

1 **Title: Fine-mapping the *CYP2A6* regional association with nicotine metabolism among**
2 **African American smokers**

3

4 Running Title: *CYP2A6* predicts nicotine metabolism across populations

5

6 Jennie G. Pouget*^{1,2}, Haidy Giratallah^{1,3}, Alec W.R. Langlois^{1,3}, Ahmed El-Boraie^{1,3}, Caryn
7 Lerman⁴, Jo Knight⁵, Lisa Sanderson Cox⁶, Nikki L. Nollen⁶, Jasjit S. Ahluwalia⁷, Christian
8 Benner⁸, Meghan J. Chenoweth^{1,2,3}, Rachel F. Tyndale^{1,2,3}

9

10 ¹Campbell Family Mental Health Research Institute, Centre for Addiction and Mental Health

11 ²Department of Psychiatry, University of Toronto

12 ³Department of Pharmacology & Toxicology, University of Toronto

13 ⁴Department of Psychiatry, University of Pennsylvania

14 ⁵Data Science Institute and Medical School, Lancaster University

15 ⁶Department of Population Health, University of Kansas School of Medicine, Kansas City,
16 Kansas, USA

17 ⁷Departments of Behavioral and Social Sciences and Medicine, Brown University, Providence,
18 Rhode Island, USA

19 ⁸Institute for Molecular Medicine Finland (FIMM), University of Helsinki, Helsinki, Finland

20

21 Corresponding author: Rachel F Tyndale, 416 978-6374, r.tyndale@utoronto.ca

22

23

24 **Abstract**

25 The nicotine metabolite ratio (NMR; 3’hydroxycotinine/cotinine) is a stable biomarker for
26 CYP2A6 enzyme activity and nicotine clearance, with demonstrated clinical utility in
27 personalizing smoking cessation treatment. Common genetic variation in the *CYP2A6* region is
28 strongly associated with NMR in smokers. Here, we investigated this regional association in
29 more detail. We evaluated the association of *CYP2A6* single-nucleotide polymorphisms (SNPs)
30 and * alleles with NMR among African American smokers (N=953) from two clinical trials of
31 smoking cessation. Stepwise conditional analysis and Bayesian fine-mapping were undertaken.
32 Putative causal variants were incorporated into an existing African ancestry-specific genetic risk
33 score (GRS) for NMR, and the performance of the updated GRS was evaluated in both African
34 American (n=953) and European ancestry smokers (n=933) from these clinical trials. Five
35 independent associations with NMR in the *CYP2A6* region were identified using stepwise
36 conditional analysis, including the deletion variant *CYP2A6**4 (beta=-0.90, p=1.55x10⁻¹¹). Six
37 putative causal variants were identified using Bayesian fine-mapping (posterior probability,
38 PP=0.67), with the top causal configuration including *CYP2A6**4, rs116670633, *CYP2A6**9,
39 rs28399451, rs8192720, and rs10853742 (PP=0.89). Incorporating these putative causal variants
40 into an existing ancestry-specific GRS resulted in comparable prediction of NMR within African
41 American smokers, and improved trans-ancestry portability of the GRS to European smokers.
42 Our findings suggest that both * alleles and SNPs underlie the association of the *CYP2A6* region
43 with NMR among African American smokers, identify a shortlist of variants that may causally
44 influence nicotine clearance, and suggest that portability of GRSs across populations can be
45 improved through inclusion of putative causal variants.

46

47 **Introduction**

48

49 Tobacco use remains the leading cause of preventable death and disease in North America.¹

50 Nicotine (the primary addictive agent in tobacco)² is metabolized to cotinine primarily by the

51 liver enzyme CYP2A6, and then to 3'hydroxycotinine exclusively by hepatic CYP2A6.^{3,4} The

52 **nicotine metabolite ratio (NMR; 3'hydroxycotinine/cotinine)** is a stable biomarker for nicotine

53 metabolism by CYP2A6 in smokers.^{5,6} Individual differences in NMR predict total nicotine

54 clearance, and thus smoking behaviours (including cessation) as well as health outcomes. In

55 particular, higher NMR (i.e. faster nicotine inactivation and CYP2A6 activity) is associated with

56 greater nicotine dependence, cigarette consumption, and lung cancer risk along with lower

57 cessation.^{7,8} Furthermore, NMR has translational potential in personalizing cessation treatment

58 given that smokers with higher NMR show greater benefit from treatment with varenicline

59 (compared to nicotine replacement therapy).^{9,10}

60

61 The NMR can only be reliably measured in current, regular smokers. This limits its use as a

62 biomarker in longitudinal studies of smoking initiation or smoking-related disease risk in

63 occasional/non-smokers, and limits the potential clinical utility of using NMR to guide

64 personalized counselling on smoking-related risks to promote prevention efforts and behavioural

65 change. However, because NMR is highly heritable ($h^2=60-80\%$ ^{11,12}), an individual's NMR

66 could potentially be estimated using their genetic information regardless of their current smoking

67 status (i.e. using a genetic risk score that predicts NMR). To achieve this, large-scale genetic

68 studies of NMR are required to robustly identify the underlying genetic risk variants.

69

70 To date, most genetic studies of NMR have been undertaken in European ancestry smokers, and
71 the genetic architecture of NMR in non-European smokers remains only partially understood,
72 contributing to potential health disparities.¹³ In European smokers, the largest GWAS of NMR
73 conducted (n=5,185) identified a strong genome-wide association near *CYP2A6* on chromosome
74 19, and a second association near *TMPRSS11E* on chromosome 4.¹⁴ The *CYP2A6* association
75 pattern in European smokers was complex, with six independent variants identified in
76 conditional analysis and a top causal configuration including 13 variants identified in Bayesian
77 fine-mapping.¹⁴ To our knowledge we have conducted the largest GWAS of NMR in African
78 American smokers to date (n=954), finding a single genome-wide association near *CYP2A6*. The
79 association pattern in African American smokers was unique compared to that observed in
80 Europeans,¹⁵ with 58 of the 96 genome-wide significant hits not reaching genome-wide
81 significant in Europeans and a different lead variant (rs12459249) that was not in high linkage
82 disequilibrium (LD) with the top variant in Europeans ($r^2 < 0.6$).¹⁷

83
84 While GWAS provide comprehensive coverage of single nucleotide polymorphisms (SNPs),
85 there are several well characterized *CYP2A6* * alleles with known functional effects on *CYP2A6*
86 activity that are not well captured using standard GWAS approaches.¹⁸ Incorporating both
87 *CYP2A6* * alleles and common genetic variants identified by GWAS, we previously developed
88 ancestry-specific genetic risk scores (GRSs) to estimate an individual's NMR from their genetic
89 information.^{19,20} These GRSs explained 33.8% and 32.4% of variance in NMR in European¹⁹ and
90 African²⁰ ancestry populations, respectively, and showed reasonable prediction of slow vs.
91 normal nicotine metabolizer status in these populations (AUC=0.78 and 0.73, respectively).^{19,20}
92 As has been previously described for GRS more broadly,¹³ given differences in LD structure

93 across ancestral populations these ancestry-specific GRSs showed poor portability across
94 populations, with the European and African ancestry GRSs explaining only 18-20% of variance
95 in NMR in the alternate population.²⁰ Additionally, Bloom *et al.* developed an ancestry-specific
96 GRS for a different nicotine metabolism measure (D₂-cotinine:[D₂-nicotine+D₂-cotinine]) in
97 Europeans using * alleles and other variants from the literature.²¹ Development of a universal
98 GRS using multi-ancestry cohorts is another promising approach, with Baurley *et al.* reporting
99 similar predictive performance across African, Asian, and European ancestry smokers using
100 machine learning algorithms to predict NMR based on age, sex, ancestry, BMI, and a set of 263
101 SNPs prioritized from GWAS (of which 198 were located in the *CYP2A6* region).²²

102

103 In summary, previous large-scale efforts have been undertaken to fine-map the *CYP2A6* regional
104 association with NMR in European ancestry smokers.¹⁴ However, to our knowledge there has
105 been no previous study fine-mapping the genome-wide *CYP2A6* association in African ancestry
106 smokers. Given growing interest in developing genetic tools to assist with smoking counseling
107 and cessation, in the current study we address this knowledge gap and the potential health
108 disparities it creates. Building on our previous studies in a group of African Americans
109 participating in two large smoking cessation trials (**Figure S1**), here we investigated the *CYP2A6*
110 association with NMR in more detail using an updated conditional analysis and new Bayesian
111 fine-mapping approach to analyze both SNPs and * alleles (including structural variants) in the
112 region. We also evaluated whether incorporating the putative causal variants identified by fine-
113 mapping improved an existing ancestry-specific GRS to genetically predict NMR in African
114 American populations, and the portability of this GRS to predict NMR in those of European
115 ancestry.

116

117 **Materials and Methods**

118 **Participants**

119 Our study sample comprised African and European ancestry smokers from two clinical trials of
120 cessation: Pharmacogenetics of Nicotine Addiction Treatment 2 (PNAT-2; NCT01314001)¹⁰ and
121 Kick-it-at-Swope 3 (KIS-3; NCT00666978).²³ The clinical trial protocols were approved by
122 institutional review boards at all participating sites and the University of Toronto.

123

124 Study design of both PNAT-2 and KIS-3 have been described in detail elsewhere.^{10,23} Briefly,
125 PNAT-2 randomized eligible adult smokers (aged 18-65 years, smoking ≥ 10 cigarettes/day) by
126 NMR group (normal metabolizers vs. slow metabolizers) to treatment with placebo, nicotine
127 patch, or varenicline for smoking cessation; all three treatment arms received behavioural
128 counselling.¹⁰ Approximately 37% of the total PNAT-2 sample were African ancestry
129 (genetically determined based on comparison of genome-wide data to population reference
130 panels as previously described,²⁰ see **Quality Control** below for further details), and were
131 included in the primary analyses here (n=506, **Table 1**). We conducted additional analyses
132 evaluating the portability of GRSs developed to predict NMR in African populations to the
133 subset of PNAT-2 participants that were European ancestry (genetically determined as
134 previously described,¹⁹ n=933).

135

136 KIS-3 randomized eligible adult light smokers (aged ≥ 18 years, smoking ≤ 10 cigarettes/day)
137 who self-identified as African American to treatment with bupropion or placebo for smoking
138 cessation; both treatment arms received health education counselling.²³ Recruitment for KIS-3

139 was from a community-based clinic in Kansas, MO.²³ Participants who were African ancestry
140 (genetically determined, as previously described,²⁰ n=458) were included in the primary analyses
141 (**Table 1**).

142

143 **Outcome Measure**

144 *Nicotine metabolite ratio (NMR, 3’hydroxycotinine/cotinine ratio)*

145 We measured NMR as a continuous variable by determining the ratio of
146 3’hydroxycotinine/cotinine concentrations in blood samples collected at the time of clinical trial
147 enrollment, when participants were smoking regularly. Cotinine and 3’hydroxycotinine
148 concentrations were determined using liquid chromatography-tandem mass spectrometry, as
149 previously described.²⁴

150

151 **Genetic Data Collection**

152 *Genotyping*

153 To capture common SNPs, we conducted genome-wide genotyping using the Illumina
154 HumanOmniExpressExome-8 v1.2 array (Illumina, San Diego, CA, USA) at the Centre for
155 Applied Genomics, Hospital for Sick Children (Toronto, ON, Canada). We also included a
156 previously described custom iSelect® add-on, capturing an additional 2,688 variants associated
157 with nicotine metabolism and/or smoking behaviours for richer coverage of regions of interest
158 including *CYP2ABFGST* (chromosome 19), *CHRNA5-A3-B4* (chromosome 15), *OCT2*
159 (chromosome 6), and *UGT2B* (chromosome 4).¹⁵

160

161 We directly genotyped the following 12 *CYP2A6* * alleles: *CYP2A6**46 (formerly *CYP2A6**1B),
162 *CYP2A6**1x2, *CYP2A6**4, *CYP2A6**9, *CYP2A6**12, *CYP2A6**17, *CYP2A6**20, *CYP2A6**23,
163 *CYP2A6**25/*26/*27 (all tagged by rs28399440), *CYP2A6**28, *CYP2A6**31, *CYP2A6**35 as
164 previously described.^{19,20} These *CYP2A6* * alleles have demonstrated functional effects on
165 *CYP2A6* activity, and include structural variants (*CYP2A6* gene deletions and duplications) as
166 well as amino acid changes (see **Table S2** for details). Individuals with structural variants
167 (*CYP2A6**1x2, *CYP2A6**4, *CYP2A6**12, *CYP2A6**34, and *CYP2A6**53) were re-genotyped
168 using an approach with improved accuracy, as previously described.²⁵

169

170 ***Quality Control***

171 We performed quality control for samples and raw genotype data using PLINK,²⁶ following
172 standard protocols as previously described.¹⁵ Individuals with discrepant sex, genotype call
173 rate<0.98, heterozygosity rate>3 SDs from sample mean, substantial cryptic relatedness
174 (PI_HAT>0.185), or substantial non-African admixture (determined by visual inspection of
175 multidimensional scaling (MDS) plots) were excluded. Self-reported African American ancestry
176 was highly concordant with genetically determined ancestry in our sample (>95% concordance
177 rate).¹⁵ Variants with call rate<0.98, minor allele frequency (MAF)<0.01, or Hardy-Weinberg
178 equilibrium (HWE) p-value<1x10⁻⁶ were excluded.

179

180 ***Imputation***

181 We imputed chromosome 19 using the Michigan Imputation Server, which utilizes Minimac4.²⁷
182 Accurately sequencing the *CYP2A6* region is challenging due to extensive variability, regions of
183 high homology (i.e. including the pseudogene *CYP2A7*), and complex structural variation,¹⁸ poor

184 sequencing quality in this region reduces the quality of imputed genotype calls made using
185 standard reference panels. Therefore, we compared the results of imputation using two different
186 cosmopolitan reference panels: the TOPMed Version R2 reference panel (N=97,256 with ~30%
187 African ancestry from African, African Caribbean, or African American populations),²⁸ and the
188 1000 Genomes Phase 3 reference panel (N=2,504 with ~25% of African ancestry from the
189 following populations: Esan in Nigeria (ESN), Gambian in Western Division, Mandinka (GWD),
190 Luhya in Webuye, Kenya (LWK), Mende in Sierra Leone (MSL), Yoruba in Ibadan, Nigeria
191 (YRI), African Caribbean in Barbados (ACB), people with African ancestry in Southwest USA
192 (ASW)).²⁹ The TOPMED imputation was performed with pre-phasing of haplotypes using Eagle
193 v2.4 and human genome build hg38.³⁰ The 1000 Genomes Phase 3 imputation was performed
194 with pre-phasing of haplotypes using ShapeIT v2.r79034³¹ and human genome build hg37, as
195 previously described.³²

196

197 Post-imputation quality control was performed using PLINK²⁶ to exclude duplicate and multi-
198 allelic variants, as well as variants with poor imputation quality (INFO<0.6) or HWE p-
199 value<1x10⁻⁶. We then compared the density of coverage and imputation quality across the two
200 imputation methods.

201

202 **Statistical Analyses**

203 *Association Testing*

204 All statistical analyses were done using R Statistical Software unless otherwise specified.³³ We
205 used a mega-analytic approach, pooling data from both clinical trials (PNAT-2 and KIS-3) for all
206 analyses unless otherwise specified.

207
208 Based on LD patterns in our sample, and in keeping with prior *CYP2A6* fine-mapping efforts in
209 European ancestry smokers,¹⁴ we included variants within 5 Mb of *CYP2A6* in our analyses
210 (chromosome 19:38,000,000-43,000,000bp; Genome Reference Consortium Human Build 38,
211 hg38). We evaluated the association of these variants in the *CYP2A6* region with NMR. Given
212 the non-normal distribution of NMR in our sample, we applied rank-based inverse normal
213 transformation using the R package RNOmni³⁴ and used these transformed NMR values for all
214 analyses unless otherwise specified (**Figure S2**).

215
216 Association testing was done in SNPTEST v2.5.2³⁵ using linear regression to test the association
217 of imputed genotype dosages with normalized NMR using an additive genotypic model with
218 adjustment for age, sex, body mass index (BMI), and two ancestry-informative dimensions to
219 account for population substructure as covariates.

220
221 ***Stepwise Conditional Analysis***
222 To identify the number of independent associations in the *CYP2A6* region, we completed
223 stepwise conditional analysis in SNPTEST v2.5.2³⁵ by including genotype dosages for the top
224 variant as an additional covariate in the base model described above (effectively conditioning on
225 additive effects of the top variant), and repeating this procedure until no further association
226 signals reached genome-wide significance ($p < 5 \times 10^{-8}$). Regional association plots were
227 constructed using LocusZoom, with LD information from the 1000 Genomes Phase 3 African
228 populations reference panel.³⁶

229

230 ***Bayesian Fine-mapping***

231 To identify potentially causal variants in the *CYP2A6* region, we used FINEMAP v1.4 specifying
232 a maximum of 20 potential causal variants.³⁷ FINEMAP performs Bayesian fine-mapping using
233 a shotgun stochastic search method to identify the most likely causal configuration of variants,
234 given association summary statistics and local LD patterns.³⁷ We also performed exploratory
235 functionally informed fine-mapping in FINEMAP³⁷ by assigning a higher prior probability to
236 *CYP2A6* * alleles (prior probability=0.70 for these variants being causal) compared to non-*
237 allele variants (prior probability=0.50). Summary statistics were obtained as described above
238 using SNPTEST v2.5.2,³⁵ and the SNP correlation matrix was computed from genotype dosages
239 in our sample using LDstore v2.0.³⁸ Regional association plots were constructed using R.³³

240

241 ***Variant Annotation***

242 To annotate variants identified in our analyses we used RegulomeDB,³⁹ a publicly available
243 database that estimates a variant's likelihood of having a regulatory function using a probability
244 score that range from 0 to 1 (with 1 being most likely to be a regulatory variant). The probability
245 score is constructed based on a machine learning model integrating functional genomic data
246 including ChIP-seq signal, DNase-seq signal, information content change, and DeepSEA
247 scores.³⁹

248

249 We also evaluated whether variants were known to influence expression of genes encoding
250 functional proteins using publicly available expression quantitative trait loci (eQTL) data from
251 the Genotype-Tissue Expression (GTEx) Project.⁴⁰ The GTEx Project eQTL analysis was based
252 on whole genome sequencing and RNA-seq data collected from 838 donors (~13% African

253 ancestry) across 49 tissues. Given the potential misidentification of *CYP2A6* transcripts as
254 pseudogene *CYP2A7* due to high sequence homology, we considered eQTL data for pseudogene
255 *CYP2A7* along with all other protein-coding genes. The data used for the analyses described in
256 this manuscript were obtained from the GTEx Portal on 12/04/2024.

257

258 ***Incorporation of Putative Causal Variants into an Existing Genetic Risk Score (GRS) for***
259 ***NMR***

260 To investigate whether Bayesian fine-mapping improved the predictive power of genetically
261 determined NMR in African American smokers, we compared our previously described GRS for
262 this ancestral population²⁰ (referred to here as the **original GRS**) to GRSs including putative
263 causal variants identified by fine-mapping in the current study. The original GRS included eight
264 *CYP2A6* * alleles (*1x2, *4, *9, *12, *17, *20, *25/*26/*27, *35) and three LD-independent
265 genome-wide significant SNPs (rs12459249, rs111645190, rs185430475) identified in an earlier
266 conditional analysis of the *CYP2A6* region.¹⁵ The initial GRS estimation was constructed using
267 mentholated cigarette use as an additional covariate, and explained 32.4% of the variance in log-
268 NMR.²⁰ We elected to not adjust for menthol in the current study in order to maximize sample
269 size (10% of participants were missing menthol data) and because menthol adjustment did not
270 appreciably alter SNP effects on NMR.³² For harmonization with data used in the current study,
271 we therefore recalculated the weights for all variants in the original GRS using the analytic
272 approach described below (without adjustment for mentholated cigarette use), and with *CYP2A6*
273 * allele genotypes obtained using a more recent genotyping approach with improved accuracy.²⁵

274

275 The **updated GRS** included all eight *CYP2A6* * alleles from the original GRS and the six LD-
 276 independent putative causal variants identified by FINEMAP as the lead variant in their
 277 respective credible set. We did not include the three GWAS conditional hits in the *CYP2A6*
 278 region from the original GRS²⁰ in our updated GRS given that two of these SNPs (rs12459249
 279 and rs111645190) were in high LD ($r^2 > 0.80$) with putative causal variants identified by fine-
 280 mapping (rs10853742 and rs28399451, respectively) and the remaining SNP (rs185430475) did
 281 not show robust association with NMR in our updated analysis ($p > 1 \times 10^{-4}$). To construct the
 282 updated GRS, the effect size of each putative causal variant was estimated separately in KIS-3
 283 and PNAT-2 by association testing in SNPTEST v2.5.2³⁵ using linear regression to test the
 284 association of imputed genotype dosages with square-root transformed NMR as the outcome
 285 variable using an additive genotypic model with adjustment for age, sex, BMI, and two ancestry-
 286 informative dimensions to account for population substructure as covariates. Given that the
 287 overall variance in log-NMR explained was comparable for GRSs with variant weights derived
 288 from linear regression against square-root or rank-transformed NMR, square-root transformed
 289 NMR was used for comparability of weights with the original GRS.²⁰ The overall effect size for
 290 each variant was then estimated in the total sample (KIS-3 and PNAT-2) by fixed-effects meta-
 291 analysis using the meta v1.7 R package,⁴¹ followed by multiplication of the resultant β
 292 coefficient by the standard deviation of the sqrt-NMR to unstandardize the scores.²⁰ The GRS
 293 was then computed for each n individual in the total sample as follows, where d refers to the
 294 number of risk alleles and β refers to the effect size for each i variant included in the GRS:

$$wGRS = \sum_{i=1}^n \beta_i * d_i$$

296 To evaluate the performance of the updated and original GRSs,²⁰ we first calculated the variance
 297 in log-transformed NMR (log-NMR, which best represents the nicotine clearance rate⁴²)

298 explained by each GRS in linear regression models of log-NMR ~ GRS using the R function
299 `lm`.³³ We also evaluated the variance in log-NMR explained by a GRS that included only the five
300 variants identified by conditional analysis, and the six putative causal variants identified by
301 FINEMAP.

302

303 Next, we compared the transferability of the updated and original GRSs²⁰ from African to
304 European populations by calculating the variance explained in log-NMR by each GRS in the
305 European ancestry subset of PNAT-2 (N=933).

306

307 **Results**

308 Clinical characteristics of the final discovery sample are presented in **Table 1**. From PNAT-2,
309 two samples were excluded due to missing or outlying normalized NMR values. From KIS-3,
310 eight samples were excluded due to cotinine concentrations <10ng/mL (which suggest non-daily
311 smoking⁴³), and one sample was excluded due to missing BMI. After quality control, our final
312 sample therefore comprised 953 African American smokers (n=504 from PNAT-2, and n=449
313 from KIS-3).

314

315 Following imputation using the TOPMED reference panel, 104,131 variants in the *CYP2A6*
316 region (chromosome 19:38,000,000-43,000,000bp; Genome Reference Consortium Human
317 Build 38, hg38) were available for analysis. The median INFO score for variants in the *CYP2A6*
318 region was 0.97 (mean=0.92, SD=0.096), suggesting high imputation quality. After imputation
319 using the 1000 Genomes reference panel, 46,154 variants in the *CYP2A6* region were available
320 for analysis with median INFO score 0.91 (mean=0.88, SD=0.110). Given the denser coverage

321 and higher quality genotypes obtained from imputation using the TOPMED reference panel
322 (**Figure S3**), we used imputed genotype dosages from these data for our analyses along with 12
323 directly genotyped *CYP2A6* * alleles.

324

325 Within the *CYP2A6* region a total of 113 variants showed robust association ($p < 5 \times 10^{-8}$) with
326 NMR, including four of the 12 * alleles genotyped in our sample (*CYP2A6**17, *CYP2A6**9,
327 *CYP2A6**4, and *CYP2A6**25/*26/*27, **Table S2**). Overall, these *CYP2A6* * alleles were less
328 strongly associated with NMR than other variants in the region (p-values ranging from
329 $p = 2.06 \times 10^{-26}$ for *CYP2A6**17 to $p = 4.40 \times 10^{-8}$ for *CYP2A6**25/*26/*27, **Table S2**). The strongest
330 association was observed for rs11878604 (beta=-0.689, $p = 4.75 \times 10^{-44}$), a SNP located ~16kb 3'
331 of *CYP2A6* (**Figure 1**). This lead variant had a RegulomeDB probability score of 0.69 (scores
332 range from 0 to 1, with 1 most likely to represent a variant with regulatory function);³⁹
333 rs11878604 was also identified as an adrenal eQTL for *CYP2A6* in the GTEx Project, with the
334 allele associated with lower NMR (i.e. reduced *CYP2A6* activity) showing association with
335 decreased *CYP2A6* expression in adrenal gland tissue (**Table S1, Figure S4**).

336

337 Stepwise conditional analysis with SNPTTEST³⁵ identified five independent associations with
338 NMR in the *CYP2A6* region (**Figure 1, Table S1**). Only the lead variant (rs11878604) was
339 identified as an eQTL for *CYP2A6* in GTEx. After conditioning on imputed rs11878604
340 genotype dosage, a second independent association was identified with the directly genotyped
341 *CYP2A6**4 allele (beta=-1.033, $p = 8.54 \times 10^{-13}$). The *CYP2A6**4 allele confers a whole gene
342 deletion of *CYP2A6*, and individuals with this allele have correspondingly decreased *CYP2A6*
343 activity.^{44,45} Notably, in our sample *CYP2A6**4 was not in LD with any other individual variant

344 in the region (all $r^2 < 0.15$), consistent with previous literature indicating that *CYP2A6*4* cannot
345 be tagged by nearby SNPs.⁴⁶ *CYP2A6*4* was not genotyped in the 1000 Genomes Phase 3
346 African populations used as an LD reference for construction of regional association plots by
347 LocusZoom, and as such there is no LD information displayed on the *CYP2A6*4* regional
348 association plot (**Figure 1b**). Conditioning on rs11878604 and *CYP2A6*4* revealed a third
349 independent association with rs10853742 located ~9kb 3' of *CYP2A6* (beta=0.405, $p=5.65 \times 10^{-12}$),
350 a SNP with a RegulomeDB probability score of 0.61 that was identified as a skin eQTL for
351 *CYP2A7* in the GTEx Project (**Table S1, Figure S4**). Conditioning on rs11878604, *CYP2A6*4*,
352 and rs10853742 identified a fourth independent association with rs28399451 (beta=-0.3398,
353 $p=5.59 \times 10^{-10}$). Located within intron 6 of *CYP2A6*, rs28399451 had a RegulomeDB probability
354 score of 0.135 and was identified as a skin and peripheral nerve eQTL for *CYP2A7* in the GTEx
355 Project (**Table S1, Figure S4**). Conditioning on genotype dosages of these four variants
356 (rs11878604, *CYP2A6*4*, rs10853742, rs28399451) identified a fifth independent association
357 with rs116670633 (beta=-0.676, $p=6.27 \times 10^{-10}$); this SNP was located ~85kb 5' of *CYP2A6*, had
358 a RegulomeDB probability score of 0.135, and was not identified as an eQTL in the GTEx
359 Project. After conditioning on these five variants, there were no remaining genome-wide
360 associations with NMR (**Figure 1**). These findings were consistent when association testing was
361 run independently in PNAT-2 and KIS-3 and then meta-analyzed using an inverse-variance
362 weighting approach (**Table S1**).

363

364 Bayesian fine-mapping with FINEMAP³⁷ identified six causal variants contributing to the
365 *CYP2A6* region association with NMR (posterior probability of six causal variants in the region,
366 PP=0.67). The top causal configuration included *CYP2A6*4*, rs116670633, *CYP2A6*9*,

367 rs28399451, rs8192720, and rs10853742; the posterior probability of these six variants
368 representing the true causal configuration was 0.090, and together they explained 31% of the
369 heritability of NMR (**Figure 2**). In addition to the top causal configuration, Bayesian fine-
370 mapping identified six “credible sets” (**Figure 2, Table 2**); each credible set can be interpreted as
371 containing a causal variant with 95% coverage probability. The lead variants in credible sets 1-5
372 were highly likely to be causal (*CYP2A6*4*, rs116670633, *CYP2A6*9*, rs28399451, rs8192720;
373 PIP for these variants being truly causal >0.50). Four of the putative causal variants identified by
374 FINEMAP were also identified by conditional analysis (*CYP2A6*4*, rs116670633, rs28399451,
375 rs10853742). Exploratory functionally-informed FINEMAP analyses specifying a maximum of
376 six causal variants and upweighting the 12 *CYP2A6* * alleles, which have well characterized
377 functional effects on CYP2A6 activity (summarized in **Table S2**), provided consistent results
378 and did not identify any alternative putative causal variants.

379
380 The six credible sets were made up of differing numbers of putatively causal variants, typically
381 in high LD with each other (**Figure S5**). **Credible set 1** included only *CYP2A6*4* (PIP=1),
382 which was not in significant LD with any other variant in the region. As described above,
383 *CYP2A6*4* is a whole-gene deletion variant conferring absent CYP2A6 activity;⁴⁵ because it is a
384 structural variant, *CYP2A6*4* eQTL data is not available in existing eQTL datasets which use
385 array-based technology for genotyping. **Credible set 2** included only rs116670633, which as
386 described above, is a SNP located ~85kb upstream of *CYP2A6* with limited evidence of
387 regulatory function (PIP=0.985); this variant was not in LD with any of the variants in other
388 credible sets, but was in low LD with *CYP2A6*35* ($r^2=0.46$). **Credible set 3** included *CYP2A6*9*
389 (PIP=0.890), a functional promoter region variant that decreases CYP2A6 activity, along with 22

390 other SNPs in linkage disequilibrium with *CYP2A6**9 that each had very low PIPs (PIP
391 range=0.001-0.02, **Table S3**). **Credible set 4** included three variants in high LD with each other
392 (**Figure S5**), with lead variant rs28399451 (PIP=0.603). The variants in credible set 4 were also
393 in moderate LD with *CYP2A6**17 ($r^2=0.67-0.70$). One variant in credible set 4 (rs28399439) was
394 an adipose eQTL for *CYP2A6* in GTEx, although unexpectedly the allele associated with lower
395 NMR (i.e. slower CYP2A6 activity) was associated with increased *CYP2A6* expression (**Table 2**,
396 **Figure S4**). The remaining two variants in credible set 4 (lead variant rs28399451 and
397 rs4803380) were skin and peripheral nerve eQTLs for *CYP2A7*. **Credible set 5** included three
398 variants in high LD with each other (**Figure S5**), with the top variant being rs8192720
399 (PIP=0.574). The variants in credible set 5 were in moderate LD with *CYP2A6**25/*26/*27
400 ($r^2=0.50-0.53$) and low LD with *CYP2A6**20 ($r^2=0.37-0.39$); these three variants were not
401 identified as eQTLs in GTEx (**Table 2**). **Credible set 6** included four variants, with lead variant
402 rs10853742 (PIP=0.448). The variants in credible set 6 were in low LD with the lead variant
403 from conditional analysis (rs11878604, $r^2=0.46$). All four variants in credible set 6 were skin
404 eQTLs for *CYP2A7* in GTEx (**Table 2**, **Figure S4**).

405

406 Incorporating the putative causal variants identified through fine-mapping into our existing
407 ancestry-specific GRS²⁰ resulted in a new “updated GRS.” As a benchmark, the “original GRS”
408 comprising eight *CYP2A6* * alleles and three SNPs (rs12459249, rs111645190, rs185430475)
409 identified in an earlier conditional analysis¹⁵ explained 33.2% of the variance in log-NMR in our
410 sample of African American smokers (**Figure 3a**, **Table 3**). The updated GRS included the same
411 eight *CYP2A6* * alleles, excluded rs185430475, and included four new SNPs identified by fine-
412 mapping (rs11667603, rs8192720, rs10853742, rs28399451). Two of these new putative causal

413 variants (rs10853742, rs28399451) were represented by tag SNPs in the original GRS in the
414 African ancestry sample (**Figure S5**), while in the European ancestry sample only rs10853742
415 was represented by a proxy variant in the original GRS ($r^2=0.95$ with rs12459249). The updated
416 GRS showed similar prediction of NMR as the original GRS within the African ancestry training
417 sample (variance in log-NMR $R^2=0.345$ vs. 0.332 for the original GRS; **Figure 3a-c, Table 3**),
418 and improved prediction of NMR in an independent European ancestry sample ($R^2=0.282$ vs.
419 0.228 for the original GRS; **Figure 3b-d**). In comparison, a GRS including the six FINEMAP
420 putative causal variants alone improved prediction of NMR to a lesser degree ($R^2=0.334$ vs.
421 0.332 for the original GRS in African and $R^2=0.251$ vs. 0.228 for the original GRS in European
422 ancestry; **Table 3**), suggesting the SNPs identified by fine-mapping provide independent
423 predictive information from *CYP2A6* * alleles.

424

425 **Discussion**

426 In this study we evaluated the strong regional association of *CYP2A6* with NMR among African
427 Americans participating in two large clinical trials of smoking cessation, performing an updated
428 conditional analysis and novel fine-mapping analyses which improved an existing tool to
429 genetically predict NMR. Importantly, our analyses focused on treatment-seeking individuals
430 participating in clinical trials of smoking cessation, which excluded individuals with serious
431 medical or psychiatric comorbidities (including comorbid substance use) and those who were
432 pregnant or breastfeeding. As such, an important future direction will be to expand these
433 analyses in community samples of smokers to evaluate external validity in the general
434 population.

435

436 Previous conditional analysis of the *CYP2A6* regional association in this sample described by
437 Chenoweth *et al* identified three independent associations (rs12459249, rs111645190,
438 rs185430475);¹⁵ this earlier work did not include *CYP2A6* * alleles, and used an older reference
439 panel for genotype imputation resulting in low-density SNP coverage. The conditional analyses
440 and fine-mapping presented here included denser SNP genotyping coverage and 12 directly
441 genotyped *CYP2A6* * alleles (several of which are structural variants with robust functional
442 effects on *CYP2A6* activity),⁴⁷⁻⁵⁶ providing a more comprehensive view of variation in the
443 *CYP2A6* region than any previous study in this population. In addition to confirming two
444 previously reported *CYP2A6* associations with NMR in African American smokers, our
445 conditional analysis identified three novel associations: rs11878604, *CYP2A6**4 (full *CYP2A6*
446 gene deletion), and rs116670633.

447
448 In this first fine-mapping effort of the *CYP2A6* regional association with NMR in African
449 populations to date, we identified six causal variants in the region (posterior probability,
450 PP=0.67). Prior fine-mapping using a similar analytic approach in European populations
451 identified 13 causal variants in the region. The variants comprising the top causal configuration
452 in our African ancestry sample were distinct from those in Europeans (*CYP2A6**4, rs116670633,
453 *CYP2A6**9, rs28399451, rs8192720, rs1085374; PP=0.090), and explained 31% of the
454 heritability of NMR. Interestingly, *CYP2A6**9 is a known functional allele conferring reduced
455 *CYP2A6* activity,⁵⁰ while the remaining four lead SNPs identified by FINEMAP were not
456 associated with altered *CYP2A6* expression in GTEx (recognizing that regulatory information in
457 publicly available databases is limited by methodological challenges inherent in measuring
458 *CYP2A6* gene expression levels due to structural and copy number variation in this region, as

459 well as high sequence homology with pseudogene *CYP2A7*). Importantly, the top putative causal
460 variant identified was *CYP2A6*4* (PIP=1), a loss-of-function mutation conferring whole gene
461 deletion of *CYP2A6*. *CYP2A6*4* is not included in the vast majority of genomic studies because
462 it cannot be genotyped accurately using array-based technologies, and is not tagged by any
463 individual SNP in the region.⁴⁶ The strong evidence we observed for a causal association
464 between *CYP2A6*4* and NMR highlights the importance of including *CYP2A6* structural variants
465 in future genetic studies of tobacco-related phenotypes. To help facilitate their inclusion we
466 recently developed a method to impute *CYP2A6* structural variants from SNP haplotypes
467 obtained using standard genotyping array data (sensitivity >60%, false positive rate <1% in both
468 African and European ancestry populations).²⁵

469
470 Finally, we demonstrated that an updated GRS including the putative causal variants identified in
471 African American smokers (versus those identified by conditional analysis in an earlier GRS)
472 captured similar amounts of variation in log-NMR in African ancestry individuals, and improved
473 the portability of the GRS to European ancestry individuals. Future work evaluating the
474 performance of our updated GRS in independent validation samples including diverse ancestry
475 smokers is needed to evaluate whether this improved portability extends across other ancestries.
476 One potential explanation for the improved performance of our African ancestry-specific
477 updated GRS within European smokers is that fine-mapping identified novel variants influencing
478 NMR that were not represented in the original GRS (i.e. rs11670633, rs8192720). Additionally,
479 prior work has demonstrated that including putative causal variants identified by fine-mapping
480 improves the transferability of GRS across diverse populations because of differences in LD
481 structure which result in tag SNPs from one ancestral population no longer being good proxies

482 for the underlying true causal variants in other ancestral populations.^{57,58} Consistent with this, the
483 LD patterns between tag SNPs included in our original GRS and the four putatively causal SNPs
484 included in the updated GRS were different in our African and European samples.

485

486 Overall, our results further elucidate the genetic architecture of the *CYP2A6* regional association
487 with NMR among African American smokers and provide a shortlist of variants that may
488 causally influence nicotine clearance in this population, which could be prioritized for
489 investigation in future functional studies of *CYP2A6* activity. In particular, the strong evidence
490 for a causal association observed between *CYP2A6*4* and NMR highlights the importance of
491 including *CYP2A6* structural variants in future genetic studies of tobacco-related phenotypes.
492 Finally, the potential utility of genomic data including genetic risk scores (GRS) in medical
493 decision making is growing and complements the utility of other biomarkers such as NMR,
494 particularly in situations where NMR measurements are not available or feasible (i.e. non-
495 smokers). Given that incorporating putative causal variants improved trans-ancestry portability
496 of an existing GRS for NMR in this study, our results demonstrate the broader value of fine-
497 mapping efforts as a tool to refine and improve the potential clinical utility of GRS across
498 diverse populations which may ultimately help address potential health disparities exacerbated
499 by existing Euro-centric GWAS data.¹³

500 **Acknowledgements**

501 Computations were performed on the CAMH Specialized Computing Cluster, funded by the
502 Canada Foundation for Innovation Research Hospital Fund. The GTEx Project was supported by
503 the Common Fund of the Office of the Director of the National Institutes of Health, and by NCI,
504 NHGRI, NHLBI, NIDA, NIMH, and NINDS. This work was funded by a Canadian Institutes of
505 Health Research (CIHR) Project grant (PJY-159710), National Institutes of Health (NIH) Grants
506 PGRN DA020830, CA091912, and R35 CA197461, and a Canada Research Chair in
507 Pharmacogenomics (Tyndale). Ahluwalia funded in part by P20GM130414, a NIH funded
508 Center of Biomedical Research Excellence (COBRE)

509

510 **Conflict of Interest**

511 The other authors declare no conflicts of interest. The funders had no role in study design, data
512 collection and analysis, decision to publish, or preparation of the manuscript. Dr. Ahluwalia
513 received sponsored funds for travel expenses as a speaker for the 2021 and 2022 annual GTNF
514 conference. Dr. Ahluwalia serves as a consultant and has equity in Qnovia, a start-up company
515 developing a prescription nicotine replacement product for FDA approval. Other authors declare
516 that they have no competing interests.

517

518 **References**

519

- 520 1. National Center for Chronic Disease Prevention and Health Promotion (US) Office on
 521 Smoking and Health & Atlanta (GA): Centers for Disease Control and Prevention (US).
 522 *The Health Consequences of Smoking—50 Years of Progress: A Report of the Surgeon*
 523 *General*. (2014).
- 524 2. Benowitz, N. L. Nicotine addiction. *N Engl J Med* **362**, 2295–2303 (2010).
- 525 3. Nakajima, M. *et al.* Role of human cytochrome P4502A6 in C-oxidation of nicotine. *Drug*
 526 *Metabolism and Disposition* **24**, 1212–1217 (1996).
- 527 4. Nakajima, M. *et al.* Characterization of CYP2A6 involved in 3'-hydroxylation of cotinine
 528 in human liver microsomes. *Journal of Pharmacology and Experimental Therapeutics*
 529 **277**, 1010–1015 (1996).
- 530 5. Benowitz, N. L., St Helen, G., Dempsey, D. A., Jacob III, P. & Tyndale, R. F. Disposition
 531 kinetics and metabolism of nicotine and cotinine in African American smokers: Impact of
 532 CYP2A6 genetic variation and enzymatic activity. *Pharmacogenet Genomics* **26**, 340–350
 533 (2016).
- 534 6. Dempsey, D. *et al.* Nicotine metabolite ratio as an index of cytochrome P450 2A6
 535 metabolic activity. *Clin Pharmacol Ther* **76**, 64–72 (2004).
- 536 7. Chenoweth, M. J. & Tyndale, R. F. Pharmacogenetic optimization of smoking cessation
 537 treatment. *Trends Pharmacol Sci* **38**, 66 (2017).
- 538 8. Murphy, S. E. Biochemistry of nicotine metabolism and its relevance to lung cancer.
 539 *Journal of Biological Chemistry* **296**, 1–16 (2021).
- 540 9. Lerman, C. *et al.* Nicotine metabolite ratio predicts efficacy of transdermal nicotine for
 541 smoking cessation. *Clin Pharmacol Ther* **79**, 600–608 (2006).
- 542 10. Lerman, C. *et al.* Use of the nicotine metabolite ratio as a genetically informed biomarker
 543 of response to nicotine patch or varenicline for smoking cessation: A randomised, double-
 544 blind placebo-controlled trial. *Lancet Respir Med* **3**, 131–138 (2015).
- 545 11. Loukola, A. *et al.* A genome-wide association study of a biomarker of nicotine
 546 metabolism. *PLoS Genet* **11**, e1005498 (2015).
- 547 12. Swan, G. E. *et al.* Genetic and environmental influences on the ratio of 3'-hydroxycotinine
 548 to cotinine in plasma and urine. *Pharmacogenet Genomics* **19**, 398 (2009).
- 549 13. Martin, A. R. *et al.* Current clinical use of polygenic scores will risk exacerbating health
 550 disparities. *Nat Genet* **51**, 584–591 (2019).
- 551 14. Buchwald, J. *et al.* Genome-wide association meta-analysis of nicotine metabolism and
 552 cigarette consumption measures in smokers of European descent. *Mol Psychiatry* **26**, 2223
 553 (2021).
- 554 15. Chenoweth, M. J. *et al.* Genome-wide association study of a nicotine metabolism
 555 biomarker in African American smokers: impact of chromosome 19 genetic influences.
 556 *Addiction* **113**, 523 (2018).
- 557 16. Loukola, A. *et al.* A genome-wide association study of a biomarker of nicotine
 558 metabolism. *PLoS Genet* **11**, e1005498 (2015).
- 559 17. Alexander, T. A. & Machiela, M. J. LDpop: an interactive online tool to calculate and
 560 visualize geographic LD patterns. *BMC Bioinformatics* **21**, (2020).

- 561 18. Wassenaar, C. A., Zhou, Q. & Tyndale, R. F. CYP2A6 genotyping methods and strategies
562 using real-time and end point PCR platforms. *Pharmacogenomics* **17**, 147–162 (2015).
- 563 19. El-Boraie, A. *et al.* Evaluation of a weighted genetic risk score for the prediction of
564 biomarkers of CYP2A6 activity. *Addiction Biology* **25**, e12741 (2020).
- 565 20. El-Boraie, A. *et al.* Transferability of ancestry-specific and cross-ancestry CYP2A6
566 activity genetic risk scores in African and European populations. *Clin Pharmacol Ther*
567 **110**, 975–985 (2021).
- 568 21. Bloom, J. *et al.* The contribution of common CYP2A6 alleles to variation in nicotine
569 metabolism among European-Americans. *Pharmacogenet Genomics* **21** (2011).
- 570 22. Baurley, J. W. *et al.* Predicting nicotine metabolism across ancestries using genotypes.
571 *BMC Genomics* **23** (2022).
- 572 23. Cox, L. S. *et al.* Bupropion for smoking cessation in African American light Smokers: A
573 randomized controlled trial. *J Natl Cancer Inst* **104**, 290–298 (2012).
- 574 24. Jacob, P. *et al.* Determination of the nicotine metabolites cotinine and trans-3'-
575 hydroxycotinine in biologic fluids of smokers and non-smokers using liquid
576 chromatography-tandem mass spectrometry: biomarkers for tobacco smoke exposure and
577 for phenotyping cytochrome P450 2A6 activity. *J Chromatogr B Analyt Technol Biomed*
578 *Life Sci* **879**, 267–276 (2011).
- 579 25. Langlois, A. W. R. *et al.* Genotyping, characterization, and imputation of known and
580 novel CYP2A6 structural variants using SNP array data. *J Hum Genet* (2023)
581 doi:10.1038/S10038-023-01148-Y.
- 582 26. Purcell, S. *et al.* PLINK: A tool set for whole-genome association and population-based
583 linkage analyses. *Am J Hum Genet* **81**, 559–575 (2007).
- 584 27. Das, S. *et al.* Next-generation genotype imputation service and methods. *Nat Genet* **48**,
585 1284–1287 (2016).
- 586 28. Taliun, D. *et al.* Sequencing of 53,831 diverse genomes from the NHLBI TOPMed
587 Program. *Nature* **590**, 290–299 (2021).
- 588 29. Auton, A. *et al.* A global reference for human genetic variation. *Nature* **526**, 68–74
589 (2015).
- 590 30. Loh, P. R. *et al.* Reference-based phasing using the Haplotype Reference Consortium
591 panel. *Nat Genet* **48**, 1443–1448 (2016).
- 592 31. O'Connell, J. *et al.* A general approach for haplotype phasing across the full spectrum of
593 relatedness. *PLoS Genet* **10**, e1004234 (2014).
- 594 32. Chenoweth, M. J. *et al.* Analyses of nicotine metabolism biomarker genetics stratified by
595 sex in African and European Americans. *Sci Rep* **11**, 1–12 (2021).
- 596 33. R Core Team. R: A language and environment for statistical computing. Preprint at
597 (2021).
- 598 34. McCaw Z. RNOmni: Rank Normal Transformation Omnibus Test. R package version
599 1.0.0. (2020).
- 600 35. Marchini, J., Howie, B., Myers, S., McVean, G. & Donnelly, P. A new multipoint method
601 for genome-wide association studies by imputation of genotypes. *Nat Genet* **39**, 906–913
602 (2007).
- 603 36. Boughton, A. P. *et al.* LocusZoom.js: interactive and embeddable visualization of genetic
604 association study results. *Bioinformatics* **37**, 3017–3018 (2021).
- 605 37. Benner, C. *et al.* FINEMAP: efficient variable selection using summary data from
606 genome-wide association studies. *Bioinformatics* **32**, 1493–1501 (2016).

- 607 38. Benner, C. *et al.* Prospects of fine-mapping trait-associated genomic regions by using
608 summary statistics from genome-wide association studies. *Am J Hum Genet* **101**, 539–551
609 (2017).
- 610 39. Dong, S. & Boyle, A. P. Predicting functional variants in enhancer and promoter elements
611 using RegulomeDB. *Hum Mutat* **40**, 1298 (2019).
- 612 40. Lonsdale, J. *et al.* The Genotype-Tissue Expression (GTEx) project. *Nat Genet* **45**, 580–
613 585 (2013).
- 614 41. Balduzzi, S., Rücker, G. & Schwarzer, G. How to perform a meta-analysis with R: A
615 practical tutorial. *Evid Based Ment Health* **22**, 153–160 (2019).
- 616 42. Tanner, J. A. *et al.* Nicotine metabolite ratio (3-hydroxycotinine/cotinine) in plasma and
617 urine by different analytical methods and laboratories: implications for clinical
618 implementation. *Cancer Epidemiol Biomarkers Prev* **24**, 1239–1246 (2015).
- 619 43. Vartiainen, E., Seppälä, T., Lillsunde, P. & Puska, P. Validation of self reported smoking
620 by serum cotinine measurement in a community-based study. *J Epidemiol Community*
621 *Health* **56**, 167–170 (2002).
- 622 44. Gaedigk, A. *et al.* The Pharmacogene Variation (PharmVar) Consortium: Incorporation of
623 the human cytochrome P450 (CYP) allele nomenclature database. *Clin Pharmacol Ther*
624 **103**, 401 (2018).
- 625 45. Mwenifumbo, J. C., Zhou, Q., Benowitz, N. L., Sellers, E. M. & Tyndale, R. F. New
626 CYP2A6 gene deletion and conversion variants in a population of Black African descent.
627 *Pharmacogenomics* **11**, 189 (2010).
- 628 46. Kumasaka, N. *et al.* Haplotypes with copy number and single nucleotide polymorphisms
629 in CYP2A6 locus are associated with smoking quantity in a Japanese population. *PLoS*
630 *One* **7**, e44507 (2012).
- 631 47. Mwenifumbo, J. C. *et al.* Identification of novel CYP2A6*1B variants: The CYP2A6*1B
632 allele is associated with faster in vivo nicotine metabolism. *Clin Pharmacol Ther* **83**, 121
633 (2008).
- 634 48. Rao, Y. *et al.* Duplications and defects in the CYP2A6 gene: Identification, genotyping,
635 and in vivo effects on smoking. *Mol Pharmacol* **58**, 747–755 (2000).
- 636 49. Nunoya, K. ichi *et al.* A new deleted allele in the human cytochrome P450 2A6
637 (CYP2A6) gene found in individuals showing poor metabolic capacity to coumarin and
638 (+)-cis-3,5-dimethyl-2-(3-pyridyl)thiazolidin-4-one hydrochloride. *Pharmacogenetics* **8**,
639 239–249 (1998).
- 640 50. Kiyotani, K. *et al.* Decreased coumarin 7-hydroxylase activities and CYP2A6 expression
641 levels in humans caused by genetic polymorphism in CYP2A6 promoter region
642 (CYP2A6*9). *Pharmacogenetics* **13**, 689–695 (2003).
- 643 51. Oscarson, M. *et al.* Characterization of a novel CYP2A7/CYP2A6 hybrid allele
644 (CYP2A6*12) that causes reduced CYP2A6 activity. *Hum Mutat* **20**, 275–283 (2002).
- 645 52. Fukami, T. *et al.* A novel polymorphism of human CYP2A6 gene CYP2A6*17 has an
646 amino acid substitution (V365M) that decreases enzymatic activity in vitro and in vivo.
647 *Clin Pharmacol Ther* **76**, 519–527 (2004).
- 648 53. Fukami, T. *et al.* A novel CYP2A6*20 allele found in African-American population
649 produces a truncated protein lacking enzymatic activity. *Biochem Pharmacol* **70**, 801–808
650 (2005).
- 651 54. Ho, M. K., Mwenifumbo, J. C., Zhao, B., Gillam, E. M. J. & Tyndale, R. F. A novel
652 CYP2A6 allele, CYP2A6*23, impairs enzyme function in vitro and in vivo and decreases

653 smoking in a population of Black-African descent. *Pharmacogenet Genomics* **18**, 67–75
654 (2008).

655 55. Mwenifumbo, J. C. *et al.* Novel and established CYP2A6 alleles impair in vivo nicotine
656 metabolism in a population of Black African descent. *Hum Mutat* **29**, 679–688 (2008).

657 56. Al Koudsi, N., Ahluwalia, J. S., Lin, S. K., Sellers, E. M. & Tyndale, R. F. A novel
658 CYP2A6 allele (CYP2A6*35) resulting in an amino-acid substitution (Asn438Tyr) is
659 associated with lower CYP2A6 activity in vivo. *Pharmacogenomics J* **9**, 274–282 (2009).

660 57. Weissbrod, O. *et al.* Leveraging fine-mapping and multipopulation training data to
661 improve cross-population polygenic risk scores. *Nat Genet* **54**, 450–458 (2022).

662 58. Amariuta, T. *et al.* Improving the trans-ancestry portability of polygenic risk scores by
663 prioritizing variants in predicted cell-type-specific regulatory elements. *Nat Genet* **52**,
664 1346–1354 (2020).

665 59. Austin, P. C. An introduction to propensity score methods for reducing the effects of
666 confounding in observational studies. *Multivariate Behav Res* **46**, 399–424 (2011).

667

668 **Figure Legends**

669 **Figure 1. Conditional analysis of *CYP2A6* region identified five independent associations**
670 **with NMR in African ancestry smokers (a-e), including *CYP2A6* deletion variant**
671 ***CYP2A6**4 (b).** Genomic positions based on Genome Reference Consortium build 38, hg38.

672 **Figure 2. Bayesian fine-mapping of *CYP2A6* association with NMR.** Top causal configuration
673 included *CYP2A6**4, rs116670633, *CYP2A6**9, rs28399451, rs8192720, and rs10853742;
674 posterior probability of this top configuration being truly causal=0.090; NMR heritability
675 explained by top configuration (h^2)=0.31.

676 **Figure 3. Variance in log-NMR explained by the original GRS in African American**
677 **smokers (a) and its portability to European ancestry smokers (b), as well as the updated**
678 **GRS in African American smokers (c) and its portability to European ancestry smokers**
679 **(d).** The original GRS comprised * alleles and SNPs identified in a previous conditional analysis,
680 whereas the updated GRS replaced these SNPs with putative causal SNPs identified by fine-
681 mapping (for details of the variants included in the original and updated GRS, see **Table 3**). R^2
682 represents the variance in log-NMR explained.

683 **Tables**

684

685 **Table 1.** Sociodemographic and clinical characteristics of the final study sample

	Total Sample (n=953)	PNAT-2 (n=504)	KIS-3 (n=449)	Standardized Difference^a
% Female (n)	57.9 (552)	50.4 (254)	66.4 (298)	0.33
Age ± SD (range)	47.1 ± 10.7 (19-80)	47.3 ± 9.8 (20-65)	46.8 ± 11.6 (19-80)	0.04
BMI ± SD (range)	30.8 ± 7.5 (15-68)	30.5 ± 7.1 (18-58)	31.2 ± 7.8 (15-68)	0.10
Cigarettes/day ± SD (range)	12.3 ± 6.4 (1-40)	16.3 ± 6.3 (5-40)	7.8 ± 2.6 (1-17)	1.76
Cotinine (ng/mL) ± SD (range)	260 ± 128 (14-837)	274 ± 130 (32-837)	244 ± 123 (14-681)	0.24
NMR ± SD (range)	0.35 ± 0.23 (0.01-1.79)	0.33 ± 0.20 (0.01-1.17)	0.38 ± 0.26 (0.02-1.79)	0.23

^aStandardized differences (SD) were used to evaluate differences in study covariates between the two clinical trial samples included in the current study, with SD < 0.1 generally accepted as indicating a minimal difference between groups.⁵⁹ SD compare differences in mean/prevalence in units of the pooled standard deviation, which allows for comparison of the relative balance of variables in different units, and are not influenced by sample size.⁵⁹

686

687

688

689

690

691

692

693

694 **Table 2.** Association with NMR and functional annotations for *CYP2A6* region variants identified by fine-mapping

95% Credible Set	Variant	Chromosome 19 Position (bp) ^a	Location Relative to <i>CYP2A6</i>	Ref Allele	Effect Allele	MAF ^b	INFO ^c	Beta ^d	SE ^d	PIP ^e	log ₁₀ BF ^f	GTE _x Project eQTLs ^g	RegulomeDB Probability Score ^h
1	<i>CYP2A6</i> *4	40843541-40850447	Whole gene deletion	-	Deletion	0.02	Typed	-1.033	0.143	1	14.57	Not available ^h	Not available ⁱ
2	rs116670633	40935245	84.8kb 5'	T	G	0.03	0.99	-0.409	0.129	0.989	6.53	None	0.135
3 ^j	<i>CYP2A6</i> *9 (rs28399433)	40843541-40850447	Promoter (TATA box)	A	C	0.08	Typed	-0.493	0.077	0.788	5.14	<i>CYP2A6</i> (adrenal): NES=-0.51; p=6.0x10 ⁻⁶ <i>CYP2A7</i> (lung): NES=0.45; p=6.4x10 ⁻⁶ <i>EGLN2</i> (artery): NES=-0.25; p=1.2x10 ⁻⁵	0.554
4	rs28399451	40845938	Intron 6	G	A	0.14	0.93	-0.689	0.065	0.616	4.77	<i>CYP2A7</i> (skin): NES=0.73; p=6.6x10 ⁻⁹ <i>CYP2A7</i> (peripheral nerve): NES=0.52; p=7.8x10 ⁻⁵	0.135
	rs4803380	40845264	Intron 7	C	T	0.13	0.95	-0.691	0.066	0.339	4.27	<i>CYP2A7</i> (skin): NES=0.73; p=5.3x10 ⁻⁹ <i>CYP2A7</i> (peripheral nerve): NES=0.52; p=7.8x10 ⁻⁵	0.778
	rs28399439	40849808/12	Intron 2	AC	A	0.13	0.98	-0.700	0.065	0.022	2.92	<i>CYP2A6</i> (adipose): NES=0.69; p=5.9x10 ⁻⁵	0.983
5	rs8192720	40850405	Exon 1, synonymous	G	A	0.04	0.99 ^k	-0.792	0.113	0.546	4.65	None	0.609
	rs72549439	40848131	Intron 4	G	A	0.04	0.96 ^k	-0.754	0.106	0.228	4.04	None	0.244
	rs72549445	40845791	Intron 6	T	G	0.04	0.93	-0.775	0.110	0.195	3.95	None	0.981
6	rs10853742	40834668	8.9kb 3'	G	C	0.33	0.99 ^k	0.623	0.043	0.433	4.45	<i>CYP2A7</i> (skin): NES=0.23; p=1.7x10 ⁻⁵	0.609
	rs7251570	40835845	7.7kb 3'	A	G	0.34	0.95	0.636	0.044	0.300	4.20	<i>CYP2A7</i> (skin): NES=0.22; p=2.7x10 ⁻⁵	0.590

rs11667314	40835078	8.5kb 3'	T	C	0.34	0.95	0.634	0.044	0.160	3.85	<i>CYP2A7</i> (skin): NES=0.22; p=2.7x10 ⁻⁵	0.507
rs3865454	40836554	7.0kb 3'	T	G	0.34	0.95	0.635	0.044	0.087	3.55	<i>CYP2A7</i> (skin): NES=0.22; p=2.7x10 ⁻⁵	0.729

^aHuman genome reference hg38; ^bMinor allele frequency (MAF) observed in our sample; ^cImputation quality INFO scores were using R^2 values representing the estimated true correlation between imputed and real genotypes based on sample allele frequencies, as implemented in Minimac4²⁷; ^dbeta and standard error (SE) reported are from association testing using linear regression in SNPTEST of genotype dosage ~ NMR with adjustment for age, sex, BMI, and two ancestry-informative dimensions; ^eFINEMAP output, marginal Posterior Inclusion Probabilities (PIP) for each SNP represent the posterior probability that this SNP is causal; ^fFINEMAP output, the Bayes factor quantifies the evidence that a particular SNP is causal, with log10 Bayes factors greater than 2 suggesting considerable evidence for causality; ^gPublicly available expression quantitative trait loci (eQTL) data from the Genotype-Tissue Expression (GTEx) Project⁴⁰ was used to evaluate whether variants were known to influence gene expression of protein coding genes. eQTL effect alleles correspond to the effect alleles for NMR in our study, allowing for direct comparison of the directions of effect on NMR (beta) and gene expression (normalized effect size, NES); ^hRegulomeDB is a publicly available database that estimates a variant's likelihood of regulatory function using a probability score ranging from 0 to 1 (with 1 being most likely to be a regulatory variant). The score is constructed based on a machine learning model integrating functional genomic data including ChIP-seq signal, DNase-seq signal, information content change, and DeepSEA scores;³⁹ ⁱBecause *CYP2A6*4* is a structural variant (whole gene deletion), *CYP2A6*4* genotypes are not available in existing eQTL datasets which use array-based technology for genotyping; ^jCredible set 3 also included 22 SNPs with low PIPs (mean PIP=0.003, range=0.001 – 0.02) which tagged *CYP2A6*9* to varying degrees (mean D' =0.91, range=0.41 – 1) and were therefore not included in the main table above but are detailed in **Table S2**; ^kThese variants were directly genotyped in our sample, but imputed genotype dosages were used for association testing (mean correlation between direct genotyping and imputed genotype dosages=0.88, range=0.62-0.97).

695

Table 3. Effects of incorporating top putative causal variants identified by fine-mapping into an existing genetic risk score (GRS) to predict NMR in African American smokers

Model	Variants Included	Ref Allele	Effect Allele	Beta ^c	GRS Weight ^d	African American		European	
						Effect Allele Freq ^b	R ² ^a	Effect Allele Freq ^b	R ² ^a
1 - Original GRS	<i>CYP2A6</i> *4 ^{e,f}	-	Deletion	-0.935	-0.169	0.023	0.332	0.003	0.228
	<i>CYP2A6</i> *1x2 ^e	-	Duplication	0.686	0.124	0.013		0.008	
	<i>CYP2A6</i> *9 (rs28399433) ^{e,f}	A	C	-0.473	-0.086	0.083		0.066	
	<i>CYP2A6</i> *12 ^e	-	<i>CYP2A6/2A7</i> hybrid	-0.570	-0.103	0.006		0.023	
	<i>CYP2A6</i> *17 (rs28399454) ^e	C	T	-0.699	-0.127	0.107		0.001	
	<i>CYP2A6</i> *20 (rs568811809) ^e	TT	-	-0.704	-0.127	0.015		0.000	
	<i>CYP2A6</i> *25/*26/*27 (rs28399440) ^e	A	G	-0.782	-0.142	0.022		0.000	
	<i>CYP2A6</i> *35 (rs143731390) ^e	T	A	-0.345	-0.062	0.020		0.000	
	rs12459249 ^{e,g}	T	C	0.578	0.105	0.674		0.670	
	rs111645190 ^{e,g}	G	A	-0.633	-0.115	0.139		0.000	
rs185430475 ^{e,g}	C	G	0.735	0.133	0.013	0.000			
2 - Conditional analysis variants	rs11878604	T	C	-0.651	-0.118	0.232	0.295	0.077	0.224
	<i>CYP2A6</i> *4	-	Deletion	-0.935	-0.169	0.023		0.003	
	rs10853742	G	C	0.591	0.107	0.669		0.664	
	rs28399451	G	A	-0.611	-0.111	0.139		0.024	
	rs116670633	T	G	-0.407	-0.074	0.031		0.002	
3 - FINEMAP top causal variants	<i>CYP2A6</i> *4 ^{e,f}	-	Deletion	-0.935	-0.169	0.023	0.334	0.003	0.251
	<i>CYP2A6</i> *9 (rs28399433) ^{e,f}	A	C	-0.473	-0.086	0.083		0.066	
	rs10853742 ^f	G	C	0.591	0.107	0.669		0.664	
	rs28399451 ^f	G	A	-0.611	-0.111	0.139		0.024	
	rs8192720 ^f	G	A	-0.743	-0.134	0.039		0.003	
	rs116670633 ^f	T	G	-0.407	-0.074	0.031		0.002	
4 - Updated GRS <i>Original GRS</i> *	<i>CYP2A6</i> *4 ^{e,f}	-	Deletion	-0.935	-0.169	0.023	0.345	0.003	0.282
	<i>CYP2A6</i> *1x2 ^e	-	Duplication	0.686	0.124	0.013		0.008	

<i>alleles + FINEMAP top causal variants</i>	<i>CYP2A6*9</i> (rs28399433) ^{e,f}	A	C	-0.473	-0.086	0.083	0.066
	<i>CYP2A6*12</i> ^c	-	<i>CYP2A6/2A7</i> hybrid	-0.570	-0.103	0.006	0.023
	<i>CYP2A6*17</i> (rs28399454) ^c	C	T	-0.699	-0.127	0.107	0.001
	<i>CYP2A6*20</i> (rs568811809) ^c	TT	-	-0.704	-0.127	0.015	0.000
	<i>CYP2A6*25/*26/*27</i> (rs28399440) ^c	A	G	-0.782	-0.142	0.022	0.000
	<i>CYP2A6*35</i> (rs143731390) ^c	`	A	-0.345	-0.062	0.020	0.000
	rs10853742 ^f	G	C	0.591	0.107	0.669	0.664
	rs28399451 ^f	G	A	-0.611	-0.111	0.139	0.024
	rs8192720 ^f	G	A	-0.743	-0.134	0.039	0.003
rs116670633 ^f	T	G	-0.407	-0.074	0.031	0.002	

Bold font indicates novel putative causal variants identified in the present study that are not in linkage disequilibrium with variants identified in previous non-Bayesian analyses. ^aVariance in log-NMR explained (R^2) by the GRS, estimated using linear regression of log-NMR ~ GRS; ^bEffect allele frequency observed in our sample; ^cBeta reported is from fixed-effects meta-analysis of association testing results in PNAT-2 and KIS-3 samples using linear regression in SNPTEST of genotype dosage ~ sqrt-NMR with adjustment for age, sex, BMI, and two ancestry-informative dimensions; ^dGRS weights were calculated as $\beta * SD(\text{sqrt-NMR})$ to unstandardize the scores; ^eThese variants were included in the original GRS for NMR in African American smokers described by El-Boraie *et al.*,²⁰ with beta and GRS weights updated in the current study as described in **Methods**; ^fThese variants were identified as the top putative causal variants by fine-mapping in the current study; ^gThese variants were identified by earlier conditional analysis of the *CYP2A6* regional association with NMR conducted in the current study sample, described by Chenoweth *et al.*¹⁵