may lose some sensitivity while controlling type I errors. There is also  
237 the potential for reduced power to detect COPD associated genetic variants here due to the  
238 missingness of lung function phenotypes in COGEND public controls, with power being reduced

239 relative to the amount of COPD status misclassification among these controls. Third, when  
240 GAWMerge has been tested as an application of GWAS, it is limited by the MAF and genomic  
241 coverage on array genotyping technologies. Since GAWMerge extracts only SNPs within the  
242 array technology, the complete coverage of WGS (over 410 million variants within TOPMed  
243 WGS data<sup>22</sup>) is not fully utilized. Therefore, those rare variants and large insertions/deletions  
244 only detected in WGS data were lost during the extraction and merging processes  
245 (Supplementary Table 3). However, coming from a case-only dataset with array-based  
246 genotyping, the dominant scenario for use of GAWMerge, the WGS is a substantial strength,  
247 accounting for all the array genotyped variants except for technology based regional loss of  
248 variants. With our strategy of WGS data as public controls for GWAS, there will be regional loss  
249 in specific areas depending on the array technology design and quality control of the  
250 sequencing. A complete analysis of different regional genetic variants covered specifically by  
251 array-genotyping platforms or sequencing will be beneficial to calibrate the application of  
252 GAWMerge in the future.<sup>26,27</sup>

253 Overall, GAWMerge presents a practical application of integrating case-only array-  
254 genotyped data with WGS data as public controls to enable new GWAS and enhance the  
255 potential for discovering novel genetic loci. It is a general approach for integrating array and  
256 WGS genotyping technologies, breaking any barriers in such integration. The substantial  
257 availability of case-only datasets in public repositories and collected across many consortia  
258 makes the protocol broadly applicable. With >155,000 samples with WGS data within the  
259 TOPMed program, this an ample resource for selecting public controls for a variety of case-only  
260 disease datasets. With WGS data, the overlap of measured variants across genotyping  
261 platforms is overcome. Furthermore, the diversity of individuals within the TOPMed (>47,000  
262 African, >23,000 Hispanic/Latino, and >13,000 Asian ancestries) and increasing representation  
263 in other resources make widespread use of non-European public controls realistic. With many

264 other WGS resources being launched and released, the potential to use public controls to  
265 increase sample size and leverage case-only cohorts is just beginning.

266

## 267 **Methods**

268 **Dataset Descriptions.** The Trans-Omics for Precision Medicine (TOPMed) program aims to  
269 improve understanding of the diseases through the integration of Whole Genome Sequencing  
270 (WGS) and other omics data from pre-existing parent studies having large samples of human  
271 subjects. The two studies used in this work, Genetic Epidemiology of Chronic Obstructive  
272 Pulmonary Disease (COPDGene) and Evaluation of COPD Longitudinally to Identify Predictive  
273 Surrogate Endpoints (ECLIPSE), are both part of TOPMed. As of February 2020, TOPMed has  
274 gathered data from ~155k participants with rich phenotypic data. TOPMed prioritizes to increase  
275 ancestral and ethnic diversity, so ~60% of the sequenced participants are of non-European  
276 ancestry (31% African, 16% Hispanic, 9% Asian, and 4% Others).

277 COPDGene (ClinicalTrials.gov: NCT00608764) is an ongoing study of over 10,000 non-  
278 Hispanic White and African American cigarette smokers. It was designed to investigate COPD  
279 and other smoking-related lung diseases<sup>17</sup>. COPDGene subjects were initially genotyped for ~1  
280 million single nucleotide polymorphisms (SNPs) using the HumanOmniExpress array (Illumina,  
281 San Diego, CA). As part of TOPMed freeze 6a, WGS was conducted on 10,372 subjects.  
282 Among them, 9,732 subjects are overlapped with the subjects in the parent study having array  
283 genotyped data, and thus were used in our analyses.

284 ECLIPSE was an observational study launched in 2006<sup>18</sup>. It recruited 2,164 COPD  
285 subjects, 337 smoking controls, and 245 nonsmoking controls. The genotype data with Illumina  
286 HumanHap550v3.0 array (~550,000 SNPs) included 1,764 COPD subjects, 217 smoking  
287 controls, and 178 non-smoking controls. In TOPMed freeze 6a, WGS was conducted on 1,271  
288 COPD subjects and 190 smoking controls.

289 COGEND was initiated in 2001 as a genetic study of nicotine dependence<sup>15,16</sup>. Nicotine  
290 dependent cases and non-dependent smoking controls were identified and recruited from  
291 Detroit and St. Louis. Over 2,900 donated blood samples were collected and used to genotype

292 ~2.5 million SNPs using the HumanOmni2.5 array. After QC, 2,673 subjects were kept for  
293 following analyses.

294

295 **GAWMerge development.** Below we provide further details on the protocol steps, and  
296 iterations used to devise the recommended thresholds.

297

298 Quality control (QC). We performed standard QC steps for both array genotyped data and the  
299 subset of WGS data extracted in step 2 using PLINK<sup>28</sup>. Samples failing sex check or with >3%  
300 missing data were excluded. SNPs with missing rate >3% or that failed Hardy-Weinberg  
301 Equilibrium check ( $p < 1e-4$ ) were excluded from the study. A structure analysis was conducted  
302 to match ancestries to 1000 genomes reference haplotypes and mis-classified samples were  
303 excluded. In addition, we adopted standard TOPMed filters (<https://topmed.nhlbi.nih.gov/>) for  
304 variant selection. The variants that were labeled as follows were excluded: SVM (support vector  
305 machine score more negative than -0.5 and hence fails the SVM filter), CEN (falls in a  
306 centromeric region with inferred reference sequence), DISC (more than 5 percent Mendelian  
307 inconsistencies), EXHET (has excessive heterozygosity with HWE p-value less than  $1e-6$ ) or  
308 CHRXHET (has excessive heterozygosity in male chrX).

309

310 Combining Array and WGS Data. **GAWMerge**, a protocol for integrating array and WGS data is  
311 shown in Figure 1 and described in more detail in the **Results**. The WGS data were first  
312 prepared by extracting the selected control samples and the variants available within the array  
313 genotyping data. Utilizing the intersection of variants was important, as many false positives  
314 were introduced without this step<sup>9</sup>. This extraction of samples and variants was performed by  
315 BCFtools<sup>29</sup>. After QC, the intersection of SNPs between the array and WGS data was extracted,  
316 and the datasets were phased independently using SHAPEIT2<sup>30,31</sup>. The datasets were then  
317 merged using BCFtools<sup>29</sup>.



318

319 Imputation strategy. The merged array and WGS data were first imputed using Minimac4<sup>14</sup>  
320 using the thousand genomes phase 3 version 5 EUR and AFR super populations for EA and AA  
321 samples, respectively. The reference panel includes 503 EUR and 661 AFR samples with data  
322 on GRCh37 genome version. TOPMed WGS data was converted from genome version  
323 GRCh38 to GRCh37 to match the reference and array-genotyped data. Besides applying the  
324 standard imputation quality measurement  $R^2$ , we also observed poorly imputed variants  
325 indicated by Empirical  $R^2$  ( $ER^2$ ).  $ER^2$  was defined only for genotyped variants as the squared  
326 correlation between leave-one-out imputed dosages and the true, observed genotypes. Under  
327 our first test for controlling type I error (Figure 2b), array data from COPDGene EA1 (N=3,251)  
328 and WGS data from ECLIPSE EA (N=1,461), we expected no genome-wide significant  
329 associations since all individuals were smokers and no disease was being tested between the  
330 datasets. Without the  $ER^2$  filter, we found many false positives (Supplementary Figure 1a)  
331 based around the variant on chromosome 10 (chr10:32370743,  $ER^2 = 0.391$ , MAF=0.068). We  
332 recommend removing such genotyped SNPs with  $ER^2 < 0.9$  from the analysis and re-running  
333 imputation without these variants included. With this and other low-quality variants removed,  
334 false positives were controlled (Supplementary Figure 1b). With the  $ER^2$  filter of 0.9, we found  
335 that 81.1% of SNPs met this criterion (Supplementary Figure 1c) and these removed SNPs  
336 were scattered across the genome (Supplementary Figure 1d).

337 Filtering association test results. Association analysis was conducted using rvTest<sup>32</sup> with ten  
338 principal components included to account for population substructure. Besides the common  
339 filters for minor allele frequency (MAF>0.01) and imputation quality ( $R^2 > 0.8$ ), we also  
340 investigated the imputation quality difference between array-genotyped samples and WGS-  
341 genotyped samples by comparing the imputation quality within each sample type,  $R^2_{array}$  and  
342  $R^2_{WGS}$ . We verified that the imputation quality between the two types of data were similar.

343 However, some outliers ( $|R_{array}^2 - R_{WGS}^2| \geq 0.1$ ) were a major source of false positives, and  
344 were removed from the results as a post-association testing filter. Using the same test between  
345 COGEND and COPDGene EA sample comparison, inflation of GWAS P-values was apparent  
346 when  $|R_{array}^2 - R_{WGS}^2| \geq 0.1$ , but otherwise no inflation was observed (Supplementary Figure 2a).  
347 An imputation quality difference of  $\geq 0.1$  only filtered out about 5% of variants (Supplementary  
348 Figure 2b), and the removed variants were scattered throughout the genome (Supplementary  
349 Figure 2c).

350

### 351 **GAWMerge implementation.**

352 GAWMerge was developed within the DNANexus computing environment  
353 (<https://www.dnanexus.com/>) and the BioData Catalyst ecosystem<sup>33</sup>. The protocol within the  
354 DNANexus computing environment used docker images, which have been packaged together  
355 into DNANexus applications. The BioData Catalyst ecosystem<sup>33</sup> protocol was implemented in  
356 the common workflow language (CWL); therefore, it is interoperable in other computing  
357 ecosystems. Both implemented workflows are built using the same docker images of the  
358 underlying software programs ([https://github.com/RTIInternational/biocloud\\_docker\\_tools](https://github.com/RTIInternational/biocloud_docker_tools) and  
359 <https://hub.docker.com/u/rtibiocloud>). The protocol has been written to easily adapt to plink or  
360 vcf formats of the genotype files, therefore either are acceptable. The BioData Catalyst workflow  
361 leverages key services, tools, and workflows available within the ecosystem including BioData  
362 Catalyst Powered by Gen3, BioData Catalyst Powered by PIC-SURE, and BioData Catalyst  
363 Powered by Seven Bridges. These tools make discovery of data for use as public controls easy  
364 with their easy-to-use web interface.

365 To discover optimal controls to combine with available cases, TOPMed phenotypic data  
366 were easily accessible using the Gen3 and PIC-SURE tools within the BioData Catalyst  
367 ecosystem. With these tools, users identify which studies were comparable for use as public

368 controls, urge the access request for these studies within dbGaP, and then use as public

369 controls with the protocol.

370 Computation of GAWMerge is comparable to other GWAS efforts. For example, in the  
371 analysis comparing ECLIPSE WGS data and COPDGene EA array data, phasing the 10,302  
372 variants on chromosome 10 (overlapped with the array data) of the 1,461 samples in ECLIPSE  
373 WGS data took ~9 hours using a machine with 32GB memory and 16 CPUs. The following  
374 imputation ran on a machine with 16GB memory and 4 CPUs for 2 hour and 37 minutes. Then  
375 the re-imputation runs for similar amount of time.

376

377 **Table 1.** Dataset characteristics.

		COGEND	COPDGene	ECLIPSE	
Array type		Illumina HumanOmni 2.5	Illumina HumanOmni1-Quad_v1-0_B	Illumina HumanHap550v3.0	
Array-genotyped data	N, SNPs	2,443,179	1,051,295	561,466	
	Participants, total N	2,673	9,962	2,159	
	Ancestry group, N (%)	European	1,961 (73%)	6,664 (67%)	2,159 (100%)
		African American	712 (27%)	3,298 (33%)	NA
	Sex, N (%)	Males	1,019 (38%)	5,333 (54%)	1,367 (63%)
		Females	1,654 (62%)	4,629 (46%)	792 (37%)
	COPD diagnosis, N (%)	Yes	NA	4,280 (43%)	1,764 (82%)
No			3,632 (36%)	395 (14%)	
Age (mean±SD)		36.6±5.6	59.6±9.0	62.2±8.2	
WGS-genotyped data*	Participants, total N	NA	9,737	1,484	
	Ancestry group, N (%)	European	NA	6,502 (67%)	1,461 (98%)
		African American		3,235 (33%)	23 <sup>¶</sup> (2%)
	Sex, N (%)	Males		5,213 (54%)	933 (64%)
		Females		4,524 (46%)	528 (36%)
	COPD diagnosis, N (%)	Yes		4,186 (43%)	1,271 (87%)
		No		3,549 (36%)	190 (13%)
Age (mean±SD)				59.6±9.0	62.7±7.7

378 \* All WGS genotyped data are from TOPMed freeze6a.

379 <sup>¶</sup> The number of African American in ECLIPSE is too small and excluded from following analysis.

380

381

382 **Table 2.** Recovery of GWAS-identified variants, following application of our protocol to each of 3 GWAS  
 383 and their meta-analysis, compared to published risk loci for COPD with combined data from COPDGene,  
 384 ECLIPSE, NETT/NAS, and GenKOLS (Norway)<sup>19</sup>.

SNP	Position	Risk Allele	Related gene	Reported (N=12,337)		Current meta-analysis (N=10,461)		
				OR	P-value	OR	Direction	P-value
rs12914385	chr15:78898723	T	<i>CHRNA3</i>	1.36	2.70E-16	1.28	+++	3.35E-16
rs4416442	chr4:89866713	C	<i>FAM13A</i>	1.36	9.44E-15	1.21	+++	2.66E-10
rs7937 <sup>23</sup>	chr19:41302706	C	<i>CYP2A6</i>	0.74	2.88E-09	0.84	---	1.91E-08
rs4846480	chr1:218598469	A	<i>TGFB2</i>	1.26	1.25E-07	1.19	+++	9.37E-08
rs13141641	chr4:145506456	T	<i>HHIP</i>	1.39	3.66E-15	1.23	?++*	2.64E-07
rs754388	chr14:93115410	C	<i>RIN3</i>	1.33	6.69E-08	1.12	?++*	0.020
rs626750	chr11:102720945	G	<i>MMP3/12</i>	1.36	5.35E-09	1.14	?++*	0.005

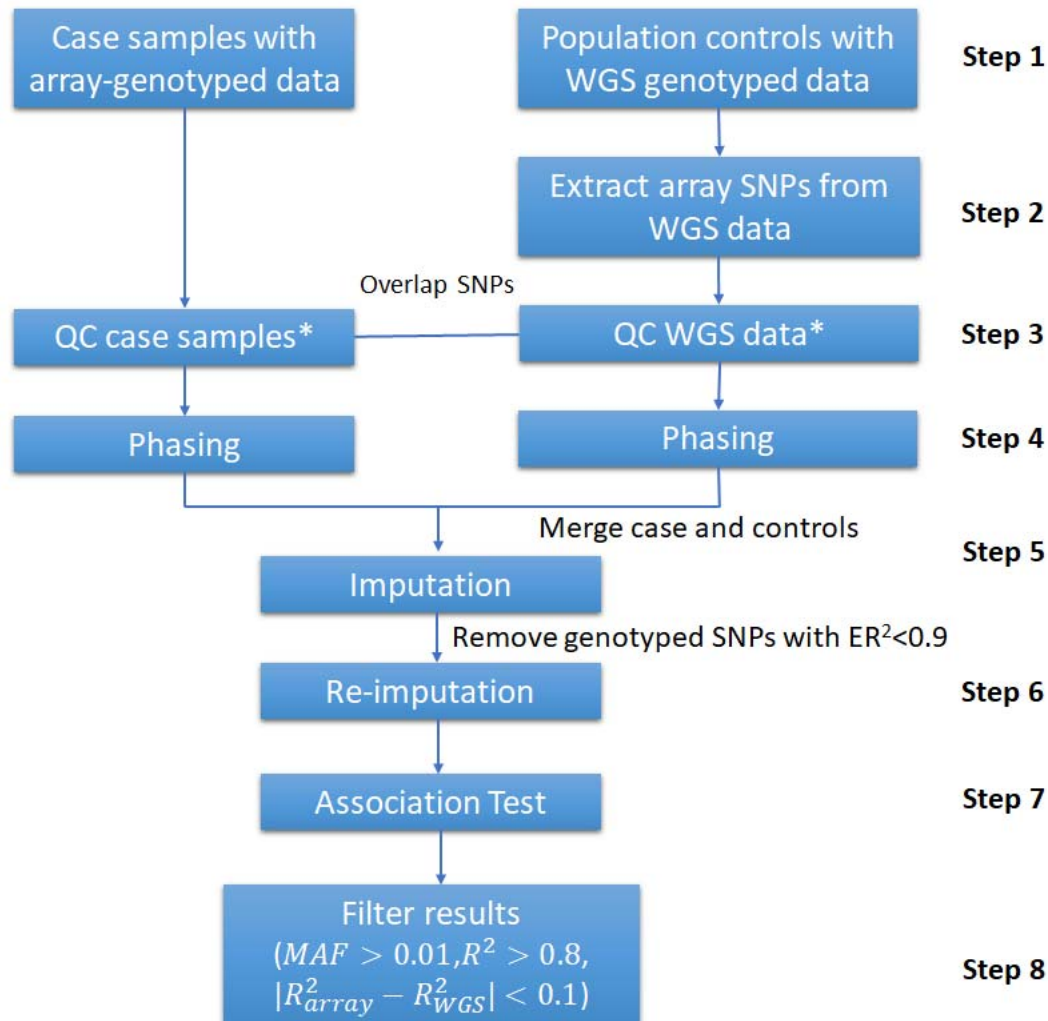
385 \* The question mark “?” means the SNP is missing from the first analysis, and it may result in reduced power in the final meta-analysis.

386

387

388 **Figures**

389



391

392 **Figure 1:** Overview of the protocol to use whole-genome sequencing (WGS) data as public control in

393 GWAS. \*The quality control (QC) of the case and public control data is conducted independently

394 according to the steps outlined in the methods.

395

A

Technical Comparison		
Array data	WGS data	Results
COPDGene EA (N=6,501)	COPDGene EA (N=6,501)	Tech EA (Suppl Figure 5a)
COPDGene AA (N=3,235)	COPDGene AA (N=3,235)	Tech AA (Suppl Figure 5b)

B

Control for Type 1 Error		
Array Data	WGS Data	Results
COGENE EA (N=1,961)	COPDGene EA1* (N=3,251)	COGENE EA GWAS (Suppl Figure 6a)
COPDGene EA2* (N=3,251)	ECLIPSE EA (N=1,461)	ECLIPSE EA GWAS (Suppl Figure 6b)
COGENE AA (N=712)	COPDGene AA (N=1,710)	COGENE AA GWAS (Suppl Figure 6c)

Integration via  
Meta-Analysis  
(Figure 3)

C

Replication of Known GWAS Hits		
Array data	WGS data	Results
COGENE EA (N=1,961)	COPDGene EA COPD cases (N=2,736)	COGENE EA GWAS (Suppl Figure 7a)
ECLIPSE EA COPD case (N=1,764)	COPDGene EA COPD controls (N=2,475)	ECLIPSE EA GWAS (Suppl Figure 7b)
COGENE AA (N=712)	COPDGene AA COPD cases (N=813)	COGENE AA GWAS (Suppl Figure c)

Integration via  
Meta-Analysis  
(Figure 4)

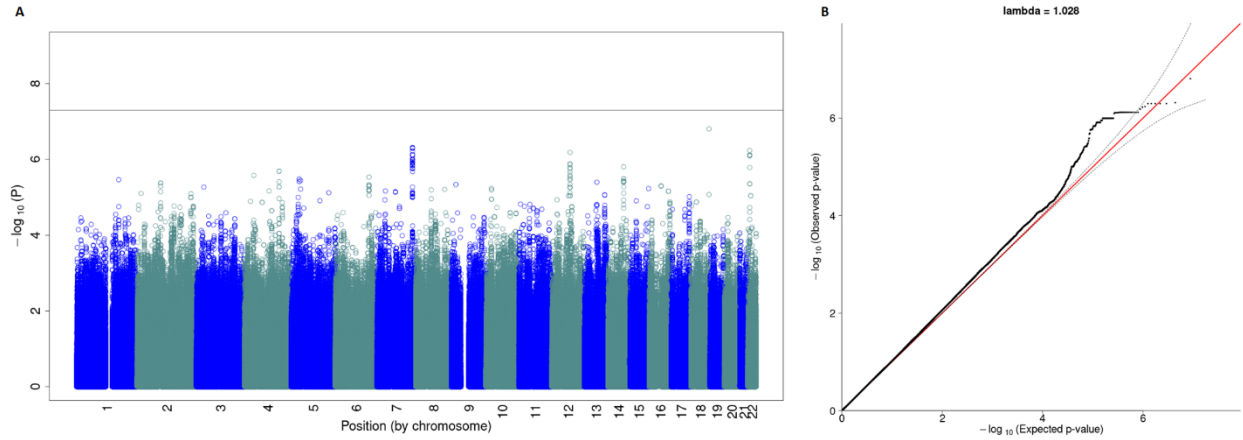
395

396 **Figure 2.** Evaluation design for (a) technical comparison, (b) type-I error assessment, and (c)

397 known GWAS hits. \*The samples with European ancestry in COPDGene were evenly divided to two

398 subsets of samples. EA1 includes all COPD cases and some COPD controls to match the COPD prevalence

399 in ECLIPSE. EA2 has all the rest COPD free samples.

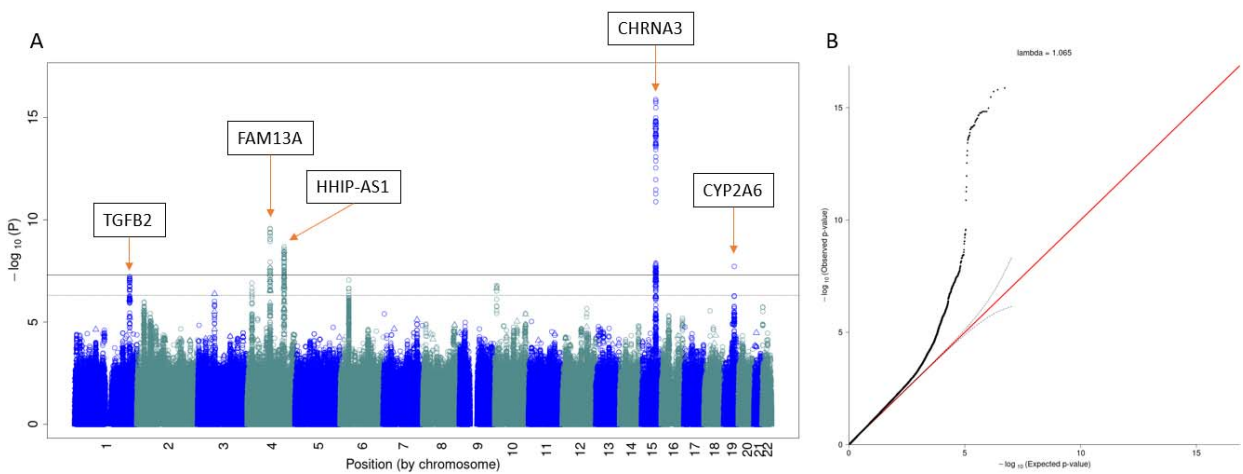


400

401 **Figure 3.** Meta-analysis results from evaluation for type-I error. The Manhattan plot (a) shows

402 the expected no signal, while the QQ-plot (b) shows no inflation.

403



404

405 **Figure 4.** Meta-analysis results for replication of GWAS hits for COPD. The Manhattan plot (a)

406 shows the replicated signals, while the QQ-plot (b) shows inflation due to the true signal.



## 407 Data availability

408 The individual-level genotype and phenotype data used are all available through dbGap. The dbGap  
409 study accession number for COGEND is phs000404, for COPDGene are phs000179 (parent study with  
410 array genotype data) and phs000951 (WGS data generated by TOPMed), and for ECLIPSE are phs001252  
411 (parent study with array genotype data) and phs001472 (WGS data generated by TOPMed).

## 412 Code availability

413 The codes to run the protocol can be found at <https://github.com/RTIInternational/GAWMerge>.

## 414 References

- 415 1. Luca, D. *et al.* On the use of general control samples for genome-wide association studies:  
416 genetic matching highlights causal variants. *Am J Hum Genet* **82**, 453-63 (2008).
- 417 2. Cooper, J.D. *et al.* Meta-analysis of genome-wide association study data identifies additional  
418 type 1 diabetes risk loci. *Nat Genet* **40**, 1399-401 (2008).
- 419 3. Rao, D.C. An overview of the genetic dissection of complex traits. *Adv Genet* **60**, 3-34 (2008).
- 420 4. Todd, J.A. *et al.* Robust associations of four new chromosome regions from genome-wide  
421 analyses of type 1 diabetes. *Nat Genet* **39**, 857-64 (2007).
- 422 5. Johnson, E.O. *et al.* KAT2B polymorphism identified for drug abuse in African Americans with  
423 regulatory links to drug abuse pathways in human prefrontal cortex. *Addiction biology* **21**, 1217-  
424 1232 (2016).
- 425 6. Ho, L.A. & Lange, E.M. Using public control genotype data to increase power and decrease cost  
426 of case-control genetic association studies. *Human Genetics* **128**, 597-608 (2010).
- 427 7. Mukherjee, S. *et al.* Including Additional Controls from Public Databases Improves the Power of  
428 a Genome-Wide Association Study. *Human Heredity* **72**, 21-34 (2011).
- 429 8. Zhuang, J.J. *et al.* Optimizing the power of genome-wide association studies by using publicly  
430 available reference samples to expand the control group. *Genet Epidemiol* **34**, 319-26 (2010).
- 431 9. Johnson, E.O. *et al.* Imputation across genotyping arrays for genome-wide association studies:  
432 assessment of bias and a correction strategy. *Hum Genet* **132**, 509-22 (2013).
- 433 10. Lindstrom, S. *et al.* A comprehensive survey of genetic variation in 20,691 subjects from four  
434 large cohorts. *PLoS One* **12**, e0173997 (2017).
- 435 11. Kowalski, M.H. *et al.* Use of >100,000 NHLBI Trans-Omics for Precision Medicine (TOPMed)  
436 Consortium whole genome sequences improves imputation quality and detection of rare variant  
437 associations in admixed African and Hispanic/Latino populations. *PLoS Genet* **15**, e1008500  
438 (2019).
- 439 12. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data. *Nature* **562**,  
440 203-209 (2018).
- 441 13. Wall, J.D. *et al.* The GenomeAsia 100K Project enables genetic discoveries across Asia. *Nature*  
442 **576**, 106-111 (2019).
- 443 14. Das, S. *et al.* Next-generation genotype imputation service and methods. *Nat Genet* **48**, 1284-  
444 1287 (2016).
- 445 15. Bierut, L.J. *et al.* Novel genes identified in a high-density genome wide association study for  
446 nicotine dependence. *Hum Mol Genet* **16**, 24-35 (2007).

- 447 16. Saccone, S.F. *et al.* Cholinergic nicotinic receptor genes implicated in a nicotine dependence  
448 association study targeting 348 candidate genes with 3713 SNPs. *Hum Mol Genet* **16**, 36-49  
449 (2007).
- 450 17. Regan, E.A. *et al.* Genetic epidemiology of COPD (COPDGene) study design. *COPD* **7**, 32-43  
451 (2010).
- 452 18. Vestbo, J. *et al.* Evaluation of COPD Longitudinally to Identify Predictive Surrogate End-points  
453 (ECLIPSE). *Eur Respir J* **31**, 869-73 (2008).
- 454 19. Cho, M.H. *et al.* Risk loci for chronic obstructive pulmonary disease: a genome-wide association  
455 study and meta-analysis. *Lancet Respir Med* **2**, 214-25 (2014).
- 456 20. Verlouw, J.A.M. *et al.* A comparison of genotyping arrays. *European Journal of Human Genetics*  
457 (2021).
- 458 21. Danilov, K.A., Nikogosov, D.A., Musienko, S.V. & Baranova, A.V. A comparison of BeadChip and  
459 WGS genotyping outputs using partial validation by sanger sequencing. *BMC Genomics* **21**, 528  
460 (2020).
- 461 22. Taliun, D. *et al.* Sequencing of 53,831 diverse genomes from the NHLBI TOPMed Program.  
462 *Nature* **590**, 290-299 (2021).
- 463 23. Cho, M.H. *et al.* A genome-wide association study of COPD identifies a susceptibility locus on  
464 chromosome 19q13. *Hum Mol Genet* **21**, 947-57 (2012).
- 465 24. Hobbs, B.D. *et al.* Genetic loci associated with chronic obstructive pulmonary disease overlap  
466 with loci for lung function and pulmonary fibrosis. *Nat Genet* **49**, 426-432 (2017).
- 467 25. Wojcik, G.L. *et al.* Genetic analyses of diverse populations improves discovery for complex traits.  
468 *Nature* **570**, 514-518 (2019).
- 469 26. Abel, H.J. & Duncavage, E.J. Detection of structural DNA variation from next generation  
470 sequencing data: a review of informatic approaches. *Cancer Genet* **206**, 432-40 (2013).
- 471 27. Gudbjartsson, D.F. *et al.* Large-scale whole-genome sequencing of the Icelandic population. *Nat*  
472 *Genet* **47**, 435-44 (2015).
- 473 28. Purcell, S. *et al.* PLINK: a tool set for whole-genome association and population-based linkage  
474 analyses. *Am J Hum Genet* **81**, 559-75 (2007).
- 475 29. Danecek, P. *et al.* Twelve years of SAMtools and BCftools. *Gigascience* **10**(2021).
- 476 30. Delaneau, O., Marchini, J. & Zagury, J.F. A linear complexity phasing method for thousands of  
477 genomes. *Nat Methods* **9**, 179-81 (2011).
- 478 31. Delaneau, O., Marchini, J., Genomes Project, C. & Genomes Project, C. Integrating sequence and  
479 array data to create an improved 1000 Genomes Project haplotype reference panel. *Nat*  
480 *Commun* **5**, 3934 (2014).
- 481 32. Zhan, X., Hu, Y., Li, B., Abecasis, G.R. & Liu, D.J. RVTESTS: an efficient and comprehensive tool for  
482 rare variant association analysis using sequence data. *Bioinformatics* **32**, 1423-6 (2016).
- 483 33. National Heart, L., and Blood Institute, National Institutes of Health, U.S. Department of Health  
484 and Human Services. The NHLBI BioData Catalyst. *Zenodo* (2020).

485

## 486 Acknowledgements

487 Primary support for developing GAWMerge, conducting analyses, and preparing the

488 manuscript was provided by National Institute on Drug Abuse grants to Dr. Johnson: R01

489 DA044014 (PI: Johnson); R01 DA043980 (M-PIs: Scacheri, Johnson, Akbarian); R01

490 DA051908 (M-PIs: Johnson and Jacobson)

491 Support for this work was provided by the National Institutes of Health, National Heart,  
492 Lung, and Blood Institute, through the BioData Catalyst program (award 1OT3HL142479-01,  
493 1OT3HL142478-01, 1OT3HL142481-01, 1OT3HL142480-01, 1OT3HL147154-01). Any opinions  
494 expressed in this document are those of the author(s) and do not necessarily reflect the views of  
495 NHLBI, individual BioData Catalyst team members, or affiliated organizations and institutions.

496 Molecular data for the Trans-Omics in Precision Medicine (TOPMed) program was  
497 supported by the National Heart, Lung and Blood Institute (NHLBI). Genome sequencing for  
498 NHLBI TOPMed: COPDGene (phs000179.v6.p2) was performed at NWGC (3R01HL089856-  
499 08S1, HHSN268201600032I, and HHSN268201600032I), and Broad Genomics  
500 (HHSN268201500014C and HHSN268201500014C). Genome sequencing for NHLBI TOPMed:  
501 ECLIPSE (phs001252.v1.p1) was performed at MDI (HHSN268201600037I). Core support  
502 including centralized genomic read mapping and genotype calling, along with variant quality  
503 metrics and filtering were provided by the TOPMed Informatics Research Center (3R01HL-  
504 117626-02S1; contract HHSN268201800002I). Core support including phenotype  
505 harmonization, data management, sample-identity QC, and general program coordination were  
506 provided by the TOPMed Data Coordinating Center (R01HL-120393; U01HL-120393; contract  
507 HHSN268201800001I). We gratefully acknowledge the studies and participants who provided  
508 biological samples and data for TOPMed.

509 The COPDGene project described was supported by Award Number U01 HL089897 and  
510 Award Number U01 HL089856 from the National Heart, Lung, and Blood Institute. The content  
511 is solely the responsibility of the authors and does not necessarily represent the official views of  
512 the National Heart, Lung, and Blood Institute or the National Institutes of Health. The  
513 COPDGene project is also supported by the COPD Foundation through contributions made to  
514 an Industry Advisory Board that has included AstraZeneca, Bayer Pharmaceuticals, Boehringer-

515 Ingelheim, Genentech, GlaxoSmithKline, Novartis, Pfizer, and Sunovion. A full listing of  
516 COPDGene investigators can be found at: <http://www.copdgene.org/directory>.

517       The ECLIPSE study (NCT00292552) was sponsored by GlaxoSmithKline. The ECLIPSE  
518 investigators included: Bulgaria: Y. Ivanov, Pleven; K. Kostov, Sofia. Canada: J. Bourbeau,  
519 Montreal; M. Fitzgerald, Vancouver, BC; P. Hernandez, Halifax, NS; K. Killian, Hamilton, ON; R.  
520 Levy, Vancouver, BC; F. Maltais, Montreal; D. O'Donnell, Kingston, ON. Czech Republic: J.  
521 Krepelka, Prague. Denmark: J. Vestbo, Hvidovre. The Netherlands: E. Wouters, Horn-  
522 Maastricht. New Zealand: D. Quinn, Wellington. Norway: P. Bakke, Bergen. Slovenia: M.  
523 Kosnik, Golnik. Spain: A. Agusti, J. Sauleda, P. de Mallorca. Ukraine: Y. Feschenko, V.  
524 Gavrisyuk, L. Yashina, Kiev; N. Monogarova, Donetsk. United Kingdom: P. Calverley, Liverpool;  
525 D. Lomas, Cambridge; W. MacNee, Edinburgh; D. Singh, Manchester; J. Wedzicha, London.  
526 United States: A. Anzueto, San Antonio, TX; S. Braman, Providence, RI; R. Casaburi, Torrance  
527 CA; B. Celli, Boston; G. Giessel, Richmond, VA; M. Gotfried, Phoenix, AZ; G. Greenwald,  
528 Rancho Mirage, CA; N. Hanania, Houston; D. Mahler, Lebanon, NH; B. Make, Denver; S.  
529 Rennard, Omaha, NE; C. Rochester, New Haven, CT; P. Scanlon, Rochester, MN; D. Schuller,  
530 Omaha, NE; F. Sciruba, Pittsburgh; A. Sharafkhaneh, Houston; T. Siler, St. Charles, MO; E.  
531 Silverman, Boston; A. Wanner, Miami; R. Wise, Baltimore; R. ZuWallack, Hartford, CT.  
532 ECLIPSE Steering Committee: H. Coxson (Canada), C. Crim (GlaxoSmithKline, USA), L.  
533 Edwards (GlaxoSmithKline, USA), D. Lomas (UK), W. MacNee (UK), E. Silverman (USA), R.  
534 Tal-Singer (Co-chair, GlaxoSmithKline, USA), J. Vestbo (Co-chair, Denmark), J. Yates  
535 (GlaxoSmithKline, USA). ECLIPSE Scientific Committee: A. Agusti (Spain), P. Calverley (UK),  
536 B. Celli (USA), C. Crim (GlaxoSmithKline, USA), B. Miller (GlaxoSmithKline, USA), W. MacNee  
537 (Chair, UK), S. Rennard (USA), R. Tal-Singer (GlaxoSmithKline, USA), E. Wouters (The  
538 Netherlands), J. Yates (GlaxoSmithKline, USA).

539       COGEND was supported by grants from the National Cancer Institute (NCI; grant  
540 number P01 CA089392, PI: Laura Bierut) and NIDA (R01 DA036583 and R01 DA025888, PI:

541 Laura Bierut), both of the National Institutes of Health (NIH). Genotype data are available via  
542 dbGaP as part of the “Genetic Architecture of Smoking and Smoking Cessation” (accession  
543 number phs000404.v1.p1) and “Study of Addiction: Genetics and Environment (SAGE)”  
544 (accession number phs000092.v1.p1). Funding support for genotyping, which was performed at  
545 CIDR, was provided by 1 X01 HG005274-01 and by the NIH Genes, Environment and Health  
546 Initiative [GEI] (U01 HG004422). CIDR is fully funded through a federal contract from the NIH to  
547 The Johns Hopkins University, contract number HHSN268200782096C. Assistance with  
548 genotype cleaning, as well as with general study coordination, was provided by the GENEVA  
549 Coordinating Center (U01 HG004446).

550

### 551 **Author contributions**

552 RM and FF co-led the development of GAWMerge, all analyses, and the writing of the  
553 manuscript. NG, DBH, GPP, and EOJ conceptualized and helped develop GAWMerge. MHC,  
554 JEH, LJB, SML, KY, and EKS provided valuable cohort and TOPMed datasets expertise used in  
555 the development of GAWMerge. All authors edited the manuscript.

556

557

### 558 **Competing interests**

559 Edwin Silverman has received institutional grant support from GlaxoSmithKline and Bayer.  
560 Michael H. Cho has received grant support from GSK and Bayer, and consulting or speaking  
561 fees from Illumina, Genentech, and AstraZeneca. All other authors have no competing interests.

562

563

564

565

566