

1 **Title Page**

2 **Title:** Polymethylation Scores for Prenatal Maternal Smoke Exposure Persist Until Age 15 and
3 Are Detected in Saliva

4 **Authors:** Freida A. Blostein¹, Jonah Fisher², John Dou¹, Lisa Schenper³, Erin B. Ware² Daniel
5 A. Notterman^{3*}, Colter Mitchell^{2*}, Kelly M. Bakulski^{1*+}

6 **Affiliations:**

7 ¹Department of Epidemiology, University of Michigan School of Public Health, Ann Arbor,
8 Michigan, United States of America

9 ²Institute for Social Research, University of Michigan, Ann Arbor, Michigan, United States of
10 America

11 ³Department of Molecular Biology, Princeton University, Princeton, New Jersey, United States
12 of America

13 *Co-senior authors

14 ⁺Corresponding author: bakulski@umich.edu 1415 Washington Heights, Ann Arbor, MI, 48109

15 **Abstract**

16 **Background:** Prenatal maternal smoking has negative implications for child health. DNA
17 methylation signatures can function as biomarkers of prenatal maternal smoking. However little
18 work has assessed how DNA methylation signatures of prenatal maternal smoking vary across
19 ages, ancestry groups, or tissues. In the Fragile Families and Child Wellbeing study, we tested

20 whether prenatal maternal smoking was associated with salivary polymethylation scores for
21 smoking in participants. We assessed the associations' consistency at ages 9 and 15, their
22 portability across participants from African, European, and Hispanic genetic ancestries and the
23 accuracy of exposure classification using area under the curve (AUC) from receiver operating
24 curve analyses.

25 **Results:** We created saliva polymethylation scores using coefficients from a meta-analysis of
26 prenatal maternal smoke exposure and DNA methylation in newborn cord blood. In the full
27 sample at age 9 (n=753), prenatal maternal smoke exposure was associated with a 0.52 (95%CI:
28 0.36, 0.67) standard deviation higher polymethylation score for prenatal smoke exposure. The
29 direction of the association was consistent when stratified by genetic ancestries. In the full
30 sample at age 15 (n=746), prenatal maternal smoke exposure was associated with a 0.46 (95%CI:
31 0.3, 0.62) standard deviation higher polymethylation score for prenatal smoke exposure, and the
32 effect size was attenuated among the European and Hispanic genetic ancestry samples. The
33 polymethylation score was reasonably accurate at classifying prenatal maternal smoke exposure
34 (AUC age 9=0.78, *P* value comparing to base model<0.001, age 15=0.77, *P* value comparing to
35 base model<0.001). The polymethylation score showed higher classification accuracy than using
36 a single *a priori* site in the *AHRR* gene (cg05575921 AUC=0.74, *P* value comparing to
37 polymethylation score=0.03; age 15=0.73, *P* value comparing to polymethylation score=0.01).

38 **Conclusions:** Prenatal maternal smoking was associated with DNA methylation signatures in
39 saliva samples, a clinically practical tissue. Polymethylation scores for prenatal maternal
40 smoking were portable across genetic ancestries and more accurate than individual DNA
41 methylation sites. DNA polymethylation scores from saliva samples could serve as robust and
42 practical clinical biomarkers of prenatal maternal smoke exposure.

43 **Keywords**

44 DNA methylation, prenatal maternal smoking, salivary biomarker, longitudinal cohort

45 **Background**

46 Maternal prenatal smoking is a public health concern. Children exposed to cigarette smoking in
47 utero are more likely to have low-birth weight, negative neurodevelopmental outcomes, and
48 asthma [1]. In 2018, 11% of pregnant women in the United States aged 15-44 reported any past-
49 month cigarette use [2]. Measuring prenatal maternal smoking behavior is challenging. Maternal
50 prenatal smoking is underreported due to stigma [3, 4]. The gold-standard for smoking
51 measurements, serum cotinine levels, have a half-life of nine hours in pregnant women, limiting
52 the ability for accurate detection to a short window [5]. Further, serum cotinine measures during
53 pregnancy are rarely available outside of birth cohorts, limiting clinical application as a marker
54 of prenatal exposures. Developing a portable and reliable biomarker of prenatal maternal smoke
55 exposure would have important implications for research and clinical practice.

56 In infant cord blood and placental tissues, maternal cigarette smoking during pregnancy is
57 associated with reproducible DNA methylation signatures, though this has primarily been tested
58 in European cohorts [6–8]. Prenatal maternal cigarette smoking exposure is associated with
59 postnatal DNA methylation signatures in blood from cohorts of older children and adults [9, 10].
60 Few studies have considered the persistence of the association between maternal prenatal
61 cigarette smoking and DNA methylation at different ages in the same individuals. Additionally,
62 little work has considered the portability of the association between maternal prenatal cigarette
63 smoking and DNA methylation outside of cord blood and placental tissue. DNA methylation
64 drives cell differentiation and cell type proportions differ across tissue types. Epigenetic markers

65 of prenatal maternal smoking are known to differ between cord blood and placental samples [7].
66 Cord blood and placental samples are unlikely to be available in clinical settings. While
67 peripheral blood is a more realistic clinical sample, saliva is even easier to collect. To our
68 knowledge, no study has evaluated associations between maternal prenatal cigarette smoking and
69 salivary DNA methylation.

70 Additionally, DNA methylation can vary by genetics [11–14], and DNA methylation signatures
71 of own-smoking differ across genetic ancestry groups [15–18]. Only a few studies have
72 examined associations between prenatal maternal cigarette smoking and DNA methylation in
73 non-European ancestry populations. Among 954 infants, 70% of whom had mothers identifying
74 as Black, prenatal maternal smoking was associated with cord blood DNA methylation at 38
75 CpG sites [19]. The direction of the associations between prenatal maternal cigarette smoking
76 and DNA methylation at these sites was generally similar in children of Black and non-Black
77 participants [19]. Among 572 3–5-year-old children, 186 of whom were of African or admixed
78 genetic ancestry, prenatal maternal smoking was associated with DNA methylation in peripheral
79 blood samples, but ancestry-stratified analyses were not performed [20]. Similarly, among 89
80 middle-aged women, of whom 28 reported African American or Hispanic ethnicity, prenatal
81 maternal cigarette smoking was associated with DNA methylation at 17 of 190 tested CpG sites
82 [9]. Adjusting for ancestry did not substantially affect the association between prenatal maternal
83 smoking and DNA methylation at any of the sites, although the authors note that African
84 American women had on average 2–5% higher mean DNA methylation on the absolute scale as
85 compared with White and Hispanic women at 5 *CYP1A1* CpG sites [9]. Of 148 CpG sites
86 selected based on their association with prenatal maternal smoking among primarily European
87 children, 7 CpG sites were also associated with prenatal maternal cigarette smoking in a cohort

88 of 572 Latino children [6, 21]. While prenatal maternal cigarette smoking is associated with
89 DNA methylation, the consistency of this association across genetic ancestry is understudied.
90 After association testing, an important next step is to evaluate DNA methylation as a biomarker
91 of prenatal maternal smoke exposure by comparing DNA methylation-based classification of
92 maternal smoking behavior to self-reported behavior. Biomarkers can consist of single DNA
93 methylation sites or summary measures across multiple sites [17]. One option is a
94 polymethylation score, analogous to a polygenic score. In the polymethylation score, CpG sites
95 are weighted by the strength of their association with maternal smoking from a previous,
96 independent sample and then summed to a single score [24, 25]. One such DNA methylation
97 score calculated from cord blood classified maternal smoking behavior with an area under the
98 curve (AUC) of 0.82 in a primarily European cohort [23]. Among 572 3–5-year-old children,
99 186 of whom were of African or admixed ancestry, prenatal maternal smoking behavior was
100 classified with an AUC of 0.87 [20]. Among middle-aged adults, a score calculated using DNA
101 methylation in peripheral blood samples from a single time point could even predict prenatal
102 smoke exposure (30 years previously) with an AUC of 0.72 (95% confidence interval 0.69, 0.76)
103 [10]. Ideally, a biomarker for prenatal maternal smoking would be consistent over the life course
104 and portable across genetic ancestry groups. However, the consistency of DNA methylation as a
105 biomarker for prenatal maternal smoking across age, tissue, and genetic ancestry has not been
106 evaluated in the same study.

107 DNA methylation differences hold promise as a biomarker for prenatal maternal smoking.
108 Before translation to clinical applications, we must understand how and if the signal varies
109 across age, tissue sample, and genetic ancestry. In the Fragile Families and Child Wellbeing
110 study, a diverse longitudinal birth cohort of children, we aimed to assess the potential of salivary

111 DNA methylation biomarkers for prenatal maternal cigarette smoking. We tested associations
112 between prenatal maternal smoke exposure and saliva DNA methylation, including
113 polymethylation scores for smoke exposure, individual *a priori* CpG sites, global DNA
114 methylation, and epigenetic clocks. We then performed age- and ancestry-stratified analyses to
115 test the hypothesis that DNA methylation could serve as a persistent and consistent biomarker of
116 prenatal maternal smoking.

117 **Results**

118 **Study sample descriptive statistics**

119 In a sample of the Fragile Families and Child Wellbeing Study, saliva DNA methylation was
120 measured on 1806 samples from 897 unique participants with the Illumina 450K array. Complete
121 data on covariates of interest was available on 1499 samples from 809 participants who were
122 included in the analysis (Figure 1). There were 690 participants with both age 9 and age 15 DNA
123 methylation samples. Excluded samples were similar to included samples, except that included
124 samples were slightly more likely to be from children of African genetic ancestry and less likely
125 to be from children of Hispanic genetic ancestry (Supplemental Table 1). In the included sample,
126 20% percent of the mothers reported any prenatal maternal smoking, 12% reported prenatal
127 alcohol use and 5% reported prenatal drug use. The mean income to poverty ratio of mothers at
128 birth was 2.2. Of the children in the analytic sample, 50% were male, 60% were of African
129 genetic ancestry, and 24% were of Hispanic genetic ancestry.

130 We calculated several summary measures of DNA methylation, including polymethylation
131 scores for smoke exposure, global DNA methylation, and epigenetic clocks. Many of the
132 methylation summary measures were correlated with each other (Supplemental Figure 1). For

133 example, as expected, the percent of estimated immune cells was perfectly inversely correlated
134 with the percent of estimated epithelial cells (Pearson $\rho=-1$; P value <0.0001). Our
135 polymethylation score for prenatal smoke exposure, which was calculated using coefficients
136 from a cell-type corrected regression, was weakly correlated with estimated cell-type proportion
137 of immune cells (Pearson $\rho=0.05$; P value $=0.04$).

138 We compared epigenetic measures between the age 9 and age 15 visits among the 690
139 individuals with data from both visits (Supplemental Figure 2). The correlation across ages was
140 strongest for the polymethylation score for prenatal smoke exposure (Pearson $\rho=0.9$, P
141 value <0.0001). Global DNA methylation (0.5, P value <0.0001) and epigenetic clocks were less
142 strongly correlated across ages (Pediatric: 0.51; P value <0.0001 , GRIM: 0.47; P value <0.0001)
143 The distribution of epigenetic ages shifted higher between the age 9 and age 15 visits, as
144 expected (Supplemental Figure 3). The pediatric clock more closely reflected chronological age
145 than the GRIM clock. Among age 9 samples, the mean estimated age from the pediatric clock
146 was 9 and the mean estimated age from the GRIM clock was 25. Among age 15 samples, the
147 mean estimated age from the pediatric clock was 12 and the mean estimated age from the GRIM
148 clock was 30. The distribution of estimated cell-type proportions also shifted between visits
149 while the distribution of global methylation and polymethylation scores were more consistent
150 (Supplemental Figure 3).

151 In bivariate analyses and multivariable models, we focused on the polymethylation score for
152 prenatal smoke exposure and the single top CpG site from prior research, cg05575921 in the
153 *AHRR* gene, as hypothesized biomarkers, and used global methylation and the pediatric clock as
154 negative controls.

155 **Bivariate associations between prenatal maternal smoking and DNA methylation summary**
156 **measures**

157 Mothers who reported smoking during pregnancy had lower income-to-poverty ratios (1.49) than
158 those who did not (2.38). Mothers who smoked were more likely to report prenatal alcohol use
159 (30% vs 7%), prenatal drug use (18% vs 2%) and postnatal smoking (96% vs 26%) than mothers
160 who did not report prenatal maternal smoking (Table 1). Children of mothers who reported
161 smoking during pregnancy were more likely to be of European genetic ancestry (23%) than
162 children of mothers who did not (14%, P value= <0.001). Children of mothers who reported
163 smoking during pregnancy had higher prenatal maternal smoking polymethylation scores than
164 children of mothers who did not at both the age 9 (mean-centered scores of 0.08 vs -0.04, P
165 value <0.001) and age 15 (0.11 vs -0.01, P value <0.001) visits (Table 1, Supplemental Table 2).
166 At age 9, children exposed to prenatal maternal smoking had lower DNA methylation at
167 cg05575921 (76.82%) than children of non-smoking mothers (77.81%, P value=0.05). At age
168 15, children exposed to prenatal smoking had lower DNA methylation at cg05575921 (76.13%)
169 than children of non-smoking mothers, although this difference was not significant (76.82% P
170 value=0.24). Epigenetic age from the pediatric and GRIM clocks and global DNA methylation
171 did not differ between children exposed vs unexposed to prenatal maternal smoking.

172 **Multivariable associations between prenatal maternal smoking and DNA methylation**
173 **summary measures**

174 The association between prenatal maternal smoking and the polymethylation score for prenatal
175 smoke exposure was also observed in multivariable models, adjusting for base model covariates
176 of child sex, maternal income to poverty ratio at baseline, indicator variable for city of residence

177 being Detroit, Chicago or Toledo, proportion of salivary immune cells, sample plate from DNA
178 methylation analysis, child age and the first two principal components of genetic ancestry
179 (Figure 3; Supplemental Table 3). At age 9, prenatal maternal smoke exposure was associated
180 with a 0.52 (95% CI: 0.36, 0.67) standard deviation higher polymethylation score for prenatal
181 smoke exposure. At age 15, prenatal maternal smoke exposure was associated with a 0.46
182 (95% CI: 0.3, 0.62) standard deviation higher prenatal smoke exposure polymethylation score for
183 prenatal smoke exposure. A consistent association was observed when stratifying by genetic
184 ancestry. In the African genetic ancestry sample (n= 488) at age 9, prenatal maternal smoking
185 was associated with a 0.55 (95% CI: 0.35, 0.75) standard deviation higher polymethylation score
186 for prenatal smoke exposure. The direction of the association was similar in the European and
187 Hispanic genetic ancestry samples, though the effect size was slightly attenuated and no longer
188 statistically significant (Figure 3; Supplemental Table 4).

189 At age 9, prenatal maternal smoking was associated with 1 percent lower DNA methylation at
190 cg055975921 (95% CI: -1.65, -0.34) after adjusting for base model covariates. Similarly, at age
191 15 prenatal maternal smoking was associated with 0.8 percent lower DNA methylation at
192 cg055975921 (95% CI: (-1.58, -0.02)). Prenatal maternal smoking remained significantly
193 associated with a decrease in cg05575921 DNA methylation in the African genetic ancestry
194 sample at age 9 (-1.3 (95% CI: -2.15, -0.46)) although the association was not significant at age
195 15 (-0.97 (95% CI: -1.98, 0.04); Figure 3). In the European genetic ancestry sample, the
196 association was consistent in direction at age 9 (-0.81 (95% CI: -2.45, 0.82)) and age 15 (-1.79
197 (95% CI: -3.43, -0.15)). In the Hispanic genetic ancestry sample, the association was no longer
198 significant and not consistent in direction (Age 9: -0.05 (95% CI: -1.69, 1.59); Age 15: 0.82
199 (95% CI: -1.28, 2.92)) (Figure 3).

200 Prenatal maternal smoking was not associated with global DNA methylation or epigenetic age
201 (Figure 3).

202 **Sensitivity analyses for association testing**

203 After additionally adjusting for other prenatal exposures and postnatal smoke exposure, the
204 association between prenatal maternal smoking and the polymethylation score was robust
205 (Supplemental Figure 4). Similarly, after adjusting for other prenatal exposures and postnatal
206 smoke exposure, the direction of the association between prenatal maternal smoking and
207 cg055975921 was consistent, although no longer significant at age 15. Very few children who
208 were exposed to maternal smoking prenatally were unexposed to postnatal smoke (Supplemental
209 Figure 5). However, children exposed only to postnatal smoke did not have higher
210 polymethylation scores than children unexposed to both postnatal and prenatal smoke
211 (Supplemental Figure 5). Results from linear mixed effect models were similar to age-stratified
212 and base covariate adjusted linear models (Supplemental Tables 3 & 4). Results were also
213 similar when controlling for surrogate variables instead of known covariates (Supplemental
214 Figures 4 and 6; Supplemental Table 5).

215 Results were also similar when using polymethylation scores constructed with alternative
216 regression coefficients (see Methods, Supplemental Figure 7). For our main analysis we used
217 coefficients from a regression of sustained smoking exposure and DNA methylation in newborn
218 cord blood with cell-type control. As sensitivity analyses, we used coefficients from regressions
219 of: sustained smoking exposure and DNA methylation in newborn cord blood without cell-type
220 control, sustained smoking exposure and DNA methylation in peripheral blood from older
221 children without cell-type control, and any smoking exposure and DNA methylation in newborn

222 cord blood without cell-type control [6]. The association between prenatal maternal smoking and
223 the polymethylation score using regression coefficients from newborn cord blood with cell-type
224 controls was the strongest.

225 In addition to cg055975921 in the *AHHR* gene, we tested the association of four other *a priori*
226 probe sites identified from previous meta-analyses. We replicated the association of prenatal
227 maternal smoking with these CpG sites (Supplemental Table 2, Supplemental Figure 8).

228 **Accuracy of polymethylation scores as a biomarker of prenatal maternal smoking**

229 Next, we compared the accuracy of different DNA methylation summary measures for
230 classifying prenatal maternal smoking using receiver operating curves (Figure 4A; Supplemental
231 Table 6). We estimated classification of prenatal maternal smoking when using DNA
232 methylation summary measures alone. We also estimated prenatal maternal smoking
233 classification when using DNA methylation summary measures in addition to base model
234 covariates (child sex, maternal income to poverty ratio at baseline, child age at DNA methylation
235 measurement, indicator variable for city of residence being Detroit, Chicago or Toledo,
236 estimated immune cell proportion, plate from methylation processing and the first two genetic
237 principal components). When used without DNA methylation measures, these base model
238 variables had an area under the curve (AUC) of 0.73 at age 9 and 0.72 at age 15.

239 At age 9, including the polymethylation score for prenatal smoke exposure significantly
240 improved the base model covariates classification (AUC: 0.78, *P* value comparing to base
241 model<0.001). At age 15, similarly, including the polymethylation score for prenatal smoke
242 exposure significantly improved the base model covariates classification (AUC:0.77, *P*
243 value<0.001).

244 At age 9, the base model covariates with the polymethylation score for prenatal smoke exposure
245 had a larger AUC than the base model covariates with cg05575921 (AUC base model covariates
246 + cg05575921: 0.74; P value = 0.02). At age 15, this was also true (AUC for base model
247 covariates + cg05575921 = 0.73, P value = 0.01; Supplemental Table 6). Accurate classification
248 was not improved by including polymethylation scores from age 9 and age 15 together (Figure
249 4B).

250 **Sensitivity analyses for biomarker accuracy assessment**

251 Classification of prenatal maternal smoke exposure when using other coefficients as the weights
252 in construction of the polymethylation scores was similar to the results when using coefficients
253 from cell-type controlled regressions of sustained prenatal maternal smoking and DNA
254 methylation in cord blood (Supplemental Table 7).

255 **Discussion**

256 In the longitudinal Fragile Families and Child Wellbeing birth cohort, we observed that prenatal
257 maternal smoking was associated with several characterizations of DNA methylation in
258 children's saliva samples from ages 9 and 15. Prenatal maternal smoking was associated with
259 polymethylation scores for prenatal smoke exposure across strata of both child age and genetic
260 ancestry. Global methylation and epigenetic clocks were not associated with maternal smoking
261 exposure. Polymethylation scores for prenatal smoke exposure had reasonable accuracy for
262 classifying prenatal maternal smoking (AUC > 0.7). Classification when using polymethylation
263 scores for prenatal maternal smoke exposure was better than when using a single *a priori* CpG
264 site, cg05575921 in the *AHRR* gene.

265 Our findings are consistent with the previous literature on associations between prenatal maternal
266 smoking and DNA methylation. We replicated the top hit from a previous epigenome wide
267 association analysis [6]. Further, we found evidence of association in several additional hits from
268 epigenome wide association studies of DNA methylation in cord and peripheral blood [6, 26].
269 There are no previous studies of prenatal maternal smoking and saliva DNA methylation.
270 However, previous work has shown that the majority of CpG sites are similarly methylated in
271 blood and saliva (reviewed in [27]).

272 We advance prenatal smoking - DNA methylation literature by evaluating the persistence of the
273 association between prenatal maternal smoking and DNA methylation as children age and its
274 portability across tissue and ancestry. Polymethylation scores built using coefficients from meta-
275 analysis of cord blood DNA methylation from primarily European-ancestry newborns were still
276 associated with prenatal maternal smoking in our independent saliva samples from a diverse
277 cohort at ages 9 and 15 [6]. The portability of other risk scores, such as polygenic risk scores,
278 across genetic ancestries is a complex research area [28] and the portability of epigenetic
279 summary measures has been identified as a key area for evaluation [29]. In this case, the
280 polymethylation score for prenatal maternal smoking appears to be portable across genetic
281 ancestry groups. Effect estimates were consistent across ancestries at age 9, although Hispanic
282 and European genetic ancestry samples had lower, non-significant effect estimates at age 15. The
283 decreased precision may be the result of the reduced sample size of the Hispanic and European
284 genetic ancestry samples. The reduction in effect estimate magnitude could reflect higher
285 unreported smoking initiation in European and Hispanic ancestry teens in the United States than
286 African ancestry teens. Genetic ancestry correlates with race, and White and Latino teens have
287 much higher rates of teen smoking and earlier ages at initiation than Black children [30]. Though

288 we excluded children who reported own-smoking from our analytic sample, under-reporting of
289 own-smoking could influence DNA methylation at age 15, creating outcome misclassification.
290 This could also explain the observed attenuation of the effect on *AHRR*:cg05575921 methylation
291 by age 15 in our sample, as DNA methylation at this probe is known to vary by personal
292 cigarette smoking [22, 31].

293 Our findings are also consistent with previous work on DNA methylation as a biomarker for
294 prenatal maternal smoke exposure. Prenatal maternal smoking classification using saliva DNA
295 methylation in our sample of 9- and 15-year-olds performed comparably to classification using
296 peripheral blood from older adults (AUC 0.72) [10]. However, previous prenatal maternal smoke
297 exposure biomarkers created using LASSO regression and cord blood from newborns performed
298 better (AUCs ranging from 0.82-0.97) [20] A support vector machine approach performed on
299 peripheral blood from 3–5-year-old children also more accurately classified prenatal maternal
300 smoke exposure (AUC=0.87) [23]. Differences in DNA methylation patterns across tissues and
301 over time may influence the performance of DNA methylation biomarkers. Different methods
302 for summarizing across multiple DNA methylation sites may also influence biomarker
303 performance. The similarity in performance between saliva DNA methylation in our study and
304 peripheral blood DNA methylation in older adults is encouraging for the use of saliva as a
305 readily accessible clinical sample.

306 Our analysis contributes to the development of an accurate methylation biomarker for prenatal
307 maternal smoking by evaluating the impact of specific methodological choices on accuracy of
308 classification. DNA methylation biomarkers may be especially susceptible to confounding by
309 subject age at methylation measurement and cell-type proportion [29]. Prenatal maternal smoke
310 exposure classification accuracy of polymethylation score was similar when using coefficients

311 from cord blood vs peripheral blood samples from older children. Classification accuracy was
312 not improved by using two time-points of methylation measurement. Classification accuracy was
313 improved when polymethylation scores used coefficients which incorporated cell-type control.
314 Our analysis therefore suggests that cell-type control when both generating (in the epigenome
315 wide association study) and applying coefficients for polymethylation scores may positively
316 influence outcome classification accuracy, although more work comparing coefficients from
317 different populations is still needed.

318 Our analysis suggests that polymethylation scores may be more accurate than using single CpG
319 sites as biomarkers. The site cg05575921 in the *AHRR* gene is a consistent marker of prenatal
320 maternal smoke exposure in meta-analyses of newborn cord blood [6]. DNA methylation at
321 cg05575291 in the *AHRR* gene can accurately classify own-smoking behavior in both blood
322 (area under the curve 0.995) and saliva (area under the curve 0.971) [22, 31]. However, salivary
323 *AHRR*:cg05575921 was not a persistent marker of prenatal maternal smoking in an analysis of
324 middle-aged adult women [10]. In our sample, salivary *AHRR*:cg05575921 methylation was less
325 accurate at predicting prenatal maternal smoke exposure than polymethylation scores. The
326 accuracy of *AHRR*:cg05575921 and other single CpG site biomarkers may be influenced by
327 time-since-exposure and new environmental exposures. Incorporating information across
328 multiple sites of DNA methylation may yield a biomarker more robust to these influences.

329 Our analysis is not without its limitations. We used maternal self-report of prenatal maternal
330 smoke exposure, as serum cotinine levels were not available. Due to social desirability bias, this
331 could result in exposure misclassification. We would expect this to bias our results towards the
332 null. In any analysis of a prenatal exposure and postnatal outcome, there is the possibility of
333 selection bias into the cohort due to live birth bias. Selection bias is also possible due to loss-to-

334 follow-up between birth and age 15. Additionally, while we controlled for postnatal secondhand
335 smoke exposure and excluded children who reported any own-smoking, we cannot exclude the
336 possibility of residual confounding in this observational cohort.

337 However, our analysis also has several strengths. We analyzed samples from a large cohort of
338 diverse participants underrepresented in genetic and epigenetic research [32, 33]. While our
339 exposure measurement of prenatal maternal smoking was self-reported, it was assessed
340 prospectively and preceded outcome measurements. We analyzed repeated measures of DNA
341 methylation with reproducible array measures conducted in a single batch. We tested
342 associations between prenatal maternal smoking and multiple DNA summary measures to
343 evaluate the specificity of the polymethylation scores. In sensitivity analyses, we adjusted for
344 other prenatal exposures and postnatal smoke exposure to examine the specificity of the
345 biomarker to nature and timing of exposure.

346 **Conclusions**

347 In a large, prospective study of diverse participants, we showed that DNA methylation in
348 children's saliva had strong associations with and reasonable classification accuracy for prenatal
349 maternal smoke exposure. Further, we demonstrated that polymethylation scores could be
350 applied as a biomarker of prenatal maternal smoke exposure across genetic ancestry groups, an
351 important consideration for the equitable biomarker development. The development and
352 application of biomarkers for prenatal maternal smoke exposure has important implications for
353 epidemiological research and clinical practice. Given the difficulty of measuring prenatal
354 maternal smoke exposure, such a biomarker could allow for confounder control in research areas
355 where such control is currently impossible. Prenatal maternal smoke exposure is prevalent and

356 has negative health consequences, thus an exposure biomarker could be used to provide support
357 and health interventions for children.

358 **Methods**

359 **Cohort**

360 The Fragile Families and Child Wellbeing Study is a birth cohort of nearly 5,000 children born
361 in 20 cities in the United States between 1998 and 2000 [34]. Participants were selected at
362 delivery using a three-stage stratified random sample design which oversampled unmarried
363 mothers by a ratio of 3:1 [34]. Participants were excluded on the following criteria: those with
364 parents who planned to place the child for adoption, those where the father was deceased, those
365 who did not speak English or Spanish well enough to be interviewed, births where the mothers or
366 babies were too ill to complete the interview, and those where the baby died before the interview
367 could take place. Children were followed longitudinally with assessments at ages 1, 3, 5, 9 and
368 15; additional follow up is ongoing. Assessments included medical record extraction, biosample
369 collection, in-home assessments, and surveys of the mother, father, primary caregiver, teacher
370 and child. At ages nine and fifteen a saliva sample was taken from the child [34]. A subsample of
371 the Fragile Families cohort was selected for saliva DNA and DNA methylation processing. To be
372 eligible for the DNA methylation assessment children had to have participated and given saliva
373 at age 9 and 15; after which a random sample of these eligible participants was then selected.
374 One exception to this rule is that all participants in the Study of Adolescent Neurodevelopment
375 (SAND) were assayed even if they only provided a sample at one time-point. All analyses
376 account for this oversample of these participants.

377 **Covariates and exposure measurement**

378 Demographic and prenatal maternal substance use variables were derived from maternal self-
379 report questionnaire data at baseline (child's birth).

380 Maternal covariates included maternal income to poverty ratio at baseline, prenatal smoking,
381 prenatal maternal drug and alcohol use, and postnatal maternal or primary caregiver smoking.

382 Maternal income to poverty ratio is a constructed variable of the ratio of total household income
383 (as self-reported by the mother) to the official poverty thresholds designated by the United States

384 Census Bureau for the year preceding the interview. At baseline, mother's answered categorical
385 questions about their prenatal smoking, drug, and alcohol use. For maternal prenatal smoking,

386 mothers were asked

387 During your pregnancy, how many cigarettes did you smoke? Did you smoke...:

388 2 or more packs a day

389 1 or more but less than 2 packs per day

390 Less than 1 pack a day

391 None

392 Few participants reported smoking a pack or more a day, thus we dichotomized to any vs no
393 prenatal maternal smoking. For maternal prenatal drug and alcohol use the mothers were asked

394 During your pregnancy how often did you use drugs/drink alcohol (respectively):

395 Never

396 Less than 1 time per month

397 Several times per month

398 Several times per week

399 Every day

400 When the child was 1, 5, 9 and 15 years of age, primary caregivers responded to questions about
401 maternal and in-home smoking. To encapsulate general early childhood smoke exposure, we
402 created a binary variable for postnatal exposure at ages 1 or 5. To encapsulate recent postnatal
403 smoke exposure, we used a categorical variable for packs per day (no smoking, less than one
404 pack/day, one or more packs/day) in the month prior to the age 9 and 15 interview. To control for
405 the oversampling of the Detroit, Toledo and Chicago locations, an indicator variable denoting
406 residence in one of these three cities was created.

407 Child covariates included child sex, child report of personal cigarette smoking, and child genetic
408 ancestry. Mothers reported sex (male/female) of their child at baseline. At ages 9 children were
409 asked if they had ever smoked a cigarette or used tobacco (yes/no) and at age 15 they were asked
410 if they had ever smoked an entire cigarette (yes/no).

411 Child genetic ancestry was calculated. Grants R01 HD36916, R01 HD073352, and R01
412 HD076592 provided support for the collection, assay, and analysis of the genetic data in the
413 Fragile Families cohort. Principal components (PC) of child genetic ancestry were calculated
414 from genetic data using Illumina PsychChip_v1-1 with child saliva samples. Genetic ancestry
415 was assigned by comparing PC loadings to 1000 Genomes super population clusters. Samples
416 with $PC1 > 0.018$ and $PC2 > -0.0075$ were assigned to European ancestry. Samples with $PC1 < -$
417 0.005 and $PC2 > 0.007 + 0.75(PC1)$ were assigned to African ancestry. Samples with $PC1 > 0.018$
418 and $-0.055 < PC2 < 0.025$ were assigned to Hispanic ancestry [35].

419 DNA methylation measurement

420 Salivary samples from the children were collected at ages nine and fifteen using the
421 Oragene•DNA sample collection kit (DNA Genotek Inc., Ontario). Saliva DNA was extracted
422 manually using DNA Genotek's purification protocol using prepIT L2P. DNA was bisulfite
423 treated and cleaned using the EZ DNA Methylation kit (Zymo Research, California). Samples
424 were randomized and plated across slides by demographic characteristics. Saliva DNA
425 methylation was measured using the Illumina HumanMethylation 450k BeadArray [36] and
426 imaged using the Illumina iScan system. All samples were run in a single batch to minimize
427 technical variability.

428 DNA methylation image data were processed in R statistical software (3.5) using the minfi
429 package [37]. The red and green image pairs (n=1811) were read into R and the minfi
430 *preprocessNoob* function was used to normalize dye bias and apply background correction.
431 Further quality control was applied using the ewastools packages [38]. We dropped samples with
432 >10% of sites have detection p-value >0.01 (n=43), sex discordance between DNA methylation
433 predicted sex and recorded sex (n=20), and abnormal sex chromosome intensity (n=3). CpG sites
434 were removed if they had detection p-value >0.01 in 5% of samples (n=26,830) or were
435 identified as cross-reactive (n=27,782) (Figure 1) [39]. We used the November 2021 data freeze
436 of the Fragile Families and Child Wellbeing DNA methylation data. Relative proportions of
437 immune and epithelial cell types were estimated from DNA methylation measures using a
438 childhood saliva reference panel [40].

439 We created polymethylation scores for prenatal maternal smoke exposure. From an independent
440 meta-analysis of prenatal smoke exposure and newborn DNA methylation, we extracted the
441 regression coefficients of 6,074 CpG sites associated with prenatal maternal smoking at a false

442 discovery ratio corrected P value <0.05 [6]. We mean-centered the DNA methylation beta values
443 in our study, weighted them by the independent regression coefficients and took the sum. We
444 calculated polymethylation scores using regression coefficients from 4 different regressions. For
445 our main analysis we used coefficients from a regression of sustained smoking exposure and
446 DNA methylation in newborn cord blood with cell-type control. As sensitivity analyses, we used
447 sustained smoking exposure and DNA methylation in newborn cord blood without cell-type
448 control, sustained smoking exposure and DNA methylation in peripheral blood from older
449 children, any smoking exposure and DNA methylation in newborn cord blood) [6].

450 Global DNA methylation was calculated for each sample as the mean methylation value of each
451 sample across the cleaned probe set. Mean DNA methylation restricted to probes in genomic
452 regions (CpG island, shore, shelf or open sea, as identified in the R package
453 `IlluminaHumanMethylation450kanno.ilmn12.hg19 v 0.6.0`) was also calculated.

454 Pediatric epigenetic age was calculated for each sample using the coefficients and methods
455 provided by the creators (see <https://github.com/kobor-lab/Public-Scripts>) [41]. The GRIM
456 age clock, including the smoking pack-years sub-scale, was calculated as previously described
457 [42].

458 Single *a priori* CpG sites were selected from previous large meta-analyses of prenatal maternal
459 smoking and DNA methylation in children's cord and peripheral blood samples [6, 26].

460 Surrogate variables were calculated from DNA methylation data using the function `sva` from the
461 R package `sva` version 3.38.

462 **Statistical analyses**

463 *Inclusion exclusion criteria*

464 From the 1685 samples from 882 unique individuals with quality-controlled DNA methylation
465 data, we further excluded any samples missing data on: maternal prenatal smoking (0 samples),
466 alcohol (4 samples) or other drug use (6 samples), maternal income to poverty ratio (0 samples),
467 maternal postnatal smoking data (105 samples). We also excluded samples with missing child
468 sex (0 samples), child age (0 samples) or genetic data (22 samples). Finally, if a child reported
469 ever smoking a cigarette or using tobacco age 9, we excluded all of their available samples. Age
470 15 samples from children who reported ever smoking a cigarette at age 15 were also excluded.
471 Children who were missing a response to the question at age 9 but answered that they had never
472 smoked a whole cigarette at age 15 were kept in the sample.

473 *Base model*

474 In our base model we adjusted for maternal income to poverty ratio at baseline, child sex, child
475 age, an indicator variable for residence in Detroit, Chicago or Toledo, plate from DNA
476 methylation processing, and estimated immune cell proportion estimated from DNA methylation.
477 In nonstratified models, we adjusted for the first two components of genetic ancestry from
478 principal component analysis. In ancestry-stratified models, the first two principal components
479 from principal component analysis run within each ancestry strata were used. While child sex,
480 child age, plate from DNA methylation processing and immune cell proportions are not
481 confounders (as they cannot casually affect prenatal maternal smoke exposure), these variables
482 can strongly affect DNA methylation and so were adjusted for as precision variables.

483 *Receiver operator curve analysis*

484 To evaluate the accuracy of the DNA methylation summary measures as biomarkers of prenatal
485 maternal smoking, we used a receiver operating curve. First, we regressed exposure to prenatal
486 maternal smoking (outcome) against the DNA methylation summary measures in individual
487 logistic regressions, while adjusting for the base model variables listed above. Next, we
488 calculated receiver operating curves (ROC) and area under the curves (AUCs) using the function
489 roc from the R library pROC version 1.18.0. We compared ROC curves and AUC using the
490 Delong method and the function roc.test from the R library pROC version 1.18.0.

491 *Sensitivity analyses*

492 In sensitivity analyses we performed additional adjustments: 1) models adjusting for prenatal
493 drug and alcohol use, 2) models adjusting for prenatal drug and alcohol use and postnatal
494 maternal/primary caregiver smoking, 3) models adjusting for surrogate variables calculated from
495 DNA methylation data.

496 In addition to the cross-sectional models within visit-strata, we performed a linear mixed effect
497 model with a random intercept for child.

498 Code to perform all analyses is available (www.github.com/bakulskilab)

499 **Declarations**

500 **Ethics approval and consent to participate**

501 Participants provided written informed consent for the study. The data used in this manuscript
502 were prepared by the Fragile Families and Childhood Wellbeing Study administrators following

503 approval of the manuscript proposal. These secondary data analyses were approved by the
504 University of Michigan Institutional Review Board (IRB, HUM00129826)

505 **Consent for publication**

506 Not applicable

507 **Availability of data and materials**

508 Many of the variables used in this analysis are publicly available in the Fragile Families dataset.
509 Some of the data used in this analysis, including genetic and epigenetic data, is restricted use but
510 available to researchers upon reasonable application.

511 **Competing interests**

512 The authors have no competing interests to declare.

513 **Funding**

514 This research was made possible through several grants (Fragile Families Core Data Collection:
515 R01 HD036916, genetic data processing R01 HD036916, R01 HD073352, R01 HD076592 and
516 R01 MD011716). In addition, FB was supported by the National Institutes of Health National
517 Institute of Dental and Craniofacial Research (F31 DE029992), EB was supported by R01
518 AG055406 and R01 AG067592, and CM was supported by R01 AG071071. KB and JD were
519 supported by R01 AG067592 (MPI: Bakulski, Ware), R01 ES025531 (PI: Fallin), R01
520 ES025574 (PI: Schmidt), and R01 MD013299 (PI: Hicken).

521 **Authors' contributions**

522 FB performed the analysis with contributions from JD, JF, and EW. FB and KB wrote the paper
523 with all authors contributing to revisions. CM, DN and LS contributed to the design and
524 processing of the DNA methylation data subcohort. DN, CM, KM, EW, and FB all contributed
525 to the design and conception of the analytic plan. All authors reviewed and revised the
526 manuscript.

527 **Acknowledgements**

528 Not applicable

529

530 *Table 1: Bivariate associations between prenatal maternal smoke exposure and selected DNA*
531 *methylation summary measures and important covariates among a diverse sample of 809*
532 *children in the Fragile Families and Child Wellbeing study*

Characteristic

Age 9

Age 15

	N	Not exposed, N = 598 ^f	Smoke exposed, N = 155 ^f	p-value ²	N	Not exposed, N = 603 ^f	Smoke exposed, N = 143 ^f	p-value
Child characteristics								
Polymethylation score for smoke exposure*	753	-0.04 (0.22)	0.08 (0.22)	<0.001	746	-0.01 (0.22)	0.11 (0.23)	<0.001
<i>AHRR</i> gene: percent cg05575921 methylation	753	77.8 (5.2)	76.8 (5.6)	0.047	746	77 (6)	76 (6)	0.2
Pediatric epigenetic clock (years)	753	8.98 (0.75)	8.94 (0.83)	0.6	746	11.62 (1.75)	11.59 (1.70)	0.9
Percent global DNA methylation	753	47.51 (0.74)	47.59 (0.78)	0.3	746	47.48 (0.85)	47.58 (0.83)	0.2
Immune cell proportion (saliva)	753	0.96 (0.10)	0.96 (0.08)	0.4	746	0.92 (0.15)	0.92 (0.14)	>0.9
Epithelial cell proportion (saliva)	753	0.04 (0.10)	0.04 (0.08)	0.4	746	0.08 (0.15)	0.08 (0.14)	>0.9
Ancestry categorization from child principal components of genetic data	753			0.001	746			0.001
African ancestry		355 (59%)	97 (63%)			364 (60%)	88 (62%)	
European ancestry		82 (14%)	35 (23%)			84 (14%)	34 (24%)	
Hispanic ancestry		161 (27%)	23 (15%)			155 (26%)	21 (15%)	
Child gender	753			0.8	746			0.9

Characteristic	N	Age 9			Age 15			p-value ²
		Not exposed, N = 598 ¹	Smoke exposed, N = 155 ¹		Not exposed, N = 603 ¹	Smoke exposed, N = 143 ¹	p-value	
Boy		300 (50%)	76 (49%)		295 (49%)	69 (48%)		
Girl		298 (50%)	79 (51%)		308 (51%)	74 (52%)		
Maternal characteristics								
Maternal prenatal alcohol use	753	44 (7.4%)	47 (30%)	<0.001	746	40 (6.6%)	42 (29%)	<0.00
Maternal prenatal any drug use	753	10 (1.7%)	26 (17%)	<0.001	746	9 (1.5%)	28 (20%)	<0.00
Maternal income to poverty threshold (at baseline)	753	2.41 (2.57)	1.52 (1.56)	<0.001	746	2.41 (2.58)	1.49 (1.57)	<0.00
Maternal city of residence (at baseline)	753			0.12	746			0.6
Detroit, Chicago, or Toledo		148 (25%)	48 (31%)		171 (28%)	44 (31%)		
Not Detroit, Chicago, or Toledo		450 (75%)	107 (69%)		432 (72%)	99 (69%)		
Postnatal maternal smoking at ages 1 or 5	753			<0.001	746			<0.00
No maternal smoking at age 1 and 5		447 (75%)	7 (4.5%)		448 (74%)	7 (4.9%)		
Maternal smoking at age 1 or age 5		151 (25%)	148 (95%)		155 (26%)	136 (95%)		

Characteristic	Age 9				Age 15			
	N	Not exposed, N = 598 ^l	Smoke exposed, N = 155 ^l	p-value ²	N	Not exposed, N = 603 ^l	Smoke exposed, N = 143 ^l	p-value
Postnatal primary caregiver smoking in past month prior to visit	753			<0.001	746			<0.00
No smoking		476 (80%)	23 (15%)			511 (85%)	44 (31%)	
Less than pack a day		95 (16%)	85 (55%)			80 (13%)	77 (54%)	
Pack or more a day		27 (4.5%)	47 (30%)			12 (2.0%)	22 (15%)	

Mean (SD); n (%)

Welch Two Sample t-test; Pearson's Chi-squared test

Poymethylation score constructed using regression coefficients from a model of prenatal maternal smoking and DNA methylation in newborn cord blood as weights. DNA methylat eta values from Fragile Families and Child Wellbeing study were mean-centered, weighted by the regression coefficients and summed.

533

534

535

536

537

538 *Figure 1 - Selection of samples from the Fragile Families and Child Wellbeing study into*
539 *analytic subset. N represents the number of individuals at each step in the selection procedure,*
540 *M represents the number of samples. Individuals with repeated measures can have more than*
541 *one sample.*

542 *Figure 2 - Differences in selected DNA methylation summary measures by self-report of*
543 *prenatal maternal smoking among 809 children in the Fragile Families and Child Wellbeing*
544 *study at ages 9 and 15. Samples from children exposed to prenatal maternal smoke in grey,*
545 *samples from children unexposed to prenatal maternal smoke in black. From top-left, clockwise:*
546 *Polymethylation scores for prenatal maternal smoke exposure, constructed using regression*
547 *coefficients for prenatal smoke exposure predicting DNA methylation in newborn cord blood*
548 *samples, accounting for cell-type control. DNA methylation values from samples in the Fragile*
549 *Families and Child Wellbeing study were mean-centered, then multiplied by these regression*
550 *coefficients and summed. Pediatric epigenetic clock (years). AHRR gene: percent cg05575921*
551 *methylation. Percent global DNA methylation. **** P value<0.0001 *P value<0.05, ns P*
552 *value>0.05*

553 *Figure 3 - Prenatal maternal smoke exposure is consistently associated with polymethylation*
554 *scores at ages 9 and 15 and is portable across genetic ancestry groups in a sample of 809*
555 *children in the Fragile Families and Child Wellbeing study All models shown controlled for:*
556 *first two principal components of child genetic ancestry (from ancestry-stratified principal*
557 *components for ancestry stratified models), child sex, child age (months), maternal income-to-*
558 *poverty ratio at birth, indicator variable for residence in Detroit, Chicago, or Toledo, immune*
559 *cell proportion estimated from methylation data, yes/no other maternal prenatal drug use, yes/no*
560 *maternal prenatal alcohol use, postnatal maternal smoking when child 1 or 5 years of age,*
561 *postnatal maternal/primary care give smoking packs/day in month prior to saliva sample.*
562 *Poylmethylation score constructed using regression coefficients from a model of prenatal*
563 *maternal smoking and DNA methylation in newborn cord blood as weights. DNA methylation*
564 *beta values from Fragile Families and Child Wellbeing study were mean-centered, weighted by*
565 *the regression coefficients and summed. The resulting scores were z-score standardized and used*
566 *as outcomes in these models. Red dotted line denotes the null.*

567 *Figure 4 - Polymethylation scores accurately classify prenatal maternal smoke exposure at*
568 *ages 9 and 15 among 809 children in the Fragile Families and Child Wellbeing study A)*
569 *Receiver operator curve for select DNA methylation measures for predicting prenatal smoke*
570 *exposure using no other variables (light colors) or using base model variables (dark colors,*
571 *other variables included: child sex, child age, maternal income-poverty ratio at birth, indicator*
572 *variable for residence in Detroit, Chicago, or Toledo, immune cell proportion and batch of*
573 *methylation data processing). B) Receiver operator curve for including polymethylation scores*
574 *individually at each visit (black & light grey) or jointly at both visits (dark grey).*
575 *Poymethylation score constructed using regression coefficients from a model of prenatal*
576 *maternal smoking and DNA methylation in newborn cord blood as weights. DNA methylation*
577 *beta values from Fragile Families and Child Wellbeing study were mean-centered, weighted by*
578 *the regression coefficients and summed. The resulting scores were z-score standardized*

579

580 **References**

- 581 1. Banderali G, Martelli A, Landi M, Moretti F, Betti F, Radaelli G, Lassandro C, Verduci E
582 (2015) Short and long term health effects of parental tobacco smoking during pregnancy and
583 lactation: A descriptive review. *J Transl Med* 13:327
- 584 2. (2019) National survey on drug use and health 2018.
- 585 3. Kvalvik LG, Nilsen RM, Skjærven R, Vollset SE, Midttun O, Ueland PM, Haug K
586 (2012) Self-reported smoking status and plasma cotinine concentrations among pregnant women
587 in the norwegian mother and child cohort study. *Pediatr Res* 72:101–7

- 588 4. Dietz PM, Homa D, England LJ, Burley K, Tong VT, Dube SR, Bernert JT (2011)
589 Estimates of nondisclosure of cigarette smoking among pregnant and nonpregnant women of
590 reproductive age in the united states. *Am J Epidemiol* 173:355–9
- 591 5. Dempsey D, Jacob P 3rd, Benowitz NL (2002) Accelerated metabolism of nicotine and
592 cotinine in pregnant smokers. *J Pharmacol Exp Ther* 301:594–8
- 593 6. Joubert BR, Felix JF, Yousefi P, et al (2016) DNA methylation in newborns and maternal
594 smoking in pregnancy: Genome-wide consortium meta-analysis. *Am J Hum Genet* 98:680–696
- 595 7. Everson TM, Vives-Usano M, Seyve E, et al (2021) Placental DNA methylation
596 signatures of maternal smoking during pregnancy and potential impacts on fetal growth. *Nat*
597 *Commun* 12:5095
- 598 8. Cardenas A, Lutz SM, Everson TM, Perron P, Bouchard L, Hivert M-F (2019) Mediation
599 by Placental DNA Methylation of the Association of Prenatal Maternal Smoking and Birth
600 Weight. *American Journal of Epidemiology* 188:1878–1886
- 601 9. Tehranifar P, Wu H-C, McDonald JA, Jasmine F, Santella RM, Gurvich I, Flom JD,
602 Terry MB (2018) Maternal cigarette smoking during pregnancy and offspring DNA methylation
603 in midlife. *Epigenetics* 13:129–134
- 604 10. Richmond RC, Suderman M, Langdon R, Relton CL, Davey Smith G (2018) DNA
605 methylation as a marker for prenatal smoke exposure in adults. *Int J Epidemiol* 47:1120–1130
- 606 11. Galanter JM, Gignoux CR, Oh SS, et al (2017) Differential methylation between ethnic
607 sub-groups reflects the effect of genetic ancestry and environmental exposures. *eLife*.
608 <https://doi.org/10.7554/elife.20532>
- 609 12. Fraser HB, Lam LL, Neumann SM, Kobor MS (2012) Population-specificity of human
610 DNA methylation. *Genome Biology* 13:R8

- 611 13. Moen EL, Zhang X, Mu W, Delaney SM, Wing C, McQuade J, Myers J, Godley LA,
612 Dolan ME, Zhang W (2013) Genome-Wide Variation of Cytosine Modifications Between
613 European and African Populations and the Implications for Complex Traits. *Genetics* 194:987–
614 996
- 615 14. Rahmani E, Shenhav L, Schweiger R, et al (2017) Genome-wide methylation data mirror
616 ancestry information. *Epigenetics & Chromatin*. <https://doi.org/10.1186/s13072-016-0108-y>
- 617 15. Barcelona V, Huang Y, Brown K, Liu J, Zhao W, Yu M, Kardia SLR, Smith JA, Taylor
618 JY, Sun YV (2019) Novel DNA methylation sites associated with cigarette smoking among
619 African Americans. *Epigenetics* 14:383–391
- 620 16. Sun YV, Smith AK, Conneely KN, et al (2013) Epigenomic association analysis
621 identifies smoking-related DNA methylation sites in African Americans. *Human Genetics*
622 132:1027–1037
- 623 17. Philibert RA, Beach SRH, Lei M-K, Brody GH (2013) Changes in DNA methylation at
624 the aryl hydrocarbon receptor repressor may be a new biomarker for smoking. *Clinical*
625 *Epigenetics*. <https://doi.org/10.1186/1868-7083-5-19>
- 626 18. Philibert RA, Beach SRH, Brody GH (2012) Demethylation of the aryl hydrocarbon
627 receptor repressor as a biomarker for nascent smokers. *Epigenetics* 7:1331–1338
- 628 19. Xu R, Hong X, Zhang B, Huang W, Hou W, Wang G, Wang X, Igusa T, Liang L, Ji H
629 (2021) DNA methylation mediates the effect of maternal smoking on offspring birthweight: a
630 birth cohort study of multi-ethnic US mothernewborn pairs. *Clinical Epigenetics*.
631 <https://doi.org/10.1186/s13148-021-01032-6>
- 632 20. Ladd-Acosta C, Shu C, Lee BK, et al (2016) Presence of an epigenetic signature of
633 prenatal cigarette smoke exposure in childhood. *Environ Res* 144:139–148

- 634 21. Neophytou AM, Oh SS, Hu D, Huntsman S, Eng C, Rodriguez-Santana JR, Kumar R,
635 Balmes JR, Eisen EA, Burchard EG (2019) In utero tobacco smoke exposure, DNA methylation,
636 and asthma in latino children. *Environmental Epidemiology* 3:
- 637 22. Philibert R, Dogan M, Beach SRH, Mills JA, Long JD (2019) AHRR methylation
638 predicts smoking status and smoking intensity in both saliva and blood DNA. *American Journal*
639 *of Medical Genetics Part B: Neuropsychiatric Genetics* 183:51–60
- 640 23. Reese SE, Zhao S, Wu MC, et al (2017) DNA methylation score as a biomarker in
641 newborns for sustained maternal smoking during pregnancy. *Environ Health Perspect* 125:760–
642 766
- 643 24. Chen J, Zang Z, Braun U, et al (2020) Association of a Reproducible Epigenetic Risk
644 Profile for Schizophrenia With Brain Methylation and Function. *JAMA Psychiatry* 77:628
- 645 25. Bakulski KM, Fisher JD, Dou JF, Gard A, Schneper L, Notterman DA, Ware EB,
646 Mitchell C (2021) Prenatal Particulate Matter Exposure Is Associated with Saliva DNA
647 Methylation at Age 15: Applying Cumulative DNA Methylation Scores as an Exposure
648 Biomarker. *Toxics* 9:262
- 649 26. Wiklund P, Karhunen V, Richmond RC, et al (2019) DNA methylation links prenatal
650 smoking exposure to later life health outcomes in offspring. *Clinical Epigenetics*.
651 <https://doi.org/10.1186/s13148-019-0683-4>
- 652 27. Langie SAS, Moisse M, Declerck K, Koppen G, Godderis L, Vanden Berghe W, Drury
653 S, De Boever P (2017) Salivary DNA Methylation Profiling: Aspects to Consider for Biomarker
654 Identification. *Basic & Clinical Pharmacology & Toxicology* 121:93–101
- 655 28. Martin AR, Kanai M, Kamatani Y, Okada Y, Neale BM, Daly MJ (2019) Clinical use of
656 current polygenic risk scores may exacerbate health disparities. *Nature Genetics* 51:584–591

- 657 29. Hüls A, Czamara D (2019) Methodological challenges in constructing DNA methylation
658 risk scores. *Epigenetics* 15:1–11
- 659 30. El-Toukhy S, Sabado M, Choi K (2016) Trends in Susceptibility to Smoking by Race and
660 Ethnicity. *Pediatrics* 138:e20161254
- 661 31. Dawes K, Andersen A, Vercande K, Papworth E, Philibert W, Beach SRH, Gibbons FX,
662 Gerrard M, Philibert R (2019) Saliva DNA Methylation Detects Nascent Smoking in
663 Adolescents. *Journal of Child and Adolescent Psychopharmacology* 29:535–544
- 664 32. Sirugo G, Williams SM, Tishkoff SA (2019) The Missing Diversity in Human Genetic
665 Studies. *Cell* 177:26–31
- 666 33. Bentley AR, Callier S, Rotimi CN (2017) Diversity and inclusion in genomic research:
667 why the uneven progress? *Journal of Community Genetics* 8:255–266
- 668 34. Reichman NE, Teitler JO, Garfinkel I, McLanahan SS (2001) Fragile families: Sample
669 and design. *Children and Youth Services Review* 23:303–326
- 670 35. Ware EB, Fisher J, Schneper L, Notterman D, Mitchell C (2021) FFCWS polygenic
671 scores - release 1. Survey Research Center, Institute for Social Research, University of
672 Michigan; Department of Molecular Biology, Princeton University
- 673 36. Sandoval J, Heyn H, Moran S, Serra-Musach J, Pujana MA, Bibikova M, Esteller M
674 (2011) Validation of a DNA methylation microarray for 450,000 CpG sites in the human
675 genome. *Epigenetics* 6:692–702
- 676 37. Aryee MJ, Jaffe AE, Corrada-Bravo H, Ladd-Acosta C, Feinberg AP, Hansen KD,
677 Irizarry RA (2014) Minfi: a flexible and comprehensive Bioconductor package for the analysis
678 of Infinium DNA methylation microarrays. *Bioinformatics* 30:1363–1369

- 679 38. Heiss JA, Just AC (2018) Identifying mislabeled and contaminated DNA methylation
680 microarray data: an extended quality control toolset with examples from GEO. *Clinical*
681 *Epigenetics*. <https://doi.org/10.1186/s13148-018-0504-1>
- 682 39. Chen Y, Lemire M, Choufani S, Butcher DT, Grafodatskaya D, Zanke BW, Gallinger S,
683 Hudson TJ, Weksberg R (2013) Discovery of cross-reactive probes and polymorphic CpGs in the
684 Illumina Infinium HumanMethylation450 microarray. *Epigenetics* 8:203–209
- 685 40. Middleton LYM, Dou J, Fisher J, et al (2021) Saliva cell type DNA methylation
686 reference panel for epidemiological studies in children. *Epigenetics* 1–17
- 687 41. McEwen LM, O'Donnell KJ, McGill MG, et al (2019) The PedBE clock accurately
688 estimates DNA methylation age in pediatric buccal cells. *Proceedings of the National Academy*
689 *of Sciences* 117:23329–23335
- 690 42. Lu AT, Quach A, Wilson JG, et al (2019) DNA methylation GrimAge strongly predicts
691 lifespan and healthspan. *Aging* 11:303–327
- 692

Fragile Families cohort
n=4898

Cohort with 450K methylation data
n=897 individuals with m=1806 samples

Methylation data QC: m=43 samples & n=36 individuals with with >10% of probes failing dropped,
m=18 samples & n=13 individuals with with discordant sex samples dropped,
m=3 samples & n=2 individuals with with abnormal sex chromosome intensity dropped &
m=4 samples & n=4 individuals with with contamination/mislabeled &
53 duplicate samples dropped

Cohort with quality controlled 450K data
n=882 individuals with m=1685 samples

Missing maternal data:
m=0 samples & n=0 individuals with with missing prenatal maternal smoking
m=4 samples & n=2 individuals with with missing prenatal alcohol
m=6 samples & n=3 individuals with with missing prenatal other drug use
m=0 samples & n=0 individuals with with missing maternal income to poverty ratio data
m=88 samples & n=47 individuals with with missing postnatal maternal smoking data
m=17 samples & n=17 individuals with with missing past month maternal/primary care giver smoking packs/day

m=1570 samples & n=828 individuals with maternal covariate data

Missing child data:
m=22 samples & n=13 individuals with with missing ancestry/genetic data
m=0 samples & n=0 individuals with with missing child sex data
m=0 samples & n=0 individuals with with missing child age at sample

m=1548 samples & n=815 individuals with child covariate data

n=759 individuals with a 9 year visit

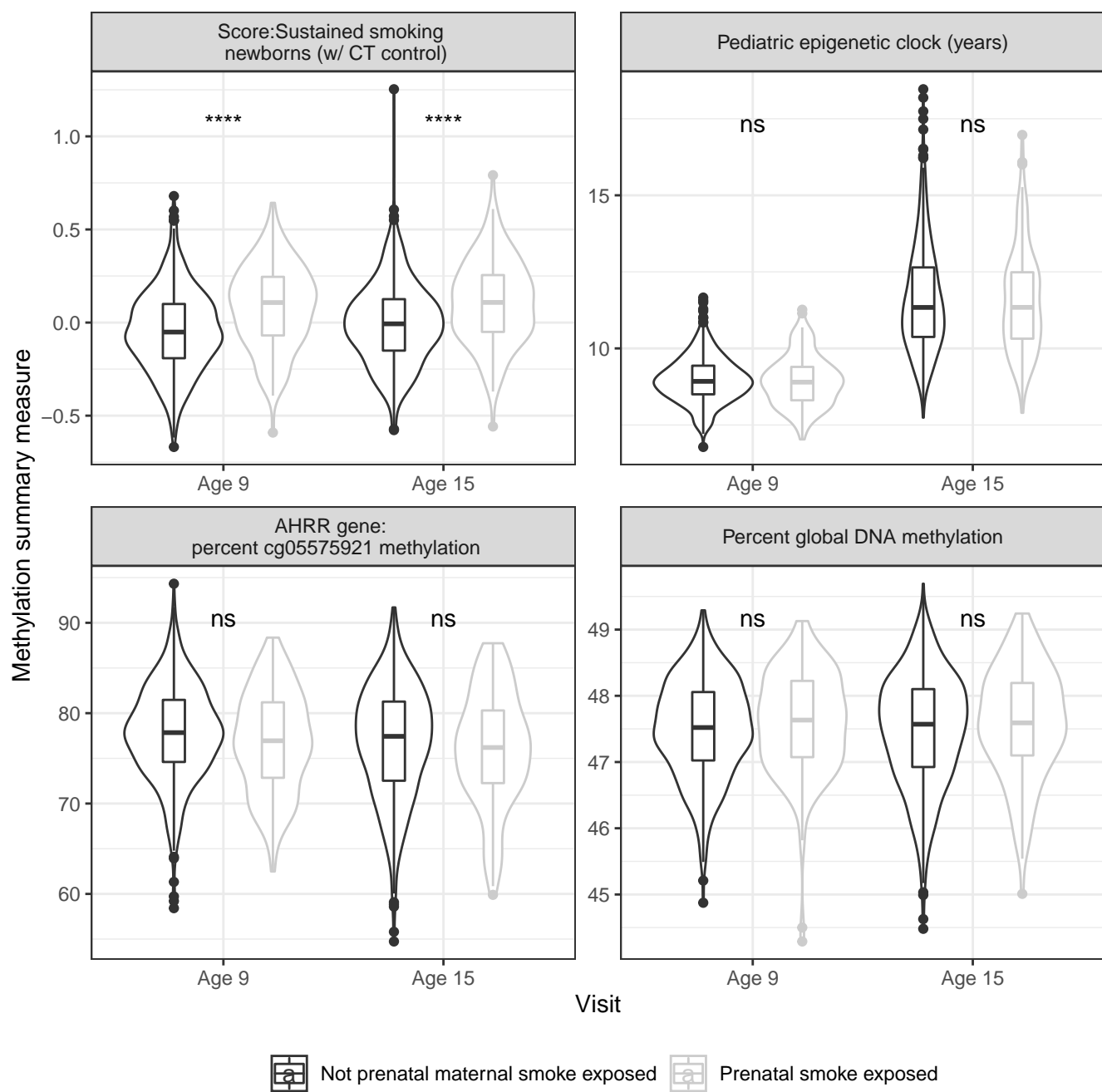
n=789 individuals with a 15 year visit

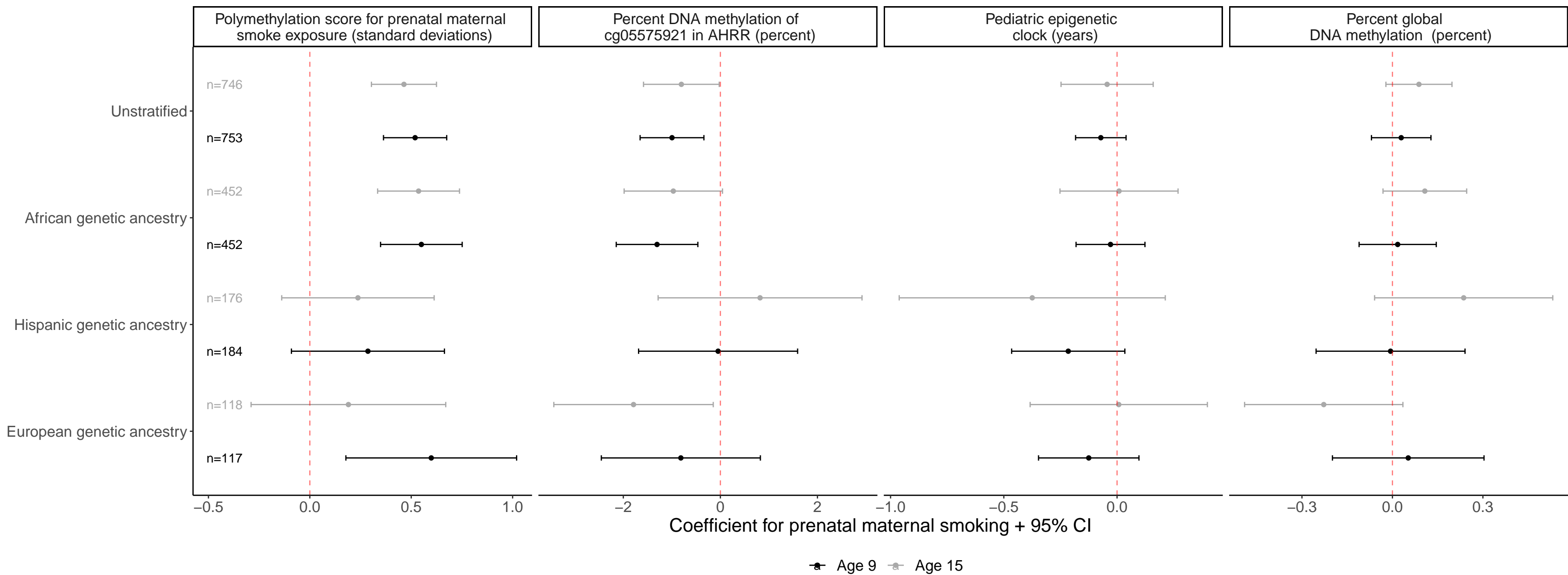
n=6 individuals smoking
or missing child smoking data
at 9 year visit

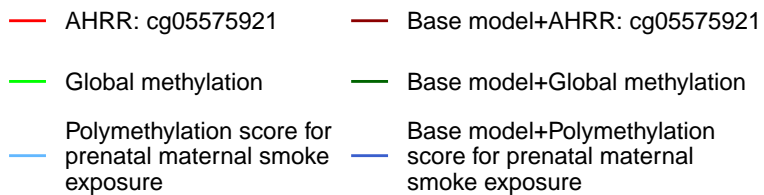
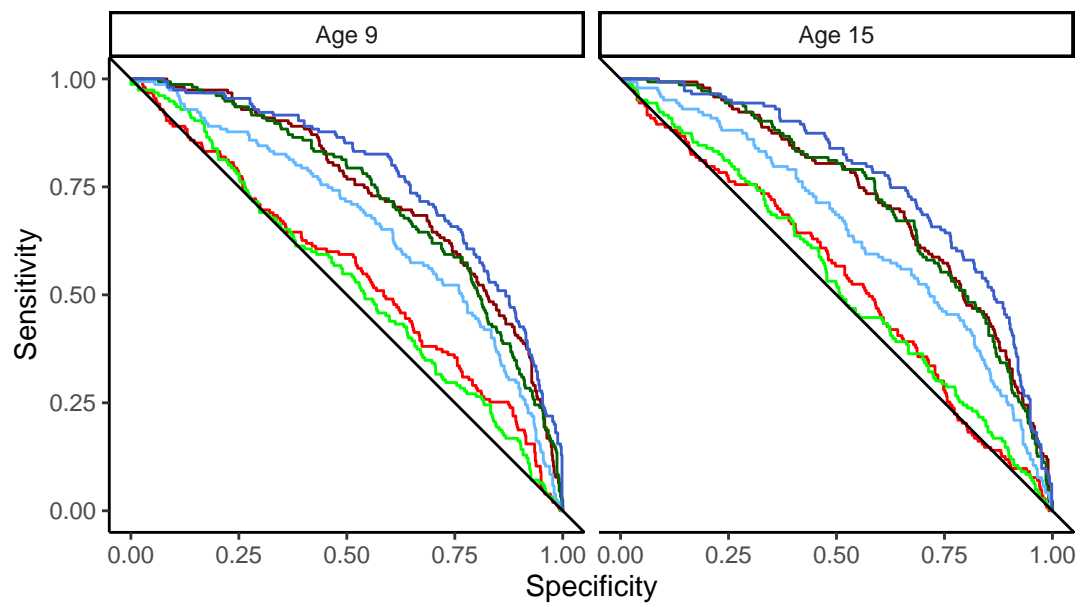
n=43 individuals smoking
or missing focal child smoking data
at 15 year visit

Exclude children who smoke:
n=753 individuals with a 9 year visit

Exclude children who smoke:
n=746 individuals with a 15 year visit





A**B**