

1 **The Immune Factors Driving DNA Methylation Variation in Human** 2 **Blood**

3
4 Jacob Bergstedt^{1,2,3,*}, Sadoune Ait Kaci Azzou¹, Kristin Tsuo¹, Anthony Jaquaniello¹, Alejandra
5 Urrutia⁴, Maxime Rotival¹, David T. S. Lin⁵, Julia L. MacIsaac⁵, Michael S. Kobor⁵, Matthew L.
6 Albert⁴, Darragh Duffy⁶, Etienne Patin^{1,8,*}, Lluís Quintana-Murci^{1,7,8,*}, for the Milieu Intérieur
7 Consortium

8
9 ¹Institut Pasteur, Université de Paris, CNRS UMR2000, Human Evolutionary Genetics Unit, Paris,
10 France

11 ²Institute of Environmental Medicine, Karolinska Institutet, Stockholm, Sweden

12 ³Department of Medical Epidemiology and Biostatistics, Karolinska Institutet, Stockholm, Sweden

13 ⁴HI-Bio, South San Francisco, CA, USA

14 ⁵Edwin S.H. Leong Healthy Aging Program, Centre for Molecular Medicine and Therapeutics,
15 Department of Medical Genetics, University of British Columbia, Vancouver, Canada

16 ⁶Institut Pasteur, Université de Paris, Translational Immunology Unit, Institut Pasteur, Paris, France

17 ⁷Chair of Human Genomics and Evolution, Collège de France, Paris, France

18 ⁸These authors contributed equally

19 *Correspondence: jacob.bergstedt@ki.se (J.B.), epatin@pasteur.fr (E.P.), quintana@pasteur.fr
20 (L.Q.-M.)

25 **Keywords:**

26 epigenetics, DNA methylation, humans, blood, immune cells, aging, cytomegalovirus, infection,
27 sex, meQTLs, EWAS, gene-by-environment interactions, cellular deconvolution

30 **Abstract**

31 Epigenetic changes are required for normal development, yet the nature and respective contribution
32 of factors that drive epigenetic variation in humans remain to be fully characterized. Here, we as-
33 sessed how the blood DNA methylome of 884 adults is affected by DNA sequence variation, age,
34 sex and 139 factors relating to life habits and immunity. Furthermore, we investigated whether these
35 effects are mediated or not by changes in cellular composition, measured by deep immunopheno-
36 typing. We show that DNA methylation differs substantially between naïve and memory T cells,
37 supporting the need for adjustment on these cell-types. By doing so, we find that latent cytomegalo-
38 virus infection drives DNA methylation variation and provide further support that the increased dis-
39 persion of DNA methylation with aging is due to epigenetic drift. Finally, our results indicate that
40 cellular composition and DNA sequence variation are the strongest predictors of DNA methylation,
41 highlighting critical factors for medical epigenomics studies.

42

43 **Introduction**

44 Epigenetic research has improved our understanding of the existing links between environmental
45 risk factors, aging, genetic variation and human disease^{1,2}. Epigenome-wide association studies
46 (EWAS) have shown that DNA methylation (i.e., 5-methylcytosine, 5mC), the most studied
47 epigenetic mark in humans, is associated with a wide range of environmental exposures along the
48 life course, such as chemicals³ or past socioeconomic status⁴⁻⁷. Changes in DNA methylation have
49 also been associated with non-communicable diseases, such as Parkinson's and Alzheimer's
50 diseases, multiple sclerosis, systemic lupus erythematosus, type 2 diabetes and cardiovascular
51 disease⁸⁻¹¹. These studies collectively suggest that DNA methylation marks could be of tremendous
52 value as gauges of the exposome and as clinical biomarkers^{12,13}.

53 However, interpretation of EWAS remains limited. First, because the epigenome of a cell
54 reflects its identity^{14,15}, a risk factor or a disease that alters cellular composition also alters 5mC
55 levels measured in the tissue¹⁶. It is thus necessary to determine if an exposure affects cellular
56 composition or DNA methylation states of cell-types, in order to better understand the link between
57 such an exposure, DNA methylation and disease¹⁷. Previous studies have accounted for cellular
58 heterogeneity in blood by using cell sorting experiments, or cellular proportions estimated from
59 5mC profiles through in-silico cell mixture deconvolution techniques^{18,19}, but these approaches
60 focus on a subset of frequent cell-types that capture only a part of blood cellular composition.
61 Second, the strong links between DNA methylation and DNA sequence variation, attested by the
62 numerous DNA methylation quantitative trait loci (meQTLs) detected so far²⁰⁻²³, suggest that
63 environmental effects on the epigenome may operate through gene-by-environment interactions, but
64 evidence for such interactions remains circumstantial. Finally, environmental risk factors with a yet-
65 unknown effect on DNA methylation, such as common infections, could confound associations
66 between other risk factors, DNA methylation and human phenotypes. Thus, a detailed study of the
67 factors that impact DNA methylation at the population level, and the extent to which their effects
68 are mediated by changes in cellular composition, is required to understand the role of epigenetic
69 variation in health and disease.

70 To address this gap, we generated whole blood-derived DNA methylation profiles at >850,000
71 CpG sites for 884 healthy adults of the Milieu Intérieur cohort. We leveraged the deep
72 characterization of the cohort, including high-resolution immunophenotyping by flow
73 cytometry^{24,25}, to determine whether and how cellular composition, intrinsic factors (i.e., age and
74 sex), genetic variation and 139 health- and immunity-related variables and environmental exposures
75 affect the blood DNA methylome. We first assessed differences in the DNA methylation profiles of
76 16 different immune cell-types. We then performed EWAS, adjusted or not for the measured

77 proportions of the 16 immune cell subsets, and mediation analyses to robustly delineate effects on
78 DNA methylation that are direct, i.e., acting through changes within cells, from those that are
79 mediated, i.e., acting through subtle changes in cellular composition²⁶. We show that adjusting
80 EWAS for 16 measured cell proportions better accounts for cellular heterogeneity than current cell
81 mixture deconvolution methods. We identify latent cytomegalovirus (CMV) infection as a key
82 factor affecting population variation in 5mC levels, through the dysregulation of human
83 transcription factors and profound changes in the proportion of differentiated T cells. We show that
84 the increased dispersion of DNA methylation with aging is independent of cellular composition,
85 supporting instead a decrease in the fidelity of the epigenetic maintenance machinery. Furthermore,
86 we show that a large part of the effects on DNA methylation of aging, smoking, CMV serostatus
87 and chronic low-grade inflammation is due to subtle changes in blood cell composition, and
88 characterize the DNA methylation signature of cell-types affected by these factors. Finally, we find
89 that the largest effects on DNA methylation are due to DNA sequence variation, whereas the most
90 widespread differences among individuals are the result of blood cellular heterogeneity. This work
91 generates new hypotheses about mechanisms underlying DNA methylation variation in the human
92 population and highlights critical factors to be considered in medical epigenomics studies.
93

94 **Results**

95 **Proportions of naïve and differentiated T cells markedly contribute to DNA methylation** 96 **variation**

97 To investigate the non-genetic and genetic factors that affect population variation in DNA
98 methylation, we quantified 5mC levels at >850,000 CpG sites, with the Illumina Infinium
99 MethylationEPIC array, in the 1,000 healthy donors of the Milieu Intérieur cohort (Fig. 1a). The
100 cohort includes individuals of Western European origin, equally stratified by sex (i.e., 500 women
101 and 500 men) and age (i.e., 200 individuals from each decade between 20 and 70 years of age), who
102 were surveyed for detailed demographic and health-related information²⁴, including factors that are
103 known to affect DNA methylation (i.e., age, sex, smoking, BMI and socioeconomic status), that
104 have been proposed to affect DNA methylation (e.g., dietary habits, upbringing) or that pertain to
105 the immune system (e.g., past and latent infections, past vaccinations, antibody levels;
106 Supplementary Data 1). All donors were genotyped at 945,213 single-nucleotide polymorphisms
107 (SNPs), yielding 5,699,237 accurate SNPs after imputation²⁵. After quality control filtering, high-
108 quality measurements of DNA methylation were obtained at 644,517 CpG sites for 884 unrelated
109 individuals²⁷ (Supplementary Fig. 1; Methods). We found that 5mC levels well reproduce expected
110 patterns across chromatin states¹⁵, supporting the good quality of the data (Supplementary Fig. 1
111 and Supplementary Notes).

112 Whereas most epigenome-wide studies adjust on estimated cellular composition to detect direct
113 effects on DNA methylation (i.e., acting through changes within cells), we sought to assess both
114 direct effects and effects that are mediated by changes in cellular composition, as the genomic
115 location and magnitude of mediated effects can inform us about how cell differentiation is regulated
116 in response to environmental exposures¹⁷. We thus measured, in all donors, the proportions of 16
117 immune cell subsets by standardized flow cytometry, including neutrophils, basophils, eosinophils,
118 monocytes, natural killer (NK) cells, dendritic cells, B cells, CD4⁻CD8⁻ T cells and naive, central
119 memory (CM), effector memory (EM) and terminally differentiated effector memory cells (EMRA)
120 CD4⁺ and CD8⁺ T cells²⁵.

121 We first determined which immune cell populations most affect DNA methylation variation, by
122 quantifying differences in 5mC levels between the 16 blood cell subsets with multivariable
123 regression models including log-ratios of cell subsets, defined according to the hierarchical and
124 compositional nature of the data²⁸ (Methods). We verified that our models are accurate, using
125 simulations and comparisons with independent DNA methylation data from sorted cellular
126 subsets²⁹. We found that our estimated effects of cell subset log-ratios on 5mC levels perform as
127 expected on simulated data (Supplementary Fig. 2 and Supplementary Notes) and are highly

128 correlated with DNA methylation differences observed between sorted immune cell fractions ($R >$
129 0.6; Supplementary Data 2). When applying these models on our data, we found that 5mC levels of
130 134,079 CpG sites (20.8% of CpG sites, Supplementary Data 2) are associated with the log-ratio of
131 myeloid vs. lymphoid lineages (Bonferroni corrected $P_{\text{adj}} < 0.05$). Furthermore, the log-ratio of
132 these subsets is the factor most associated with the first three Principal Components (PCs) of the
133 DNA methylation data (multiple linear mixed model of PC1: $P = 5.0 \times 10^{-18}$; PC2: $P = 1.6 \times 10^{-43}$;
134 PC3: $P = 6.7 \times 10^{-17}$), which respectively explain 11.4%, 7.5% and 5.5% of variation in DNA
135 methylation. Importantly, we also found that 20,758 and 44,919 CpG sites are associated with the
136 log-ratios of naïve and differentiated (CM, EM and EMRA) CD4⁺ and CD8⁺ T cell subsets,
137 respectively ($P_{\text{adj}} < 0.05$, Supplementary Data 2), supporting the view that 5mC levels differ
138 substantially among T cell sub-populations^{30,31}. Furthermore, the log-ratios of naïve and
139 differentiated CD4⁺ and CD8⁺ subsets are also associated with PC1 and PC3 ($P < 1.2 \times 10^{-4}$;
140 Fig. 1c,d). These results indicate that differences in the proportion of naïve and differentiated subsets
141 of CD4⁺ and CD8⁺ T cells contribute substantially to DNA methylation variation and may mediate
142 associations between DNA methylation and environmental exposures or diseases.

143

144 **Cell mixture deconvolution methods partially account for blood cell heterogeneity**

145 Direct effects of environmental exposures or diseases on DNA methylation are often estimated by
146 adjusting EWAS on major cell-type fractions, which are predicted in-silico from 5mC levels with
147 cell mixture deconvolution methods^{18,32}. However, standard methods only predict the overall
148 proportions of CD4⁺ and CD8⁺ T cells and may therefore overestimate the direct effects on DNA
149 methylation of factors that affect T cell composition, such as aging and viral infections^{25,33}. To test
150 this hypothesis, and to assess more generally how intrinsic and environmental factors affect the
151 DNA methylome, we conducted EWAS of 141 candidate factors, by using linear mixed models
152 adjusted on batch variables, genetic factors (i.e., associated meQTL variants), genetic ancestry,
153 smoking status, sex and a non-linear age term (Methods). Models were adjusted, or not, for the 16
154 measured cell proportions, to estimate total (i.e., direct and mediated) or direct effects, respectively.
155 Mediated effects were estimated by mediation analysis³⁴ (Methods). We considered that each
156 EWAS constitutes a separate family of association tests and used the Bonferroni correction for
157 multiple testing adjustment ($P_{\text{adj}} < 0.05$).

158 Out of the 141 candidate factors, those that have significant total effects on DNA methylation
159 include age ($n = 97,219$ CpG sites; 15.1% of CpG sites), cytomegalovirus (CMV) serostatus ($n =$
160 79,654; 12.4%), sex ($n = 23,002$; 3.6%), heart rate ($n = 2,924$; 0.5%), smoking ($n = 839$; 0.1%),
161 body temperature ($n = 175$), C-reactive protein (CRP) levels ($n = 53$), the hour of blood draw ($n =$

162 36) and traits related to lipid metabolism ($n = 3$; Fig. 1b and Supplementary Data 1). Accordingly,
163 the first PCs of DNA methylation are most strongly associated with CMV (PC1: $P = 8.3 \times 10^{-13}$;
164 PC2: $P = 7.8 \times 10^{-10}$), age (PC3: $P = 5.7 \times 10^{-29}$) and sex (PC4: $P = 2.2 \times 10^{-5}$), when not considering
165 immune cell fractions (Fig. 1c,d and Supplementary Fig. 1i,j). When adjusting on blood cell
166 composition, factors that have significant direct effects on DNA methylation include age ($n =$
167 35,701; 5.5%), sex ($n = 17,067$; 2.6%), smoking ($n = 428$; 0.07%), CMV serostatus ($n = 245$;
168 0.04%), CRP levels ($n = 39$) and lipid metabolism-related traits ($n = 3$; Fig. 1b, Supplementary Fig.
169 3 and Supplementary Notes). These results suggest that, whereas most CMV effects are mediated by
170 cellular composition, the effects of sex on DNA methylation are mainly direct, and a substantial
171 direct effect of age is also retained, even after adjusting for naïve and memory CD4⁺ and CD8⁺ T
172 cell subsets. Accordingly, first PCs of DNA methylation remain associated with sex (PC4: $P =$
173 1.3×10^{-3}) and age (PC3: $P = 1.1 \times 10^{-9}$; Fig. 1c), when considering immune cell fractions, but not
174 with CMV serostatus (PC1: $P > 0.05$; Fig. 1d). No significant direct effects of heart rate, body
175 temperature and hour of sampling were detected, indicating that the effects of these factors on DNA
176 methylation are due exclusively to changes in immune cell composition^{35,36}.

177 We then evaluated the performance of three reference-based in-silico cell mixture deconvolution
178 methods: Houseman et al.'s method, IDOL and EPIC IDOL-Ext^{18,29,32}. We observed that cell
179 proportions estimated by the three methods are substantially correlated with measured cell
180 proportions (Supplementary Fig. 4). We then compared EWAS results adjusted either on our flow
181 cytometric data or on cell proportions estimated by the three deconvolution methods. We found that
182 EWAS adjusted by the IDOL method detects more CpG sites associated with most candidate
183 factors, relative to EWAS adjusted on the measured proportions of 16 cell-types, particularly for age
184 ($n = 131,142$ vs. 35,701) and latent CMV infection ($n = 31,159$ vs. 245) (Fig. 1b,e,f). Similar results
185 were found with Houseman's method (Fig. 1b). Accordingly, the first PC of DNA methylation
186 remains strongly associated with CMV serostatus and age when adjusting on IDOL cellular
187 fractions ($P = 7.5 \times 10^{-6}$ and $P = 3.2 \times 10^{-17}$, respectively), whereas it is not when considering 16
188 measured cell proportions ($P > 0.01$). Conversely, EWAS adjusted by the EPIC IDOL-Ext method,
189 which estimates subsets of naïve and memory CD4⁺ and CD8⁺ T cell populations²⁹, provide results
190 that are similar to those of EWAS adjusted for high-resolution flow cytometric data (Fig. 1b). These
191 results suggest that first-generation deconvolution methods do not fully distinguish direct effects on
192 DNA methylation from those that are mediated by fine-grained changes in blood cell composition.

193 To further test this scenario, we conducted EWAS adjusted on flow cytometric data for only six
194 major cell-types and found results comparable to those for Houseman et al.'s and the IDOL
195 methods (Fig. 1b). Furthermore, CMV effect sizes adjusted on IDOL cellular fractions or the 6

196 major cell proportions were twice more correlated with estimated measures of DNA methylation
197 differences between naïve and differentiated CD4⁺ T cells, relative to CMV effect sizes adjusted on
198 16 measured cell proportions ($R = 0.66$, relative to $R = 0.31$, respectively; Fig. 1g,h). Together, these
199 results indicate that adjustment for the proportions of only the six major cell-types is not able to
200 fully account for blood cell heterogeneity, particularly when estimating the effects of age and CMV
201 infection on DNA methylation, two factors that are known to skew CD4⁺ and CD8⁺ T cell
202 compartments toward differentiated phenotypes²⁵.

203

204 **Cytomegalovirus infection alters the blood DNA methylome through regulation of host** 205 **transcription factors**

206 We identified CMV serostatus as one of the exposures that is associated with the largest number of
207 CpG sites (Fig. 1b). CMV is the causative agent of a latent, mainly asymptomatic, infection that
208 ranges in seroprevalence from 30% to 100% across populations³⁷. CMV is known to drastically
209 alter the composition of the CD4⁺ and CD8⁺ T cell compartments in blood^{25,33}. Accordingly, we
210 found that 85,922 CpG sites show a significant cell-composition-mediated effect of CMV serostatus
211 on DNA methylation ($P_{\text{adj}} < 0.05$; Supplementary Data 1), indicating that the effects of the latent
212 infection are mainly mediated by cellular composition. Furthermore, we observed a strong
213 correlation between mediated and total effect sizes of CMV serostatus ($R = 0.93$; Fig. 2a) and
214 99.5% of CpG sites with a significant direct effect also show a significant mediated effect ($n = 244 /$
215 245). We found that mediated effect sizes of CMV are strongly correlated with estimated measures
216 of DNA methylation differences between naïve and memory CD4⁺ and CD8⁺ T cells ($R = 0.68$ and
217 $R = 0.53$, respectively; Fig. 2b), suggesting that cell-composition-mediated effects of CMV are
218 predominantly attributable to changes in these T cell subsets.

219 One of the strongest cell-composition-mediated effects of CMV infection was observed in an
220 intron of *DNMT3A* (β value scale 95% confidence interval [CI]: [1.8%, 2.4%], $P_{\text{adj}} = 1.1 \times 10^{-23}$),
221 encoding a key DNA methyltransferase playing a role in the replication of some herpesviruses³⁸.
222 CMV⁺ donors show a substantial increase in the proportion of CD4⁺ and CD8⁺ T_{EMRA} cells ($P =$
223 6.8×10^{-35} and $P = 1.9 \times 10^{-50}$, respectively), which in turn are associated with higher 5mC levels at
224 *DNMT3A* ($P = 3.3 \times 10^{-25}$ and $P = 1 \times 10^{-53}$, respectively), supporting mediation by differentiated
225 memory T cell subsets (Fig. 2c). To test if the effects of CMV infection on 5mC levels are cell-type-
226 dependent, we derived and verified an interaction model similar to CellDMC³⁹ (Methods). We
227 restricted this analysis to interactions with the proportion of cells from the myeloid lineage, as
228 previously reported⁴⁰, and found only one CpG site where CMV effects depend on the proportion of
229 myeloid cells ($P_{\text{adj}} < 0.05$; Supplementary Data 3). These results indicate that CMV infection affects

230 a large fraction of the blood DNA methylome primarily through changes in blood cell proportions,
231 rather than through cell-type-dependent changes.

232 However, when adjusting for blood cell composition, including CD4⁺ and CD8⁺ T cell sub-
233 types, a significant direct effect of CMV serostatus was detected for 245 CpG sites. Increased 5mC
234 levels in CMV⁺ donors localize predominantly in enhancers and regions flanking transcription start
235 sites (odds ratio [OR] > 3.0, $P_{\text{adj}} < 5.3 \times 10^{-8}$; Supplementary Fig. 5), suggesting dysregulation of
236 host gene expression as a result of latent infection. The second strongest direct effect of CMV
237 infection was observed nearby the TSS of *LTBP3* (β value scale 95% CI: [1.9%, 3.1%], $P_{\text{adj}} =$
238 7.1×10^{-17} ; Fig. 2d and Supplementary Fig. 6). *LTBP3* is a regulator of transforming growth factor β
239 (TGF- β)⁴¹, which is induced in CMV latently infected cells⁴². Strikingly, CpG sites showing
240 increased 5mC levels in CMV⁺ donors are strongly enriched in binding sites for the BRD4
241 transcription factor (TF) ($n = 187 / 189$, OR = 48.0, 95% CI: [13.1, 399.0], $P_{\text{adj}} < 1.1 \times 10^{-27}$; Fig. 2e
242 and Supplementary Data 4), a bromodomain protein that plays a critical role in the regulation of
243 latent and lytic phases of CMV infection⁴³. Conversely, CpG sites showing a decrease in DNA
244 methylation in CMV⁺ donors are strongly enriched in binding sites for BATF3 (OR = 24.8, 95% CI:
245 [13.8, 42.2], $P_{\text{adj}} < 1.3 \times 10^{-14}$; Fig. 2f), which is paramount in the priming of CMV-specific CD8⁺ T
246 cells by cross-presenting dendritic cells⁴⁴. Collectively, these analyses imply that CMV infection
247 directly affects the human blood DNA methylome through the dysregulation of host TFs implicated
248 in viral latency and host immune response.

249 Finally, to motivate future research on the epigenetic effects of CMV infection, we used elastic
250 net regression and stability selection to predict CMV serostatus from DNA methylation (Methods).
251 Based on 547 CpG sites, the model predicts CMV serostatus with an out-of-sample accuracy of
252 87%, using 10-fold cross-validation. We anticipate that this model will be useful to determine if
253 latent CMV infection can confound epigenetic risk for disease^{45,46}.

254

255 **Ageing elicits DNA hypermethylation related to Polycomb repressive complexes and increased** 256 **epigenetic dispersion**

257 Although the effects of ageing on DNA methylation are well established⁴⁷⁻⁵¹; it remains unclear the
258 extent to which they are due to changes in unmeasured proportions of differentiated T cells (Fig. 1b)
259 or CMV infection, which are both strongly associated with age^{25,52}. Indeed, age has a significant
260 total effect on 5mC levels at 97,219 and 113,742 CpG sites, when adjusting or not on CMV
261 serostatus, and CMV infection mediates a substantial fraction of total age effects ($n = 10,074$ CpG
262 sites). We thus investigated how the blood DNA methylome is shaped by the intertwined processes

263 of cellular aging (i.e., direct effects) and age-related changes in blood cellular composition (i.e.,
264 mediated effects), while accounting for CMV serostatus.

265 We found that, out of the 35,701 CpG sites associated directly with age, more than 97% were
266 associated with age in a previous EWAS⁵³, indicating a strong overlap (OR 95% CI: [35.6, 40.8]).
267 In line with previous findings⁵⁴, direct effects of age are typically larger than mediated effects (Fig.
268 3a). Furthermore, the strongest direct age effects, such as those observed at *ELOVL2* and *FHL2*
269 (Supplementary Fig. 6), are not mediated by cellular composition ($P_{\text{adj}} = 1.0$), suggesting that age-
270 related changes at these CpG sites are typically shared across cell-types. We observed that 61% of
271 the CpG sites directly associated with age show a decrease in 5mC levels. Age-associated
272 demethylation predominates outside of CpG islands (CGIs) and in regions flanking transcription
273 start sites and in enhancers (Fig. 3b and Supplementary Fig. 7a,b). Conversely, DNA
274 hypermethylation was observed in 95% of age-associated CpGs within CGIs. Consistently, CpG
275 sites exhibiting increasing 5mC levels with age are mainly found in Polycomb-repressed regions,
276 bivalent TSSs and bivalent enhancers (Fig. 3b,c), which are CGI-rich regions (Supplementary Fig.
277 1M,N). Furthermore, these CpG sites are most enriched in binding sites for RING1B, JARID2,
278 RYBP, PCGF1, PCGF2 and SUZ12 TFs (OR > 10.0; Fig. 3d and Supplementary Data 4), which are
279 all part of the Polycomb repressive complexes 1 and 2. PRC1 and PRC2 mediate cellular
280 senescence and modulate longevity in invertebrates^{55,56}. Importantly, when restricting the analysis
281 to CpG sites outside of CpG islands, we found similar enrichments in Polycomb-repressed regions
282 (OR 95% CI [17.7, 20.0]) and PRC TF binding sites (RING1B OR 95% CI: [19.9, 22.4]; PCGF2
283 OR 95% CI [17.8, 20.7]). Finally, genes with age-increasing 5mC levels are strongly enriched in
284 developmental genes ($P_{\text{adj}} = 1.7 \times 10^{-48}$; Supplementary Data 5), which are regulated by PRCs⁵⁷.
285 Overall, these results confirm previously described effects of age on the blood DNA methylome,
286 while accounting more comprehensively for blood cell composition and CMV infection, and
287 support a key regulatory role of Polycomb proteins in age-related hypermethylation⁵⁸.

288 We then assessed whether age-related changes in blood cell composition or CMV seropositivity
289 could contribute to age-related changes in the variance of 5mC levels, a phenomenon known as
290 “epigenetic drift” (i.e., the divergence of the DNA methylome as a function of age owing to
291 stochastic changes)^{51,59-61}. We observed that the proportion of several cell-types in blood are
292 increasingly dispersed with aging, such as CD4⁺ T_{EMRA} cells (Fig. 3e). Therefore, we fitted models
293 parameterizing the residual variance with a linear age term, and adjusting for 16 immune cell
294 proportions, age, CMV serostatus, smoking status and sex in the mean function (Methods). We
295 observed a significant dispersion of DNA methylation with age for 3.1% of all CpG sites ($n =$
296 20,140, $P_{\text{adj}} < 0.05$). We compared these CpG sites with those previously reported to be increasingly

297 variable with age in whole blood and monocytes⁶⁰ and replicated 2,604 out of 5,075 CpG sites,
298 supporting a strong overlap between the two different approaches (OR 95% CI: [36.2, 40.8]). An
299 example of a CpG site with a large, age-increasing dispersion is found in the TSS of *MAFA* ($P_{\text{adj}} =$
300 4.4×10^{-43} ; Fig. 3f), encoding a transcription factor that regulates insulin. Strikingly, 99.4% of CpGs
301 with age-related dispersion show an increase in the variance of 5mC levels with age (Fig. 3g),
302 supporting a decrease in the fidelity of epigenetic maintenance associated with aging. In addition,
303 we found that, out of 20,140 CpG sites with age-related dispersion, 87.3% show no significant
304 changes in mean 5mC levels with age, and we detected no correlation between estimates of
305 dispersion and direct age effect sizes (Fig. 3h), implying that these results are not driven by
306 relationships between the average and variance of 5mC levels. Furthermore, when also adjusting the
307 variance function for cellular composition, we found evidence of dispersion in 8,576 CpG sites (P_{adj}
308 < 0.05), with similar effect sizes as in the previous model ($R = 0.93$; Methods). Collectively, these
309 findings indicate that aging elicits numerous DNA methylation changes in a cell-composition-
310 independent manner, including global epigenome-wide demethylation, hypermethylation of PRC-
311 associated regions and increased variance, highlighting the occurrence of different mechanisms
312 involved in epigenetic aging.

313

314 **Immunosenescence-related changes in cellular composition mediate DNA methylation** 315 **variation with age**

316 We detected a significant cell-composition-mediated effect of age at $\sim 1.1\%$ of CpG sites ($n = 7,090$;
317 Fig. 3a and Supplementary Data 1), indicating that a substantial fraction of age-related changes in
318 DNA methylation are due to age-related changes in immune cell proportions. Mediated effects are
319 most often associated with demethylation (76% of age-associated CpG sites), regardless of the
320 chromatin state or CGI density of the loci considered (Fig. 3j and Supplementary Fig. 7c,d).
321 Enhancers and regions flanking transcription start sites are enriched in CpG sites with a significant
322 cell-composition-mediated effect of age (Fig. 3i), possibly because these regions tend to be
323 regulated in a cell-type-dependent manner¹⁵. In contrast with direct age effects, CpG sites with a
324 cell-composition-mediated increase in DNA methylation are enriched in TF binding sites for
325 *RUNX1-3* (OR = 8.5, 95% CI: [4.5, 14.7], $P_{\text{adj}} < 1.2 \times 10^{-8}$), which are key regulators of
326 hematopoiesis (Fig. 2k and Supplementary Data 4). Genes with CpG sites showing a mediated
327 increase or decrease in DNA methylation with age are enriched in genes involved in lymphoid (P_{adj}
328 $= 2.0 \times 10^{-7}$) and myeloid ($P_{\text{adj}} = 6.1 \times 10^{-13}$) cell activation, respectively (Supplementary Data 5).
329 This indicates that mediated effects of age on DNA methylation are related to progressive, lifelong
330 differences in the composition of the lymphoid and myeloid cell lineages.

331 We then determined if age effects on 5mC levels depend on the proportion of cells from the
332 myeloid lineage, by using an interaction model (Methods). In line with a previous study⁵⁴, we found
333 that cell-type-dependent effects of age (Supplementary Data 3) are limited; only 10 CpG sites show
334 DNA methylation changes with age that depend on the proportion of myeloid cells ($P_{\text{adj}} < 0.05$;
335 Supplementary Data 3). Importantly, age also has a strong mediated effect on all these CpG sites
336 ($P_{\text{adj}} < 1.0 \times 10^{-10}$), implying that these loci are associated with age because of changes in blood cell
337 composition, although their relation to age is cell-type-dependent. Collectively, our findings provide
338 statistical evidence that DNA methylation variation with age results from different, non-mutually
339 exclusive mechanisms: the progressive decline of the epigenetic maintenance system that is
340 common to all cell-types, the increased heterogeneity of immune cell subsets that characterizes
341 immunosenescence⁶² and, to a lesser extent, accelerated changes within specific blood cell
342 compartments.

343

344 **Sex differences in DNA methylation are predominantly cell- and age-independent**

345 Given that substantial differences in immune cell composition have been observed between women
346 and men²⁵, we next assessed how cellular heterogeneity contributes to sex differences in DNA
347 methylation⁶³⁻⁶⁵. We found 3.6% of CpG sites ($n = 23,002$) with a significant total effect of sex,
348 2.6% ($n = 17,067$) with a significant direct effect, and only 0.2% ($n = 1,385$) with a significant cell-
349 composition-mediated effect ($P_{\text{adj}} < 0.05$; Supplementary Fig. 8a and Supplementary Data 1). Out
350 of CpG sites directly associated with sex, 96.2% were already associated with sex in a previous
351 EWAS⁵³, indicating again a strong overlap (OR 95% CI: [39.6, 46.5]). The largest direct effects of
352 sex were observed at *DYRK2*, *DNMI*, *RFTNI*, *HYDIN*, and *NABI* genes ($P_{\text{adj}} < 1.0 \times 10^{-263}$;
353 Supplementary Fig. 6). For example, the *DYRK2* promoter is 11.7% and 45.6% methylated in men
354 and women, respectively, at a CpG site that we found to be bound by the X-linked PHF8 histone
355 demethylase (Supplementary Fig. 8b,c). *DYRK2* phosphorylates amino acids and plays a key role
356 in breast and ovarian cancer development⁶⁶.

357 DNA methylation levels are higher in women at 79.7% of sex-associated autosomal CpG sites
358 (Supplementary Fig. 8d,e), a pattern also observed in newborns⁶⁴. This proportion is similar across
359 different genomic regions, based on either chromatin states or CpG density (Supplementary Fig.
360 8e,g). When quantifying how sex differences in DNA methylation vary during adulthood, by adding
361 a sex-by-age interaction term to our models (Methods), we found only 7 CpG sites with a
362 significant, sex-dependent effect of age ($P_{\text{adj}} < 0.05$; Supplementary Data 3). Confirming previous
363 findings^{53,67}, the strongest sex-by-age interaction effects were found at *FIGN* ($P_{\text{adj}} < 7.1 \times 10^{-15}$),
364 associated with risk-taking behaviors⁶⁸ and educational attainment⁶⁹, and *PRR4* ($P_{\text{adj}} < 5.6 \times 10^{-3}$),

365 associated with the dry eye syndrome, a hormone-dependent, late-onset disorder⁷⁰. Overall, our
366 findings indicate that the blood DNA methylome is widely affected by sex, but its effects are
367 typically not mediated by cellular composition and do not change during adulthood.

368

369 **Gene × cell-type and gene × environment interactions affect DNA methylation variation**

370 Gene × environment interactions are thought to underlie adaptable human responses to
371 environmental exposures through epigenetic changes⁷¹. To test if gene × environment interactions
372 affect DNA methylation, we first estimated, for each CpG site, the effects on 5mC levels of local
373 and remote DNA sequence variation, defined as genetic variants within a 100-Kb window and
374 outside a 1-Mb window centered on the CpG site, respectively (Methods). We considered local and
375 remote meQTLs to be independent families of tests and used the Bonferroni correction to adjust for
376 multiple testing. We found a significant local meQTL for 107,048 CpG sites and a significant
377 remote meQTL for 1,228 CpG sites ($P_{\text{adj}} < 0.05$; Supplementary Fig. 9 and Supplementary Data 6).
378 In agreement with previous studies^{21,23}, CpG sites with a local meQTL are enriched in enhancers
379 (OR 95% CI: [2.09, 2.21]) and depleted in TSS and actively transcribed genes (OR 95% CIs: [0.52,
380 0.56] and [0.57, 0.60]; Fig. 4a). Conversely, CpG sites under remote genetic control are enriched in
381 TSS regions (OR 95% CI: [2.10, 3.11]) and regions associated with *ZNF* genes (OR 95% CI: [1.26,
382 6.17]; Fig. 4b). Furthermore, we found that remote meQTL variants are also strongly concentrated
383 in *ZNF* genes (OR 95% CI: [14.6, 29.8]; Fig. 4c), suggesting that zinc-finger proteins (ZFPs) play a
384 role in the long-range control of DNA methylation, in line with their role in the regulation of
385 heterochromatin⁷²⁻⁷⁴.

386 We next explored whether effects of genetic variants on 5mC levels depend on the circulating
387 proportion of myeloid cells. We found evidence for cell-type-dependent meQTLs at only 249 CpG
388 sites ($P_{\text{adj}} < 0.05$; Fig. 4d and Supplementary Data 3), supporting the notion that genetic effects on
389 5mC levels are generally shared across blood cell subsets⁷⁵. The strongest signal was found between
390 5mC levels upstream of *CLEC4C* and the nearby rs11055602 variant, which has been previously
391 shown to strongly affect *CLEC4C* protein levels⁷⁶. This C-type lectin, known as CD303, is used as a
392 differentiation marker for dendritic cells, suggesting the epigenetic regulation of the locus is cell-
393 type-dependent. Accordingly, rs11055602 genotype effects on DNA methylation depend on the
394 circulating proportions of myeloid cells (β scale interaction effect, 95% CI: [0.16, 0.22], $P_{\text{adj}} =$
395 7.4×10^{-20} ; Fig. 4e), and dendritic cells (95% CI: CI: [-8.3, -5.0], $P_{\text{adj}} = 3.5 \times 10^{-15}$).

396 We then evaluated whether the main non-heritable determinants of DNA methylation variation
397 in our cohort, i.e., age, sex, CMV serostatus, smoking status and chronic low-grade inflammation
398 (CRP levels; Fig. 1b, Supplementary Fig. 3 and Supplementary Notes), can affect 5mC levels in a

399 genotype-dependent manner. We thus tested for genotype \times age, genotype \times sex, genotype \times
400 smoking jointly (Methods). Genotype \times CRP levels interactions were tested in a separate model that
401 include the other interaction effects. We found statistical evidence for genotype-dependent effects of
402 age and sex at 68 and 20 CpG sites, respectively ($P_{\text{adj}} < 0.05$, $\text{MAF} > 0.10$; Fig. 4d and
403 Supplementary Data 3), the interacting meQTL variant being local in all cases. We detected a strong
404 genotype \times age interaction for two CpG sites located in the *BACE2* gene, the 5mC levels of which
405 decrease with age only in donors carrying the nearby rs2837990 G>A allele (β scale 95% CI: [0.11,
406 0.13], $P_{\text{adj}} = 7.28 \times 10^{-10}$; Fig. 4f). *BACE2* encodes beta-secretase 2, one of two proteases involved in
407 the generation of amyloid beta peptide, a critical component in the etiology of Alzheimer's
408 disease⁷⁷. Another strong genotype \times age interaction effect was found for a CpG site upstream of
409 *FCERIA*, encoding the high-affinity IgE receptor. *FCERIA* 5mC levels decrease with age in
410 rs2251746 T>C carriers only (95% CI: [0.05,0.07], $P_{\text{adj}} = 8.6 \times 10^{-9}$), a variant known to control
411 serum IgE levels⁷⁸. Collectively, our analyses identify few, albeit strong, environment- and cell-
412 type-dependent meQTLs, supporting the relatively limited impact of gene \times cell-type and gene \times
413 environment interactions on the blood DNA methylome.

414

415 **Cellular composition and genetics drive DNA methylation variation in human blood**

416 Having established how cellular composition, intrinsic factors, genetic variation, and a broad
417 selection of non-heritable factors shape the blood DNA methylome, we next sought to compare the
418 relative impact of these factors on DNA methylation. We classified the factors into four groups: (i)
419 the cellular composition group, which consists of the 16 measured cell proportions; (ii) the intrinsic
420 group, which consists of age and sex; (iii) the genetic group, which consists of the most associated
421 local-meQTL variant around each CpG site; and (iv) the exposure group, which consists of smoking
422 status, CMV serostatus and CRP levels. Since these groups vary in their degrees of freedom, we
423 measured the relative predictive strength for each CpG site by the out-of-sample prediction
424 accuracy, estimated by cross-validation (Methods). To ensure unbiased estimates, we mapped local
425 meQTLs anew within each training set.

426 The full model that includes all groups explains $< 5\%$ of out-of-sample variance for 52.3% of
427 CpG sites (Fig. 5a), which are typically characterized by low total 5mC variance (Supplementary
428 Fig. 10). This suggests that these sites are constrained in the healthy population and that small
429 fluctuations in 5mC levels determine their variation, possibly due to measurement errors or
430 biological noise. Nevertheless, the model explains $> 25\%$ of DNA methylation variance for 20.8%
431 of CpG sites ($n = 134,305$). The strongest predictor for these CpG sites is cellular composition,
432 genetics, intrinsic factors and exposures in 74.7%, 21.5%, 3.8% and 0.01% of cases, respectively.

433 Cellular composition explains > 25% of out-of-sample variance for 1.0% of CpG sites ($n = 90,033$;
434 Fig. 5a,c and Supplementary Data 7), with the highest variance explained by cellular composition
435 for one CpG site being 71.8%. For the 2,580 CpG sites where the model explains > 75% of
436 variance, local DNA sequence variation is the strongest predictor in 99.2% of cases (Fig. 5c and
437 Supplementary Data 7). Local genetic variation explains > 25% of DNA methylation variance at
438 23,677 CpG sites, and almost as many when adjusting for cellular composition ($n = 22,865$) (Fig.
439 5a,b), indicating that genetic effects on 5mC levels are mainly cell-composition-independent.
440 Intrinsic factors explain > 25% of out-of-sample variance at 3,669 CpG sites, and > 75% at 16 sites
441 (Fig. 5c). When conditioning on cell composition, these numbers dropped to 334 and 6 CpG sites,
442 respectively, suggesting that the predictive ability of age and sex is partly mediated by immune cell
443 composition (Fig. 5b). Interestingly, environmental exposures are the weakest predictor of 5mC
444 levels, explaining > 25% of the variance at only 29 CpG sites and with a maximum variance
445 explained for a CpG site of 50.1%.

446 Finally, we estimated the proportion of variance explained by genotype \times age, genotype \times sex
447 and genotype \times exposure interactions, by considering the difference of the out-of-sample variance
448 explained by models including interaction terms and models with only main effects (Methods). We
449 found a significant increase in predictive ability when including interaction terms for 431 CpG sites
450 (ANOVA $P_{\text{adj}} < 0.05$). However, the effects were typically modest: only 13 CpG sites showed an in-
451 crease in the proportion of variance explained larger than 5% (Fig. 5b). Collectively, these results
452 show that cellular composition and local genetic variation are the main drivers of DNA methylation
453 variation in the blood of adults, reinforcing the critical need to study epigenetic risk factors and bi-
454 omarkers of disease in the context of these factors.

455 **Discussion**

456 Here, we present a rich data resource that delineates the contribution of blood cellular composition,
457 age, sex, genetics, environmental exposures and their interactions to variation in the DNA
458 methylome. All the results can be explored via a web-based browser ([MIMETH browser](#)), to
459 facilitate the exploration of the estimated effects of these factors on DNA methylation variation. We
460 found that CMV infection elicits substantial changes in the blood DNA methylome, in contrast with
461 other herpesviruses such as EBV, HSV-1, HSV-2 and VZV. Latent CMV infection is known to
462 profoundly alter the number, activation status and transcriptional profiles of immune cell
463 populations, yet its epigenetic consequences have attracted little attention. We observed that most
464 CMV effects on DNA methylation are mediated by the profound changes in blood cell
465 composition²⁵, including the CMV-driven inflation of memory CD4⁺ and CD8⁺ T cells³³. However,
466 we also detected cell-composition-independent effects of CMV infection, suggesting that the
467 herpesvirus can directly regulate the host epigenome. Notably, differentially methylated CpG sites
468 in CMV⁺ donors are strongly enriched in binding sites for BRD4, a key host regulator of CMV
469 latency⁴³, suggesting that the recruitment of BRD4 by CMV during latent infection affects BRD4-
470 regulated host genes. Furthermore, CMV⁺ donors are characterized by a strong increase in 5mC
471 levels at *LTBP3*, the product of which is involved in TGF- β secretion. TGF- β is a well-known
472 immunosuppressive cytokine induced by CMV infection⁴², which represents a possible strategy of
473 the virus to escape host immunity. These results suggest that the capacity of CMV to manipulate the
474 host epigenetic machinery results in epigenetic changes of latently infected cells.

475 Our study provides further support to the notion that three different biological mechanisms un-
476 derlie age-related changes in DNA methylation. The first elicits an increased dispersion of 5mC lev-
477 els with age that is related to epigenetic drift^{51,59-61}. We found that dispersion of DNA methylation
478 with age is not due to cellular heterogeneity, supporting instead the progressive decline in fidelity of
479 the DNA methylation maintenance machinery across cell populations. The second mechanism re-
480 sults in cell-composition-independent, global DNA demethylation and CGI-associated hypermethyl-
481 ation. Age-associated DNA demethylation could be related to the downregulation of DNMT3A/B
482 *de novo* methyltransferases, whereas CGI-associated hypermethylation may result from the down-
483 regulation of the Polycomb repressive complexes 1 and 2 and/or TET proteins, coupled with a loss
484 of H3K27me3 marks⁷⁹⁻⁸¹. Alternatively, these changes may be related to the mitotic clock, which
485 assumes a progressive accumulation of DNA methylation changes with mitotic divisions, including
486 loss of methylation at partially methylated domains (PMD) and gain of methylation at PRC2-
487 marked CpG-rich regions⁸²⁻⁸⁴. Both scenarios are supported by the enrichment of Polycomb-re-
488 pressed regions in age-associated CpG sites, and of binding sites of PRC-related TFs in CpG sites

489 methylated with age. The third mechanism elicits cell-composition-mediated demethylation at all
490 compartments of the epigenome, particularly at enhancers of myeloid activation genes. This process
491 likely reflects an increased degree of differentiation in the lymphoid compartment with age. Single-
492 cell methylomes of differentiating and dividing white blood cells will help determine the role of mi-
493 totic and post-mitotic 5mC changes during epigenetic aging.

494 Another interesting finding of our study is that environmental exposures explain a small fraction
495 of the variance of DNA methylation in healthy adults, at odds with the common view that the epige-
496 nome is strongly affected by the environment⁸⁵. Twin studies have estimated the heritability of DNA
497 methylation to range from ~20-40% (ref.⁸⁶⁻⁸⁸), suggesting that environmental effects, along with
498 gene × environment interactions, account for the remaining 60-80% (ref.⁸⁹). However, other factors,
499 including cellular composition and measurement error, may account for most of the unexplained
500 variance. Consistently, we estimated that cellular composition explains >25% of the variance for
501 ~13% of the DNA methylome, and it has been estimated that measurement error may explain >50%
502 (ref.⁹⁰). Nevertheless, a limitation of our study is that perinatal and early life exposures, which are
503 thought to contribute extensively to epigenetic variation in adulthood⁸⁵, have not been extensively
504 assessed in the Milieu Intérieur cohort. In addition, it has been hypothesized that gene × environ-
505 ment interactions are central to understand the role of epigenetics in development⁹¹, but statistical
506 evidence for interaction effects requires larger cohorts⁹², suggesting that our results might represent
507 a small, perceptible fraction of a large number of weak effects^{93,94}. Large, longitudinal cohorts ad-
508 dressing the developmental origins of disease are needed to shed new light on the role of DNA
509 methylation in the interplay between genes and the environment.

510 Collectively, our findings have broad consequences for the study and interpretation of epige-
511 netic factors involved in disease risk. First, our analyses show that first-generation cell mixture de-
512 convolution methods^{18,32} do not fully distinguish direct from cell-composition-mediated effects of
513 CMV infection and age on DNA methylation, probably because these two factors alter the propor-
514 tions of blood cell subsets that are not estimated by these methods. This reinforces the view that
515 EWAS must be interpreted with great caution, particularly when the studied diseases or conditions
516 are known to affect unmeasured immune cell fractions. Encouragingly, our findings suggest that,
517 when blood cell composition is not measured directly, high resolution cell mixture deconvolution
518 methods^{29,95} provide a more complete correction for cellular heterogeneity and are therefore ex-
519 pected to improve the interpretation of future epigenomic studies. Second, because age, sex, CMV
520 infection, smoking and chronic low-grade inflammation can influence disease risk^{45,96-99}, our results
521 emphasize the critical need to consider such factors in EWAS, as these factors can confound associ-
522 ations. Lastly, our findings reveal the epigenetic impact of aging and persistent viral infection

523 through fine-grained changes in blood cell proportions, highlighting the need to assess the respec-
524 tive role of altered cellular composition and DNA methylation in the etiology of disease¹⁷. Large-
525 scale studies using single-cell approaches will help overcome these challenges, and are anticipated
526 to further decode the epigenetic mechanisms underlying healthy aging and the environmental causes
527 of human disease.
528

529 **Methods**

530 **The Milieu Intérieur cohort**

531 The Milieu Intérieur cohort was established with the goal to identify genetic variation and
532 environmental exposures that affect phenotypes related to the immune system in the adult, healthy
533 population. The 1,000 healthy donors of the Milieu Intérieur cohort were recruited by BioTrial
534 (Rennes, France), and included 500 women and 500 men. Donors included 100 women and 100
535 men from each decade of life, between 20 and 69 years of age. Donors were selected based on
536 various inclusion and exclusion criteria that are detailed elsewhere²⁴. Briefly, donors were required
537 to have no history or evidence of severe/chronic/recurrent pathological conditions, neurological or
538 psychiatric disorders, alcohol abuse, recent use of illicit drugs, recent vaccine administration, and
539 recent use of immune modulatory agents. To avoid the influence of hormonal fluctuations in
540 women, pregnant and peri-menopausal women were not included. To avoid genetic stratification in
541 the study population, the recruitment of donors was restricted to individuals whose parents and
542 grandparents were born in Metropolitan France.

543

544 **Ethical approvals**

545 The study is sponsored by the Institut Pasteur (Pasteur ID-RCB Number: 2012-A00238-35) and was
546 conducted as a single center study without any investigational product. The Milieu Intérieur clinical
547 study was approved by the Comité de Protection des Personnes — Ouest 6 (Committee for the
548 protection of persons) on June 13, 2012 and by the French Agence Nationale de Sécurité du
549 Médicament (ANSM) on June 22, 2012. The samples and data used in this study were formally
550 established as the Milieu Intérieur biocollection (study# NCT03905993), with approvals by the
551 Comité de Protection des Personnes – Sud Méditerranée and the Commission nationale de
552 l'informatique et des libertés (CNIL) on April 11, 2018.

553

554 **DNA sampling and extraction**

555 Whole blood was drawn from the 1,000 Milieu Intérieur healthy, fasting donors every working day
556 from 8AM to 11AM, from September 2012 to August 2013, in Rennes, France. Different
557 anticoagulants were used, depending on the downstream analyses. For DNA methylation profiling,
558 blood samples were collected on EDTA, whereas samples for flow cytometry and genome-wide
559 DNA genotyping were collected on Li-heparin. Tracking procedures were established in order to
560 ensure delivery to Institut Pasteur (Paris) within 6 hours of blood draw, at a temperature between
561 18°C and 25°C. Upon receipt, samples were kept at room temperature until DNA extraction. DNA
562 for DNA methylation profiling was extracted using the Nucleon BACC3 genomic DNA extraction

563 kit (catalog #: RPN8512; Cytiva, Massachusetts, USA). High-quality genomic DNA was obtained
564 for 978 out of the 1,000 donors.

565

566 **DNA methylation profiling and data quality controls**

567 Extracted genomic DNA was treated with the EZ DNA Methylation Kit (catalog #: D5001; Zymo
568 Research, California, USA). Bisulfite-converted DNA was applied to the Infinium
569 MethylationEPIC BeadChip (catalog #: WG-317-1003; Illumina, California, USA), using the
570 manufacturer's standard conditions. The MethylationEPIC BeadChip measures 5mC levels at
571 866,836 CpG sites in the human genome. Raw IDAT files were processed with the minfi R
572 package¹⁰⁰. All samples showed average detection *P*-values < 0.005. No sample showed a mean of
573 methylated intensity signals lower than 3 × standard deviations (SD) from the cohort average.
574 Therefore, no samples were excluded based on detection *P*-values or methylated intensity signals.
575 The sex predicted from 5mC signals on sex chromosomes matched the declared sex for all samples
576 (Supplementary Fig. 1a). Using the 59 control SNPs included in the MethylationEPIC array, a
577 single sample showed high genotype discordance with the genome-wide SNP array data (see
578 'Genome-wide DNA genotyping' section) and was thus excluded (Supplementary Fig. 1b).
579 Unmethylated and methylated intensity signals were converted to M-values. A total of 2,930 probes
580 with >1% missingness (i.e., detection *P*-value > 0.05 for more than 1% of donors) were excluded
581 and remaining missing data (missingness = 0.0038%) were imputed by mean substitution. Using the
582 irlba R package, Principal Component Analysis (PCA) of M values identified nine outlier samples,
583 including eight that were processed on the same array (Supplementary Fig. 1c), which were also
584 excluded. The "noob" background subtraction method¹⁰¹ was applied on M values for the remaining
585 968 samples, which showed highly consistent epigenome-wide DNA methylation profiles
586 (Supplementary Fig. 1d,e).

587 To identify batch effects on the DNA methylation data, we searched for the factors that were the
588 most associated with the top 20 PCs of the PCA of noob-corrected M values. We used a linear
589 mixed model that included age, sex and cytomegalovirus (CMV) serostatus as fixed effects, and
590 slide position and sample plate as random effects. The models were fitted with the lme4 R
591 package¹⁰². Strong associations were observed between the first four PCs and slide position and
592 sample plate (Supplementary Fig. 1f, g). M values were thus corrected for these two batch effects
593 using the ComBat function, from the sva R package¹⁰³. After ComBat correction, the ten first PCs
594 of a PCA of M values were associated with factors known to affect DNA methylation, including
595 blood cell composition, age and sex (Supplementary Fig. 1h-j), indicating no other, strong batch
596 effect on the data (see section 'Associations with principal components of DNA methylation').

597 M-values were converted to β values, considering that $\beta = 2^M / (2^M + 1)$. Because outlier 5mC
598 values due to measurement error could inflate the type I error rate of regression models, we
599 excluded, for each CpG site, M or β values that were greater than $5 \times$ SD from the population
600 average, corresponding to $<0.1\%$ of all measures. We also excluded (i) 83,380 non-specific probes
601 that share $>90\%$ sequence identity with several genomic regions (see details in¹⁰⁴), (ii) 118,575
602 probes that overlap a SNP that is within the 50 pb surrounding the CpG site and has a MAF $>1\%$ in
603 the Milieu Intérieur cohort or in European populations from the 1,000 Genomes project¹⁰⁵, (iii) 558
604 probes that were absent from the Illumina annotations version 1.0 B4 and (iv) 16,876 probes located
605 on sex chromosomes. As a result, the final, quality-controlled data was composed of 968 donors
606 profiled at 644,517 CpG sites.

607

608 **Flow cytometry**

609 Immune cell proportions were measured using ten eight-color flow-cytometry panels²⁵. The
610 acquisition of cells was performed using two MACSQuant analyzers, which were calibrated using
611 MacsQuant calibration beads (Miltenyi, Germany). Flow cytometry data were generated using
612 MACSQuantify software version 2.4.1229.1. The mqd files were converted to FCS compatible
613 format and analyzed by FlowJo software version 9.5.3. A total of 110 cell proportions were
614 exported from FlowJo. Protocols, panels, staining antibodies and quality control filters used for
615 flow cytometry analyses are detailed elsewhere²⁵. Abnormal lysis or staining were systematically
616 flagged by trained experimenters. We removed outliers by using a scheme detailed previously²⁵. We
617 used a distance-based approach that, for each cell-type, removes observations in the right tail if the
618 distance to the closest observation in the direction of the mean is larger than 20% of the range of the
619 observations. Similarly, observations in the left tail were removed if the distance to the closest
620 observation in the direction of the mean is more than 15% than the range the observations. We
621 removed 22 observations in total, including a maximum of 8 observations for a single cell-type (i.e.,
622 for the proportion of neutrophils). Problems in flow cytometry processing, such as abnormal lysis or
623 staining, were systematically flagged by trained experimenters, which resulted in 8.7% missing
624 data. Because imputing missing data for donors who show large missingness could be inaccurate,
625 we excluded 74 donors with no data for the T cell panel. Finally, the remaining missing data were
626 imputed using the random forest-based missForest R package¹⁰⁶.

627

628 **Genome-wide DNA genotyping**

629 The 1,000 Milieu Intérieur donors were genotyped on both the HumanOmniExpress-24 and the
630 HumanExome-12 BeadChips (Illumina, California, USA), which include 719,665 SNPs and

631 245,766 exonic SNPs, respectively. Average concordance rate between the two genotyping arrays
632 was 99.9925%. The combined data set included 732,341 high-quality polymorphic SNPs. After
633 genotype imputation and quality-control filters²⁵, a total of 11,395,554 SNPs was further filtered for
634 minor allele frequencies > 5%, yielding a data set composed of 1,000 donors and 5,699,237 SNPs
635 for meQTL mapping. Ten pairs of first to third-degree related donors were detected with KING 1.9
636 (ref.¹⁰⁷). Out of the 894 donors whose blood methylome and blood cell composition were accurately
637 profiled, 884 unrelated donors were kept for subsequent analyses.

638

639 **Immune cell proportions**

640 One of the key questions in this study is whether differences in 5mC levels observed with respect to
641 different factors are due to epigenetic changes occurring within cell-types or if they in fact reflect
642 changes in blood cell composition. To answer this question, we considered the proportions of 16
643 major subsets of blood: naïve, central memory (CM), effector memory (EM) and terminally
644 differentiated effector memory (EMRA) subsets of CD4⁺ and CD8⁺ T cells, CD4⁻CD8⁻ T cells, B
645 cells, dendritic cells, natural killer (NK) cells, monocytes, neutrophils, basophils and eosinophils²⁵.
646 As these cellular proportions were measured by flow cytometry using a hierarchical gating
647 strategy²⁵, they are expected to sum to one. Yet, because of measurement errors, cell fractions do
648 not exactly sum to one in all donors. For a measure of proportion of a given cell subset in a given
649 donor, we therefore used the absolute count of the cell-type divided by the sum of absolute counts
650 of all the 16 measured cell subsets. We used the same approach when considering a reduced set of
651 six major cell-types, comprising neutrophils, monocytes, NK cells, B cells, and CD4⁺ and CD8⁺ T
652 cells, for comparison purposes.

653

654 **Compositional analysis of cellular composition**

655 We sought to study the association between 5mC levels and blood cell composition, experimentally
656 measured by flow cytometry. However, the 16 measured cellular proportions are constrained to be
657 positive and to sum to one. Consequently, a change in one cellular proportion must necessarily
658 change one or more of the other cellular proportions, complicating the interpretation of parameters
659 estimated from linear regression models with measured immune cell proportions as
660 predictors^{28,108,109}. Here, we investigated instead the effect of balances, which are transformations of
661 cell-type proportions that can be seen as a generalization of the logit-transform. These balances
662 model the effect of a relative change between two groups of cell-types. They are defined in a
663 hierarchical manner of increasing granularity, by a sequential binary partition (SBP) of the 16
664 measured cell-types, generating 15 balances in total (Supplementary Data 2). As an example, we

665 describe the first two balances. The other balances are defined in an analogous manner according to
666 the SBP and the general procedure detailed elsewhere¹⁰⁸. The first balance captures the relative
667 effect on 5mC levels of the myeloid cell-types compared to the lymphoid cell-types. Of the 16
668 measured cell-types, five are myeloid and eleven are lymphoid. Let $c_i^{M_1}, \dots, c_i^{M_5}$ be the measured
669 myeloid proportions and $c_i^{L_1}, \dots, c_i^{L_{11}}$ be lymphoid proportions for the i :th individual. The first
670 balance predictor for that individual is defined by

671

$$b_i^1 = \sqrt{\frac{5 \times 11}{5 + 11}} \log \left\{ \frac{\prod_{m=1}^5 c_i^{M_m}}{\prod_{l=1}^{11} c_i^{L_l}} \right\}, \quad (1)$$

672

673 The second balance is defined within the lymphoid group and captures the relative effect on
674 5mC levels of T cells with respect to NK cells and B cells. Let $c_i^{T_1}, \dots, c_i^{T_9}$ be the measured
675 proportions of the nine types of T cells and let c_i^B and c_i^{NK} be proportions of B cells and NK cells.
676 The balance contrasting T cells with NK cells and B cells is given by

677

$$b_i^2 = \sqrt{\frac{9 \times 2}{9 + 2}} \log \left\{ \frac{\prod_{m=1}^9 c_i^{T_m}}{c_i^B c_i^{NK}} \right\}. \quad (2)$$

678

679 All balances were computed from the SBP using the robCompositions R package¹¹⁰. To evaluate the
680 validity of our approach, we compared the estimated effects on 5mC levels of balances contrasting
681 two groups of cell-types with the measured differences in 5mC levels between the same two groups,
682 obtained from MethylationEPIC data in sorted cell-types²⁹ and found strong correlations ($R > 0.6$;
683 Supplementary Fig. 2 and Supplementary Data 2). We further evaluated the accuracy of our
684 approach by performing a simulation study. First, we simulated 5mC levels based on observed cell
685 composition data and evaluated how the balances capture 5mC differences in the relevant cell-
686 types. Second, we simulated cell composition data from a Dirichlet distribution and again evaluated
687 that regression models including the balances as predictors give the expected results
688 (Supplementary Notes).

689 The 15 balances were used to investigate the effects of immune cell composition on 5mC levels
690 at individual CpG sites (see section ‘Epigenome-wide association study of cell composition’) and
691 on principal components of epigenome-wide DNA methylation levels (see section ‘Associations
692 with principal components of DNA methylation’).

693

694 **Epigenome-wide association study of cell composition**

695 To investigate how immune cell composition affects the blood DNA methylome, we investigated
696 effects of cell-type balances on 5mC levels at each CpG site. For the p :th CpG site and the i :th
697 individual, introduce observed 5mC levels y_i^p measured on the M value scale. Let \mathbf{b}_i be a vector of
698 15 cell-type balances with corresponding parameter vector $\boldsymbol{\beta}_b^p$. Let the vector \mathbf{SNP}_i^p contain the
699 significant local SNP with the smallest P -value and all independently associated remote SNPs (see
700 section ‘Local meQTL mapping analyses’ and section ‘Remote meQTL mapping analyses’) with
701 corresponding parameter vector $\boldsymbol{\beta}_{\text{SNP}}^p$. We performed an epigenome-wide association analysis of
702 cellular composition by fitting the models,

703

$$y_i^p = \mu^p + \mathbf{b}_i^t \boldsymbol{\beta}_b^p + (\mathbf{SNP}_i^p)^t \boldsymbol{\beta}_{\text{SNP}}^p + \varepsilon_i, \quad (3)$$

704

705 where $\varepsilon_i \sim (0, \sigma_p^2)$. Models were fitted by ordinary least squares. For each balance in \mathbf{b}_i (see Eq.
706 (1) and Eq. (2) for examples), the parameters in $\boldsymbol{\beta}_b^p$ are interpreted as the change in 5mC levels for
707 an increase in the first cell-type group and the corresponding decrease in the second cell-type group.

708

709 **Associations with principal components of DNA methylation**

710 To evaluate how principal components (PCs) of DNA methylation levels are related to cell
711 composition, we first computed PCs of 5mC levels at 644,517 CpG sites, with the irlba R package.
712 Let y_i^k be the observed value of the k :th PC of the DNA methylation data and \mathbf{b}_i a vector of 15 cell-
713 type balances measured for individual i with the corresponding parameter vector $\boldsymbol{\beta}_b^k$. Given that we
714 observed variability in 5mC levels across dates of blood draw, we included them as random effects.
715 Let j be the day of blood draw for the i :th individual. The model we used to estimate the effects of
716 cellular composition on PCs of DNA methylation was,

717

$$y_i^k = \mu^k + \mathbf{b}_i^t \boldsymbol{\beta}_b^k + \text{DateOfSampling}_{j(i)} + \varepsilon_i^k, \quad (4)$$

718

719 with $\text{DateOfSampling}_{j(i)} \sim \mathcal{N}(0, \tau_k^2)$ and $\varepsilon_i \sim (0, \sigma_k^2)$. The models were fitted with the lme4 R
720 package¹⁰².

721

722 To evaluate how PCs of DNA methylation levels are related to the candidate non-heritable
factors, i.e., age, sex, smoking status, CMV serostatus, introduce the variables Age_i , Woman_i ,

723 Exsmoker_{*i*}, Smoker_{*i*} and CMV_{*i*} with corresponding parameters β_{Age}^k , β_{Woman}^k , $\beta_{\text{Exsmoker}}^k$, β_{Smoker}^k
724 and β_{CMV}^k . Let PC1_{*i*} and PC2_{*i*} be the two first PCs of the genotype matrix. Let \mathbf{c}_i be a vector of 15
725 measured cell proportions, excluding neutrophils because of the sum-to-one constraint, and $\boldsymbol{\beta}_c^k$ the
726 corresponding parameter vector. The model we used to estimate the effects of non-genetic factors
727 on PCs of DNA methylation was,

728

$$\begin{aligned} y_i^k = & \mu^k + \mathbf{c}_i^t \boldsymbol{\beta}_c^k + \text{Age}_i \beta_{\text{Age}}^k + \text{Woman}_i \beta_{\text{Woman}}^k + \text{Exsmoker}_i \beta_{\text{Exsmoker}}^k \\ & + \text{Smoker}_i \beta_{\text{Smoker}}^k + \text{CMV}_i \beta_{\text{CMV}}^k + \text{PC1}_i \beta_{\text{PC1}}^k + \text{PC2}_i \beta_{\text{PC2}}^k \\ & + \text{DateOfSampling}_{j(i)} + \varepsilon_i^k. \end{aligned} \quad (5)$$

729

730 The models were fitted with the lme4 R package¹⁰². Inference was performed using the Kenward-
731 Roger *F*-test approximation for linear mixed models, implemented in the pbkrtest R package¹¹¹.

732

733 **Epigenome-wide association studies of non-genetic factors**

734 We assessed the effects of 141 non-genetic variables (Supplementary Data 1) on the blood DNA
735 methylome of adults. The measured 5mC levels at a CpG site are the average of the DNA
736 methylation state at this CpG site of all cells in the blood sample. Many of the 141 candidate
737 variables might influence cell composition, which will cause a corresponding change in 5mC levels.
738 We denote this effect the “(cell-composition-)mediated effect”. In addition, the variable might alter
739 5mC levels within individual cells, or within cell-types. We denote this effect the direct effect (see
740 Supplementary Fig. 11 for a schematic directed acyclic graph of the system). Several factors are
741 known to have a large effect on blood cell composition in healthy donors, the most important being
742 age, sex, CMV serostatus and smoking²⁵. As an added complexity, these factors are also associated
743 with most of the other variables in the study. Based on this framework, we investigated four
744 questions, each one targeted by a separate statistical model.

745

746 The total effect

747 The total effect includes both changes in 5mC levels induced by changes in cellular composition
748 (i.e., cell-composition-mediated effects) and those induced within cell-types (i.e., direct effects). For
749 each variable of interest x and each CpG site, the total effect was estimated in a regression model
750 including, as response variable, the 5mC levels of the CpG site on the M value scale and, as
751 predictors, x_i , a nonlinear age term of 3 DoF natural splines, sex, CMV serostatus, smoking status,
752 the significant local SNP with the smallest *P*-value, independently associated remote SNPs and the

753 first two PCs of the genotype matrix. Again, since we observed variability in 5mC levels across
 754 dates of blood draw, we included them as a random effect term. For the p :th CpG site, let y_i^p be the
 755 5mC levels of the i :th individual on the M value scale, $f_{\text{Age}}^p(\text{Age}_i)$ a nonlinear age term of 3 DoF
 756 natural splines and \mathbf{SNP}_i^p a vector of the minor allele counts for the significant local SNP with the
 757 smallest P -value and independently associated remote SNPs, with corresponding parameter vector
 758 $\boldsymbol{\beta}_{\text{SNP}}^p$. The total effect of the variable x_i was estimated by the corresponding parameter β_x^p in the
 759 models,
 760

$$\begin{aligned}
 y_i^p = & \mu^p + x_i \beta_x^p + f_{\text{Age}}^p(\text{Age}_i) + \text{Woman}_i \beta_{\text{Woman}}^p + \text{Exsmoker}_i \beta_{\text{Exsmoker}}^p \\
 & + \text{Smoker}_i \beta_{\text{Smoker}}^p + \text{CMV}_i \beta_{\text{CMV}}^p + \text{PC1}_i \beta_{\text{PC1}}^p + \text{PC2}_i \beta_{\text{PC2}}^p \\
 & + (\mathbf{SNP}_i^p)^t \boldsymbol{\beta}_{\text{SNP}}^p + \text{DateOfSampling}_{j(i)} + \varepsilon_i^p,
 \end{aligned} \tag{6}$$

761
 762 where $\text{DateOfSampling}_{j(i)} \sim \mathcal{N}(0, \tau_p^2)$ and $\varepsilon_i \sim (0, \sigma_p^2)$. The effect of aging was tested in models
 763 with x removed and the non-linear age term replaced by a linear one. The effects of sex, smoking
 764 status and CMV serostatus were tested in models where we removed x . For variables relating to
 765 women only (e.g., age of menarche), we excluded men from the analysis and removed
 766 $\text{Woman}_i \beta_{\text{Woman}}^p$. The models were fitted with the lme4 R package¹⁰². Hypothesis tests were
 767 performed using the Kenward-Roger approximation of the F -test for linear mixed models,
 768 implemented in the pbkrtest R package¹¹¹.

769

770 The direct effect

771 Let the vector \mathbf{c}_i be measured proportions of the 15 immune cell-types, excluding neutrophils, for
 772 the i :th individual and $\boldsymbol{\beta}_c^p$ the corresponding parameter vector. Using the same notation as for the
 773 total effect, the direct effect of the variable x_i was estimated by β_x^p in the models,
 774

$$\begin{aligned}
 y_i^p = & \mu^p + x_i \beta_x^p + \mathbf{c}_i^t \boldsymbol{\beta}_c^p + f_{\text{Age}}^p(\text{Age}_i) + \text{Woman}_i \beta_{\text{Woman}}^p + \text{Exsmoker}_i \beta_{\text{Exsmoker}}^p \\
 & + \text{Smoker}_i \beta_{\text{Smoker}}^p + \text{CMV}_i \beta_{\text{CMV}}^p + \text{PC1}_i \beta_{\text{PC1}}^p + \text{PC2}_i \beta_{\text{PC2}}^p \\
 & + (\mathbf{SNP}_i^p)^t \boldsymbol{\beta}_{\text{SNP}}^p + \text{DateOfSampling}_{j(i)} + \varepsilon_i^p,
 \end{aligned} \tag{7}$$

775

776 We also tested the interaction effect of sex, CMV serostatus and smoking status with age by
 777 including one interaction term at a time in the model specified in Eq. (7). The models were fitted

778 with the lme4 R package¹⁰². Hypothesis tests were performed by the Kenward-Roger approximation
779 of the F -test for linear mixed models, implemented in the pbkrtest R package¹¹¹.

780

781 The mediated effect

782 The cell-composition-mediated effect was estimated as the effect on 5mC levels mediated by
783 changes in proportions of the 16 cell subsets due to the given factor. We estimated the mediated
784 effect of aging, sex, variables related to smoking, CMV serostatus and heart rate. The mediated
785 effect was estimated using a two-stage procedure. First, we fitted models with measured proportions
786 of immune cells as response variables. Let \mathbf{c}_i be a vector of measured proportions of the 15 blood
787 subsets, excluding neutrophils. Let c_i^n denote the n :th entry of the vector \mathbf{c}_i , i.e., the measured
788 proportion of the n :th cell-type for the i :th individual. Introduce the vector \mathbf{k}_i of covariate values for
789 the i :th individual, including age (3 DoF spline with an entry for each term), sex, smoking, CMV
790 serostatus and ancestry (2 PCs), but excluding the variable of interest x_i (mediated effect of aging
791 was estimated with a linear term). For the model of the n :th cell-type, let $\boldsymbol{\beta}_k^n$ be the parameter vector
792 for the covariate vector \mathbf{k}_i and β_x^n the parameter for the variable of interest x_i . In the first stage, we
793 fitted the models,

794

$$E\{c_i^n \mid x_i, \mathbf{k}_i\} = \beta_0 + x_i \beta_x^n + \mathbf{k}_i^t \boldsymbol{\beta}_k^n, \quad n = 1, \dots, 15. \quad (8)$$

795

796 Next, let y_i^p be 5mC levels in the M value scale for the p :th CpG site, θ_x^p the parameter for the
797 variable of interest, and $\boldsymbol{\theta}_c^p$ and $\boldsymbol{\theta}_k^p$ parameter vectors for the effects of cell proportions and
798 covariates. In the second stage, we fitted the models,

799

$$E\{y_i^p \mid x_i, \mathbf{c}_i, \mathbf{k}_i\} = \theta_0^p + x_i \theta_x^p + \mathbf{c}_i^t \boldsymbol{\theta}_c^p + \mathbf{k}_i^t \boldsymbol{\theta}_k^p. \quad (9)$$

800

801 The mediated effect of x_i on DNA methylation was estimated by $\boldsymbol{\beta}_x^t \boldsymbol{\theta}_c^p$ (ref.³⁴). Inference was
802 performed by the parametric bootstrap.

803

804 The direct effects adjusted by deconvolution methods

805 To compute the IDOL and Houseman-adjusted effects, we estimated proportions of CD4⁺ and CD8⁺
806 T cells, B cells, NK cells, monocytes and neutrophils, using the estimateCellCounts2 function in the
807 FlowSorted.Blood.EPIC package with either Houseman et al.'s CpG sites, or IDOL optimized CpG
808 sites¹¹². For age, sex, smoking status, CMV serostatus, heart rate, ear temperature and hour of blood

809 draw, we estimated the IDOL- and Houseman-adjusted effect by adjusting for estimated 5
810 proportions in the model specified by Eq. (7), instead of the 15 measured proportions, excluding
811 neutrophils because of the sum-to-one constraint. To compute the EPIC IDOL-Ext-adjusted effects,
812 we estimated proportions of 12 major cell-types in blood, including CD4⁺ and CD8⁺ T cells, naïve
813 and differentiated subtypes of CD4⁺ and CD8⁺ T cells, neutrophils, monocytes, basophils,
814 eosinophils, NK cells, regulatory T cells, naïve and memory B cells, using the IDOL-Ext reference
815 matrix in the estimateCellCounts2 function from the FlowSorted.BloodExtended.EPIC R package²⁹.
816 We estimated the IDOL-Ext-adjusted effect by including 11 estimated proportions in Eq. (7) instead
817 of the 15 measured proportions, excluding neutrophils because of the sum-to-one constraint.
818 Finally, for comparison purposes, we also computed the association between non-genetic factors
819 and 5mC levels by adjusting, in Eq. (7), for the proportions of the 5 major cell-types measured by
820 flow cytometry, instead of the 15 measured proportions, excluding again neutrophils.

821

822 **Prediction of CMV serostatus**

823 We built a prediction model to estimate CMV serostatus from DNA methylation data using elastic
824 net regression for binary data¹¹³, implemented in the glmnet R package¹¹⁴. We included all CpG
825 sites as predictors in the model, including those on the X and Y chromosomes. The model was built
826 from 863,906 CpG sites in 969 samples. The elastic net model has two tuning parameters that
827 determine the degree of regularization of the predictor function. We selected both tuning parameters
828 by two-dimensional five times repeated cross-validation over the two parameters. The final model
829 fitted on the full data set includes 547 CpG sites with non-zero parameters.

830

831 **Detection of the dispersion of DNA methylation with age**

832 To estimate changes in dispersion of 5mC levels with age, we fitted regression models where the
833 residual variance depends on age. Let y_i^p be 5mC levels on the M value scale for the p :th CpG site
834 and the i :th individual. Using similar notations as above, we estimated the dispersion effect of age
835 by the parameter θ^p in the models,

836

$$\begin{aligned} y_i^p = & \mu^p + \mathbf{c}_i^t \boldsymbol{\beta}_c^p + (\mathbf{SNP}_i^p)^t \boldsymbol{\beta}_{\text{SNP}}^p + f_{\text{Age}}^p(\text{Age}_i) + \text{Woman}_i \beta_{\text{Woman}}^p \\ & + \text{Exsmoker}_i \beta_{\text{Exsmoker}}^p + \text{Smoker}_i \beta_{\text{Smoker}}^p + \text{CMV}_i \beta_{\text{CMV}}^p + \text{PC1}_i \beta_{\text{PC1}}^p \\ & + \text{PC2}_i \beta_{\text{PC2}}^p + \varepsilon_i^p, \end{aligned} \quad (10)$$

837

838 where

$$\varepsilon_i^p \sim \mathcal{N}(0, \sigma_{i,p}^2), \log \sigma_{i,p} = \tau^p + \text{Age}_i \theta^p. \quad (11)$$

839

840 We devised a hypothesis test for θ by a likelihood ratio test comparing the model in Eq. (11), to a
841 model with

$$\varepsilon_i^p \sim \mathcal{N}(0, \sigma_p^2), \log \sigma_p = \tau^p. \quad (12)$$

842

843 As a sensitivity analysis, we also fitted a model with

844

$$\varepsilon_i^p \sim \mathcal{N}(0, \sigma_{i,p}^2), \log \sigma_{i,p} = \tau^p + \text{Age}_i \theta^p + \mathbf{c}_i^t \boldsymbol{\beta}_c^p. \quad (13)$$

845

846 In this case, the hypothesis test for θ was done by comparing to a model with

847

$$\varepsilon_i^p \sim \mathcal{N}(0, \sigma_{i,p}^2), \log \sigma_{i,p} = \tau^p + \mathbf{c}_i^t \boldsymbol{\beta}_c^p. \quad (14)$$

848

849 These models were fitted with the `gamlss` R package¹¹⁵.

850

851 **Local meQTL mapping analyses**

852 Local meQTL mapping was performed using the `MatrixEQTL` R package¹¹⁶. Association was tested
853 for each CpG site and each SNP in a 100-Kb window around the CpG site, by fitting a linear
854 regression model assuming an additive allele effect. Models included, as predictors, the 15 immune
855 cell proportions, a nonlinear age term encoded by 3 degrees-of-freedom (DoF) natural splines, sex,
856 smoker status, ex-smoker status and CMV serostatus. We also adjusted for the top two PCs of a
857 PCA of the genotype data. We did not include more PCs because of the low population substructure
858 observed in the cohort²⁵. For the i :th individual and the p :th CpG site, let y_i^p be the measured 5mC
859 levels on the M value scale, $\text{SNP}_i^{p,m}$ the minor allele count of the m :th tested SNP for the CpG site
860 and $f_{\text{Age}}^{p,m}(\text{Age}_i)$ a nonlinear age term of natural splines. Moreover, let the vector \mathbf{c}_i be measured
861 proportions of the 15 immune cell-types for the i :th individual, excluding neutrophils, and $\boldsymbol{\beta}_c^{p,m}$ the
862 corresponding parameter vector. The additive allele effect of the SNP was estimated by the
863 parameter $\beta_m^{p,m}$ in the models,

864

$$\begin{aligned}
 y_i^p = & \mu^{p,m} + \text{SNP}_i^{p,m} \beta_m^{p,m} + f_{\text{Age}}^{p,m}(\text{Age}_i) + \text{Woman}_i \beta_{\text{Woman}}^{p,m} + \text{Exsmoker}_i \beta_{\text{Exsmoker}}^{p,m} \\
 & + \text{Smoker}_i \beta_{\text{Smoker}}^{p,m} + \text{CMV}_i \beta_{\text{CMV}}^{p,m} + \text{PC1}_i \beta_{\text{PC1}}^{p,m} + \text{PC2}_i \beta_{\text{PC2}}^{p,m} + \mathbf{c}_i^t \boldsymbol{\beta}_c^{p,m} \\
 & + \varepsilon_i^{p,m},
 \end{aligned} \tag{15}$$

865

866 where $\varepsilon_i^{p,m}$ is a symmetrical zero-mean distribution with constant variance.

867

868 **Remote meQTL mapping analyses**

869 Testing all possible associations between 644,517 CpG sites and 5,699,237 SNPs would require
 870 performing 3,769 billion statistical tests. To reduce the multiple testing burden, remote meQTL
 871 mapping was conducted on a selection of 50,000 CpG sites with the highest residual variance in the
 872 model described in Eq. (15), but with m indexing in this case only the most associated local SNP for
 873 the p :th CpG site. For each of the 50,000 selected CpG sites, we then fitted one model per SNP
 874 located outside of a 1-Mb window around the CpG site. For each SNP-CpG pair, we estimated the
 875 additive allele effect of the remote SNP using the model specified in Eq. (15) but with m now
 876 indexing remote SNPs for the p :th CpG site. Both local and remote meQTL mapping tests were
 877 corrected for multiple testing by the Bonferroni adjustment.

878

879 **Detection of independent remote meQTLs**

880 We designed the following scheme to compute a set Φ of independently associated remote SNPs for
 881 each CpG site, where all such SNPs are associated with 5mC levels y^p at the p :th CpG site,
 882 conditional on the most associated local SNP and other SNPs in Φ . Define X_1 to be the set of SNPs
 883 with a remote association to y^p and let x^0 be the most associated significant local SNP, if it exists.
 884 The set X_1 typically includes several SNPs that are in linkage disequilibrium (LD). The algorithm
 885 uses an iterative procedure to build sets M_j of SNPs, where in the j :th iteration, SNPs that are not
 886 associated with 5mC levels at the CpG site conditional on SNPs included in M_{j-1} are discarded,
 887 while the most associated is retained in M_j . In the final step, the set Φ is constructed by elements of
 888 the final set M that are associated with 5mC levels at the CpG site conditional on all the other
 889 elements in M . Intuitively, Φ consists of the most associated SNP in each LD block. The algorithm
 890 is given in pseudocode in Algorithm (1), where the condition $\beta^p \neq 0$ is determined by an F -test on
 891 the level $\alpha = 10^{-6}$.

892

Algorithm (1): Forming a set of remote independently associated SNPs with a CpG site

If the CpG site is under local genetic control then let $M_1 = x_0$, otherwise let $M_1 = \emptyset$

Repeat for $j = 1, 2, \dots$

$$P = \{x \in X_j \setminus M_j: \beta_x^p \neq 0 \text{ in } y_i^p = \mu^p + x_i \beta_x^p + \sum_{z \in M_j} z_i \beta_z^p + \varepsilon_i^p, \varepsilon_i^p \sim (0, \sigma_p^2)\}$$

If $P = \emptyset$ Exit

$$X_{j+1} = P$$

$$M_{j+1} = M_j \cup \{x: x \text{ SNP with the smallest } P\text{-value in } P\}$$

End

$$\Phi = \{x \in M_{j+1} \setminus x_0: \beta_x^p \neq 0 \text{ in } y_i^p = \mu^p + x_i \beta_x^p + \sum_{z \in M_{j+1} \setminus \{x\}} z_i \beta_z^p + \varepsilon_i^p, \varepsilon_i^p \sim (0, \sigma_p^2)\}$$

893

894 Cell-type-dependent effects of genetic and non-genetic factors on DNA methylation

895 To investigate whether the effects of a factor on DNA methylation depend on the proportion of
 896 myeloid cells in blood, we fitted models that included an interaction term between the factor of
 897 interest (i.e., age, sex, smoking status, CMV serostatus and genetic variants) and the proportion of
 898 myeloid cells, c_i^m , defined as the sum of the proportions of cell-types from the myeloid lineage.
 899 With the same notations as above, but with y_i^p being 5mC levels on the β value scale for the p :th
 900 CpG site and the i :th individual, we estimated the cell-type-dependent effects of non-genetic factors
 901 by fitting the models,

902

$$\begin{aligned} y_i^p = & \mu^p + \text{Age}_i \beta_{\text{Age}}^p + \text{CMV}_i \beta_{\text{CMV}}^p + \text{Woman}_i \beta_{\text{Woman}}^p + \text{Smoker}_i \beta_{\text{Smoker}}^p + \text{PC1}_i \beta_{\text{PC1}}^p \\ & + \text{PC2}_i \beta_{\text{PC2}}^p + c_i^m \beta_{c^m}^p \\ & + c_i^m \left(\text{Woman}_i \theta_{\text{Woman}}^p + \text{Age}_i \theta_{\text{Age}}^p + \text{Smoker}_i \theta_{\text{Smoker}}^p + \text{CMV}_i \theta_{\text{CMV}}^p \right) \\ & + \varepsilon_i^p. \end{aligned} \quad (16)$$

903

904 We also investigated whether the effect of genotypes could be dependent on the proportion of
 905 myeloid cells in the sample. For the p :th CpG site and the i :th individual, let $\text{SNP}_i^{p,k}$ be the minor
 906 allele counts of the significant local SNP with the smallest P -value and independently associated
 907 remote SNPs. In this case, we also use 5mC levels on the β value scale. We estimated the cell-type-
 908 dependent effects of genetic factors by fitting the models,

909

$$\begin{aligned}
 y_i^p = & \mu^p + f_{\text{Age}}^p(\text{Age}_i) + \text{CMV}_i \beta_{\text{CMV}}^p + \text{Woman}_i \beta_{\text{Woman}}^p + \text{Smoker}_i \beta_{\text{Smoker}}^p \\
 & + \text{PC1}_i \beta_{\text{PC1}}^p + \text{PC2}_i \beta_{\text{PC2}}^p + \mathbf{c}_i^m \beta_{\mathbf{c}^m}^p + \sum_k \text{SNP}_i^{p,k} \beta_{\text{SNP}^{p,k}} \\
 & + \mathbf{c}_i^m \left(\sum_k \text{SNP}_i^{p,k} \theta_{\text{SNP}^{p,k}} \right) + \varepsilon_i^p.
 \end{aligned} \tag{17}$$

910

911 Inference in both cases was done by Wald tests with heteroscedasticity-consistent standard
 912 errors estimated by the sandwich R package¹¹⁷.

913

914 **Detection of gene × environment interactions**

915 We tested whether age, sex, CMV serostatus, smoking status or CRP levels could have a genotype-
 916 dependent effect on the DNA methylome. For the i :th individual and the p :th CpG site, let y_i^p be the
 917 5mC levels on the M value scale, $\text{SNP}_i^{p,k}$, $k = 1, \dots, K^p$, the minor allele counts of the significant
 918 local meQTL with the lowest P -value and the $K^p - 1$ independently associated remote meQTLs,
 919 and \mathbf{c}_i the vector of 15 measured immune cell proportions with corresponding parameter vector $\beta_{\mathbf{c}}^p$.
 920 Interaction effects were estimated for each CpG site in the model,

921

$$\begin{aligned}
 E\{y_i^p \mid \text{SNP}_i^{p,1}, \dots, \text{SNP}_i^{p,K^p}, \text{Age}_i, \text{Woman}_i, \text{Smoker}_i, \text{CMV}_i\} \\
 = \mu^p + \sum_{k=1}^{K^p} \text{SNP}_i^{p,k} \beta_{\text{SNP}^{p,k}} + \mathbf{c}_i^t \beta_{\mathbf{c}}^p + \text{PC1}_i \beta_{\text{PC1}}^p + \text{PC2}_i \beta_{\text{PC2}}^p + \text{Age}_i \beta_{\text{Age}}^p \\
 + \text{Woman}_i \beta_{\text{Woman}}^p + \text{Smoker}_i \beta_{\text{Smoker}}^p + \text{CMV}_i \beta_{\text{CMV}}^p \\
 + \sum_{k=1}^{K^p} \text{SNP}_i^{p,k} \left(\text{Age}_i \theta_{\text{Age}}^{p,k} + \text{Woman}_i \theta_{\text{Woman}}^{p,k} + \text{Smoker}_i \theta_{\text{Smoker}}^{p,k} \right. \\
 \left. + \text{CMV}_i \theta_{\text{CMV}}^{p,k} \right)
 \end{aligned} \tag{18}$$

922

923 We investigated effects of CRP levels in a separate model that simply added a log-transformed
 924 CRP term to Eq. (18). Inference was done by Wald tests with heteroscedasticity-consistent standard
 925 errors estimated by the sandwich R package¹¹⁷.

926

927 **Estimation of proportions of explained 5mC variance**

928 According to our analyses, 5mC levels in the healthy population are mainly associated with local
 929 genetic variation, blood cell composition, age, sex, smoking, CMV infection and CRP levels. We

930 grouped these variables into four categories: genetic, cell composition, intrinsic (age and sex) and
 931 exposures (smoking, CMV infection and CRP levels). For the p :th CpG site and the i :th individual,
 932 we collected observations of the minor allele count for the most associated local SNP in $x_i^{p,g}$, the
 933 proportions of the 15 cell-types, excluding neutrophils, in the vector \mathbf{x}_i^c , intrinsic factors (sex and
 934 natural spline expanded values of age) in the vector \mathbf{x}_i^{in} and exposures (smoking status, CMV
 935 serostatus and log-transformed CRP levels) in the vector \mathbf{x}_i^e , with corresponding parameters β_g^p , β_c^p ,
 936 β_{in}^p and β_e^p . We interpret here log-transformed CRP levels as a proxy measure of the exposure of
 937 chronic low-grade inflammation. For each CpG site, we define linear predictor terms by

938

$$f_g^p(x_i^{p,g}) = x_i^{p,g} \beta_g^p, \quad (19)$$

939

$$f_c^p(\mathbf{x}_i^c) = (\mathbf{x}_i^c)^t \boldsymbol{\beta}_c^p, \quad (20)$$

940

$$f_{in}^p(\mathbf{x}_i^{in}) = (\mathbf{x}_i^{in})^t \boldsymbol{\beta}_{in}^p, \quad (21)$$

941

$$f_e^p(\mathbf{x}_i^e) = (\mathbf{x}_i^e)^t \boldsymbol{\beta}_e^p \quad (22)$$

942

943 These functions vary in their degrees of freedom, so to get a fair comparison between them, we
 944 estimated group effect sizes as the out-of-sample proportion of variance explained by each group
 945 predictor. This estimation is done by indexing samples into two disjoint index groups I_1 and I_2 ,
 946 fitting the models on samples from I_1 , and evaluating the prediction accuracy on samples from I_2 .
 947 Let y_i^p be 5mC levels for the p :th CpG site on the β value scale. Take cell composition as example.
 948 To compute the total effect of cell composition on 5mC levels at the CpG site, we first fit a model
 949 with individuals in I_1 ,

950

$$y_i^{p,c} = \mu^p + (\mathbf{x}_i^c)^t \boldsymbol{\beta}_c^p, \quad i \in I_1 \quad (23)$$

951

952 with parameters $\hat{\beta}_c^p$ and $\hat{\mu}^p$ estimated by least squares. We then define the total effect size to be the
 953 squared correlation between the observations and the out-of-sample predictions in individuals in I_2 ,

954

$$(R_c^{Tot})^2 = \text{cor}(y_j, \hat{y}_j^{p,c})^2, \quad j \in I_2. \quad (24)$$

955

956 Total effects for the other predictor groups were defined analogously.

957 For groups other than the cell composition group, we also computed a direct effect. For each
958 group, it was computed as the added out-of-sample proportion of variance explained when adding
959 the group predictor term to that of the cell composition group. Take the exposures group as an
960 example, the direct effect was computed by

961

$$(R_e^D)^2 = (R_{e+c}^{\text{Tot}})^2 - (R_c^{\text{Tot}})^2, \quad (25)$$

962

963 where $(R_{e+c}^{\text{Tot}})^2$ is the total effect of the sum of the predictor terms for exposures and cell
964 composition,

965

$$f_{c+e} = f_c^p(\mathbf{x}_i^c) + f_e^p(\mathbf{x}_i^e). \quad (26)$$

966 To mitigate the impact of sampling on estimates of total and direct effects, we did four
967 independent repeats of five-fold cross-validation and averaged effect sizes across all 20 samples. To
968 have an unbiased estimation of the out-of-sample explained variance, we redid a local meQTL
969 mapping on the training set in each iteration of the cross-validation scheme. The algorithm for
970 drawing samples of the total effect is detailed in Algorithm (2).

971

Algorithm (2): Cross-validation for estimating out-of-sample group total effect size

Repeat 4 times:

For $k = 1, \dots, 5$

Index a fifth of individuals as I_k , the others are indexed as $I_{\setminus k}$

Select SNP for the predictor f_g^p by performing a local meQTL mapping on individuals in

$I_{\setminus k}$

For predictor $f_n^p \in \{f_g^p, f_c^p, f_{in}^p, f_e^p\}$

Estimate $\hat{\mu}^p, \hat{\beta}_n^p$ with $I_1 = I_{\setminus k}$

Compute $(R_n^{\text{Tot}})^2$ by Eq. (24) with $I_2 = I_k$

972 The scheme to sample the direct effects is analogous. Finally, we estimated an effect size for
973 interactions between the local SNP and non-genetic factors for each CpG site. It was computed,
974 similarly to Eq. (25), as the added out-of-sample proportion of variance explained by the regression
975 function,

976

$$\begin{aligned}
 f_{\text{Int}}^p(\text{SNP}_i^p, \text{Age}_i, \text{Woman}_i, \text{CMV}_i, \text{ExSmoker}_i, \text{Smoker}_i, \text{CRP}_i) \\
 = \mu^p + \text{SNP}_i^p \beta_{\text{SNP}}^p + \text{Age}_i \beta_{\text{Age}}^p + \text{Woman}_i \beta_{\text{Woman}}^p + \text{CMV}_i \beta_{\text{CMV}}^p \\
 + \text{ExSmoker}_i \beta_{\text{ExSmoker}}^p + \text{Smoker}_i \beta_{\text{Smoker}}^p + \log(\text{CRP}_i) \beta_{\text{CRP}}^p \\
 + \text{SNP}_i^p (\text{Age}_i \theta_{\text{Age}}^p + \text{Woman}_i \theta_{\text{Woman}}^p + \text{CMV}_i \theta_{\text{CMV}}^p \\
 + \text{ExSmoker}_i \theta_{\text{ExSmoker}}^p + \text{Smoker}_i \theta_{\text{Smoker}}^p + \log(\text{CRP}_i) \theta_{\text{CRP}}^p)
 \end{aligned} \tag{27}$$

977
 978 compared to the same regression function without interaction terms,
 979

$$\begin{aligned}
 f_{\text{Main}}^p(\text{SNP}_i^p, \text{Age}_i, \text{Woman}_i, \text{CMV}_i, \text{ExSmoker}_i, \text{Smoker}_i, \text{CRP}_i) \\
 = \mu^p + \text{SNP}_i^p \beta_{\text{SNP}}^p + \text{Age}_i \beta_{\text{Age}}^p + \text{Woman}_i \beta_{\text{Woman}}^p + \text{CMV}_i \beta_{\text{CMV}}^p \\
 + \text{ExSmoker}_i \beta_{\text{ExSmoker}}^p + \text{Smoker}_i \beta_{\text{Smoker}}^p + \log(\text{CRP}_i) \beta_{\text{CRP}}^p.
 \end{aligned} \tag{28}$$

980

981 **Biological annotations**

982 Information about the position, closest gene and CpG density of each CpG site was obtained from
 983 the Illumina EPIC array manifest v.1.0 B4. We retrieved the chromatin state of regions around each
 984 CpG site, using the 15 chromatin states inferred with ChromHMM for CD4⁺ naive T cells by the
 985 ROADMAP Epigenomics consortium¹⁵. We used peripheral blood mononuclear cells (PBMCs) as
 986 reference. The data was downloaded from the consortium webpage
 987 (https://egg2.wustl.edu/roadmap/web_portal/chr_state_learning.html). The transcription factor bind-
 988 ing site data used was public CHIP-seq data collected and processed for the 2020 release of the
 989 ReMap database¹¹⁸, including a total of 1,165 TFs. Binding sites include both direct and indirect
 990 binding. Enrichment analyses were performed by creating simple two-way tables for each target set
 991 and each annotation (i.e., chromatin states, CpG density, transcription factor binding site), and then
 992 performing Fisher's exact test. Gene ontology enrichments were computed with the gometh func-
 993 tion in the missMethyl R package¹¹⁹.

994 We tested if a set of x local or remote meQTL SNPs is enriched in disease- or trait-associated
 995 variants, by sampling at random, among all tested SNPs, 15,000 sets of x SNPs with minor allele
 996 frequencies matched to those of meQTL SNPs. For each resampled set, we calculated the
 997 proportion of variants either known to be associated with a disease or trait, or in LD (set here to $r^2 >$
 998 0.6) with a disease/trait-associated variant (P -value $< 5 \times 10^{-8}$; EBI-NHGRI Catalog of GWAS hits
 999 version e100 r2021-01-1). The enrichment P -value was estimated as the percentage of resamples for

1000 which this proportion was larger than that observed in meQTL SNPs. LD was precomputed for all
1001 5,699,237 SNPs with PLINK 1.9 (with arguments ‘–show-tags all–tag-kb 500–tag-r2 0.6’)¹²⁰.

1002

1003 **Data availability**

1004 The Infinium MethylationEPIC raw and processed data generated in this study²⁷ have been
1005 deposited in the Institut Pasteur data repository, OWEY, which can be accessed via the following
1006 link: <https://dataset.owey.io/doi/10.48802/owey.f83a-1042>. All association statistics obtained in this
1007 study (i.e., the 141 EWAS and interaction models, local meQTL mapping) can be explored and
1008 downloaded from the web browser <http://mimeth.pasteur.fr/>. The SNP array data can be accessed in
1009 the European Genome-Phenome Archive (EGA) with the accession code EGAS00001002460. All
1010 Milieu Intérieur datasets can be accessed by submitting a data access request to
1011 milieuinterieurdac@pasteur.fr, the Milieu Intérieur data access committee, which grants data access
1012 if the request is consistent with the informed consent provided by Milieu Intérieur participants.
1013 Requests are reviewed every month by the committee.

1014

1015 **Code availability**

1016 All the code supporting the current study, including the CMV estimation model, has been uploaded
1017 to GitHub¹²¹: <https://github.com/JacobBergstedt/MIMETH>.

1018 **References**

- 1019 1. Cavalli, G. & Heard, E. Advances in epigenetics link genetics to the environment and disease.
1020 *Nature* **571**, 489-499 (2019).
- 1021 2. Michalak, E.M., Burr, M.L., Bannister, A.J. & Dawson, M.A. The roles of DNA, RNA and
1022 histone methylation in ageing and cancer. *Nat Rev Mol Cell Biol* **20**, 573-589 (2019).
- 1023 3. Martin, E.M. & Fry, R.C. Environmental Influences on the Epigenome: Exposure- Associated
1024 DNA Methylation in Human Populations. *Annu Rev Public Health* **39**, 309-333 (2018).
- 1025 4. Karlsson Linner, R. *et al.* An epigenome-wide association study meta-analysis of educational
1026 attainment. *Mol Psychiatry* **22**, 1680-1690 (2017).
- 1027 5. Lam, L.L. *et al.* Factors underlying variable DNA methylation in a human community cohort.
1028 *Proc Natl Acad Sci U S A* **109 Suppl 2**, 17253-60 (2012).
- 1029 6. Stringhini, S. *et al.* Life-course socioeconomic status and DNA methylation of genes regulating
1030 inflammation. *Int J Epidemiol* **44**, 1320-30 (2015).
- 1031 7. Bush, N.R. *et al.* The biological embedding of early-life socioeconomic status and family
1032 adversity in children's genome-wide DNA methylation. *Epigenomics* **10**, 1445-1461 (2018).
- 1033 8. Hwang, J.Y., Aromolaran, K.A. & Zukin, R.S. The emerging field of epigenetics in
1034 neurodegeneration and neuroprotection. *Nat Rev Neurosci* **18**, 347-361 (2017).
- 1035 9. Mazzone, R. *et al.* The emerging role of epigenetics in human autoimmune disorders. *Clin*
1036 *Epigenetics* **11**, 34 (2019).
- 1037 10. Ling, C. & Ronn, T. Epigenetics in Human Obesity and Type 2 Diabetes. *Cell Metab* **29**, 1028-
1038 1044 (2019).
- 1039 11. van der Harst, P., de Windt, L.J. & Chambers, J.C. Translational Perspective on Epigenetics in
1040 Cardiovascular Disease. *J Am Coll Cardiol* **70**, 590-606 (2017).
- 1041 12. Wild, C.P. Complementing the genome with an "exposome": the outstanding challenge of
1042 environmental exposure measurement in molecular epidemiology. *Cancer Epidemiol*
1043 *Biomarkers Prev* **14**, 1847-50 (2005).
- 1044 13. Berdasco, M. & Esteller, M. Clinical epigenetics: seizing opportunities for translation. *Nat Rev*
1045 *Genet* **20**, 109-127 (2019).
- 1046 14. Farlik, M. *et al.* DNA Methylation Dynamics of Human Hematopoietic Stem Cell
1047 Differentiation. *Cell Stem Cell* **19**, 808-822 (2016).
- 1048 15. Roadmap Epigenomics Consortium *et al.* Integrative analysis of 111 reference human
1049 epigenomes. *Nature* **518**, 317-30 (2015).
- 1050 16. Liu, Y. *et al.* Epigenome-wide association data implicate DNA methylation as an intermediary
1051 of genetic risk in rheumatoid arthritis. *Nat Biotechnol* **31**, 142-7 (2013).

- 1052 17. Lappalainen, T. & Grealley, J.M. Associating cellular epigenetic models with human
1053 phenotypes. *Nat Rev Genet* **18**, 441-451 (2017).
- 1054 18. Houseman, E.A. *et al.* DNA methylation arrays as surrogate measures of cell mixture
1055 distribution. *BMC Bioinformatics* **13**, 86 (2012).
- 1056 19. Teschendorff, A.E., Breeze, C.E., Zheng, S.C. & Beck, S. A comparison of reference-based
1057 algorithms for correcting cell-type heterogeneity in Epigenome-Wide Association Studies.
1058 *BMC Bioinformatics* **18**, 105 (2017).
- 1059 20. Lemire, M. *et al.* Long-range epigenetic regulation is conferred by genetic variation located at
1060 thousands of independent loci. *Nat Commun* **6**, 6326 (2015).
- 1061 21. Bonder, M.J. *et al.* Disease variants alter transcription factor levels and methylation of their
1062 binding sites. *Nat Genet* **49**, 131-138 (2017).
- 1063 22. Villicana, S. & Bell, J.T. Genetic impacts on DNA methylation: research findings and future
1064 perspectives. *Genome Biol* **22**, 127 (2021).
- 1065 23. Min, J.L. *et al.* Genomic and phenotypic insights from an atlas of genetic effects on DNA
1066 methylation. *Nat Genet* **53**, 1311-1321 (2021).
- 1067 24. Thomas, S. *et al.* The Milieu Interieur study - an integrative approach for study of human
1068 immunological variance. *Clin Immunol* **157**, 277-93 (2015).
- 1069 25. Patin, E. *et al.* Natural variation in the parameters of innate immune cells is preferentially
1070 driven by genetic factors. *Nat Immunol* **19**, 302-314 (2018).
- 1071 26. Houseman, E.A., Kelsey, K.T., Wiencke, J.K. & Marsit, C.J. Cell-composition effects in the
1072 analysis of DNA methylation array data: a mathematical perspective. *BMC Bioinformatics* **16**,
1073 95 (2015).
- 1074 27. Bergstedt, J. *et al.* Whole blood DNA methylomes of 958 healthy adults from the Milieu
1075 Intérieur cohort. *OWEY* <https://dataset.owey.io/doi/10.48802/owey.f83a-1042> (2022).
- 1076 28. van den Boogaart, K.G., Filzmoser, P., Hron, K., Templ, M. & Tolosana-Delgado, R. Classical
1077 and Robust Regression Analysis with Compositional Data. *Mathematical Geosciences* **53**, 823–
1078 858 (2021).
- 1079 29. Salas, L.A. *et al.* Enhanced cell deconvolution of peripheral blood using DNA methylation for
1080 high-resolution immune profiling. *Nat Commun* **13**, 761 (2022).
- 1081 30. Jonkman, T.H. *et al.* Functional genomics analysis identifies T and NK cell activation as a
1082 driver of epigenetic clock progression. *Genome Biol* **23**, 24 (2022).
- 1083 31. Rodriguez, R.M. *et al.* Epigenetic Networks Regulate the Transcriptional Program in Memory
1084 and Terminally Differentiated CD8+ T Cells. *J Immunol* **198**, 937-949 (2017).

- 1085 32. Koestler, D.C. *et al.* Improving cell mixture deconvolution by identifying optimal DNA
1086 methylation libraries (IDOL). *BMC Bioinformatics* **17**, 120 (2016).
- 1087 33. Klenerman, P. & Oxenius, A. T cell responses to cytomegalovirus. *Nat Rev Immunol* **16**, 367-
1088 77 (2016).
- 1089 34. VanderWeele, T.J. *Explanation in Causal Inference: Methods for Mediation and Interaction*,
1090 (Oxford University Press, 2015).
- 1091 35. Inoue, T., Iseki, K., Iseki, C. & Kinjo, K. Elevated resting heart rate is associated with white
1092 blood cell count in middle-aged and elderly individuals without apparent cardiovascular
1093 disease. *Angiology* **63**, 541-6 (2012).
- 1094 36. Scheiermann, C., Kunisaki, Y. & Frenette, P.S. Circadian control of the immune system. *Nat*
1095 *Rev Immunol* **13**, 190-8 (2013).
- 1096 37. Cannon, M.J., Schmid, D.S. & Hyde, T.B. Review of cytomegalovirus seroprevalence and
1097 demographic characteristics associated with infection. *Rev Med Virol* **20**, 202-13 (2010).
- 1098 38. Rowles, D.L. *et al.* DNA methyltransferase DNMT3A associates with viral proteins and
1099 impacts HSV-1 infection. *Proteomics* **15**, 1968-82 (2015).
- 1100 39. Zheng, S.C., Breeze, C.E., Beck, S. & Teschendorff, A.E. Identification of differentially
1101 methylated cell types in epigenome-wide association studies. *Nat Methods* **15**, 1059-1066
1102 (2018).
- 1103 40. You, C. *et al.* A cell-type deconvolution meta-analysis of whole blood EWAS reveals lineage-
1104 specific smoking-associated DNA methylation changes. *Nat Commun* **11**, 4779 (2020).
- 1105 41. Morita, K. *et al.* Egr2 and Egr3 in regulatory T cells cooperatively control systemic
1106 autoimmunity through Ltbp3-mediated TGF-beta3 production. *Proc Natl Acad Sci U S A* **113**,
1107 E8131-E8140 (2016).
- 1108 42. Mason, G.M., Poole, E., Sissons, J.G., Wills, M.R. & Sinclair, J.H. Human cytomegalovirus
1109 latency alters the cellular secretome, inducing cluster of differentiation (CD)4+ T-cell
1110 migration and suppression of effector function. *Proc Natl Acad Sci U S A* **109**, 14538-43
1111 (2012).
- 1112 43. Groves, I.J. *et al.* Bromodomain proteins regulate human cytomegalovirus latency and
1113 reactivation allowing epigenetic therapeutic intervention. *Proc Natl Acad Sci U S A* **118**(2021).
- 1114 44. Torti, N., Walton, S.M., Murphy, K.M. & Oxenius, A. Batf3 transcription factor-dependent DC
1115 subsets in murine CMV infection: differential impact on T-cell priming and memory inflation.
1116 *Eur J Immunol* **41**, 2612-8 (2011).
- 1117 45. Savva, G.M. *et al.* Cytomegalovirus infection is associated with increased mortality in the
1118 older population. *Aging Cell* **12**, 381-7 (2013).

- 1119 46. Chen, S. *et al.* Associations of Cytomegalovirus Infection With All-Cause and Cardiovascular
1120 Mortality in Multiple Observational Cohort Studies of Older Adults. *J Infect Dis* **223**, 238-246
1121 (2021).
- 1122 47. Hannum, G. *et al.* Genome-wide methylation profiles reveal quantitative views of human aging
1123 rates. *Mol Cell* **49**, 359-367 (2013).
- 1124 48. Heyn, H. *et al.* Distinct DNA methylomes of newborns and centenarians. *Proc Natl Acad Sci U*
1125 *SA* **109**, 10522-7 (2012).
- 1126 49. Johansson, A., Enroth, S. & Gyllenstein, U. Continuous Aging of the Human DNA Methylome
1127 Throughout the Human Lifespan. *PLoS One* **8**, e67378 (2013).
- 1128 50. Wang, Y. *et al.* Epigenetic influences on aging: a longitudinal genome-wide methylation study
1129 in old Swedish twins. *Epigenetics* **13**, 975-987 (2018).
- 1130 51. Jones, M.J., Goodman, S.J. & Kobor, M.S. DNA methylation and healthy human aging. *Aging*
1131 *Cell* **14**, 924-32 (2015).
- 1132 52. Lachmann, R. *et al.* Cytomegalovirus (CMV) seroprevalence in the adult population of
1133 Germany. *PLoS One* **13**, e0200267 (2018).
- 1134 53. McCartney, D.L. *et al.* An epigenome-wide association study of sex-specific chronological
1135 ageing. *Genome Med* **12**, 1 (2019).
- 1136 54. Zhu, T., Zheng, S.C., Paul, D.S., Horvath, S. & Teschendorff, A.E. Cell and tissue type
1137 independent age-associated DNA methylation changes are not rare but common. *Aging (Albany*
1138 *NY)* **10**, 3541-3557 (2018).
- 1139 55. Bracken, A.P. *et al.* The Polycomb group proteins bind throughout the INK4A-ARF locus and
1140 are disassociated in senescent cells. *Genes Dev* **21**, 525-30 (2007).
- 1141 56. Siebold, A.P. *et al.* Polycomb Repressive Complex 2 and Trithorax modulate *Drosophila*
1142 longevity and stress resistance. *Proc Natl Acad Sci U S A* **107**, 169-74 (2010).
- 1143 57. Boyer, L.A. *et al.* Polycomb complexes repress developmental regulators in murine embryonic
1144 stem cells. *Nature* **441**, 349-53 (2006).
- 1145 58. Dozmorov, M.G. Polycomb repressive complex 2 epigenomic signature defines age-associated
1146 hypermethylation and gene expression changes. *Epigenetics* **10**, 484-95 (2015).
- 1147 59. Fraga, M.F. *et al.* Epigenetic differences arise during the lifetime of monozygotic twins. *Proc*
1148 *Natl Acad Sci U S A* **102**, 10604-9 (2005).
- 1149 60. Sliker, R.C. *et al.* Age-related accrual of methylomic variability is linked to fundamental
1150 ageing mechanisms. *Genome Biol* **17**, 191 (2016).
- 1151 61. Zhang, Q. *et al.* Genotype effects contribute to variation in longitudinal methylome patterns in
1152 older people. *Genome Med* **10**, 75 (2018).

- 1153 62. Nikolich-Zugich, J. The twilight of immunity: emerging concepts in aging of the immune
1154 system. *Nat Immunol* **19**, 10-19 (2018).
- 1155 63. Singmann, P. *et al.* Characterization of whole-genome autosomal differences of DNA
1156 methylation between men and women. *Epigenetics Chromatin* **8**, 43 (2015).
- 1157 64. Yousefi, P. *et al.* Sex differences in DNA methylation assessed by 450 K BeadChip in
1158 newborns. *BMC Genomics* **16**, 911 (2015).
- 1159 65. Gatev, E. *et al.* Autosomal sex-associated co-methylated regions predict biological sex from
1160 DNA methylation. *Nucleic Acids Res* **49**, 9097-9116 (2021).
- 1161 66. Correa-Saez, A. *et al.* Updating dual-specificity tyrosine-phosphorylation-regulated kinase 2
1162 (DYRK2): molecular basis, functions and role in diseases. *Cell Mol Life Sci* (2020).
- 1163 67. Yusipov, I. *et al.* Age-related DNA methylation changes are sex-specific: a comprehensive
1164 assessment. *Aging (Albany NY)* **12**, 24057-24080 (2020).
- 1165 68. Karlsson Linner, R. *et al.* Genome-wide association analyses of risk tolerance and risky
1166 behaviors in over 1 million individuals identify hundreds of loci and shared genetic influences.
1167 *Nat Genet* **51**, 245-257 (2019).
- 1168 69. Lee, J.J. *et al.* Gene discovery and polygenic prediction from a genome-wide association study
1169 of educational attainment in 1.1 million individuals. *Nat Genet* **50**, 1112-1121 (2018).
- 1170 70. Perumal, N., Funke, S., Pfeiffer, N. & Grus, F.H. Proteomics analysis of human tears from
1171 aqueous-deficient and evaporative dry eye patients. *Sci Rep* **6**, 29629 (2016).
- 1172 71. Feinberg, A.P. The Key Role of Epigenetics in Human Disease Prevention and Mitigation. *N*
1173 *Engl J Med* **378**, 1323-1334 (2018).
- 1174 72. O'Geen, H. *et al.* Genome-wide analysis of KAP1 binding suggests autoregulation of KRAB-
1175 ZNFs. *PLoS Genet* **3**, e89 (2007).
- 1176 73. Marchal, C. & Miotto, B. Emerging concept in DNA methylation: role of transcription factors
1177 in shaping DNA methylation patterns. *J Cell Physiol* **230**, 743-51 (2015).
- 1178 74. Quenneville, S. *et al.* The KRAB-ZFP/KAP1 system contributes to the early embryonic
1179 establishment of site-specific DNA methylation patterns maintained during development. *Cell*
1180 *Rep* **2**, 766-73 (2012).
- 1181 75. Hawe, J.S. *et al.* Genetic variation influencing DNA methylation provides insights into
1182 molecular mechanisms regulating genomic function. *Nat Genet* **54**, 18-29 (2022).
- 1183 76. Emilsson, V. *et al.* Co-regulatory networks of human serum proteins link genetics to disease.
1184 *Science* **361**, 769-773 (2018).
- 1185 77. Holler, C.J. *et al.* BACE2 expression increases in human neurodegenerative disease. *Am J*
1186 *Pathol* **180**, 337-50 (2012).

- 1187 78. Granada, M. *et al.* A genome-wide association study of plasma total IgE concentrations in the
1188 Framingham Heart Study. *J Allergy Clin Immunol* **129**, 840-845 e21 (2012).
- 1189 79. Beerman, I. *et al.* Proliferation-dependent alterations of the DNA methylation landscape
1190 underlie hematopoietic stem cell aging. *Cell Stem Cell* **12**, 413-25 (2013).
- 1191 80. Williams, K., Christensen, J. & Helin, K. DNA methylation: TET proteins-guardians of CpG
1192 islands? *EMBO Rep* **13**, 28-35 (2011).
- 1193 81. Li, Y. *et al.* Genome-wide analyses reveal a role of Polycomb in promoting hypomethylation of
1194 DNA methylation valleys. *Genome Biol* **19**, 18 (2018).
- 1195 82. Kim, J.Y., Tavaré, S. & Shibata, D. Counting human somatic cell replications: methylation
1196 mirrors endometrial stem cell divisions. *Proc Natl Acad Sci U S A* **102**, 17739-44 (2005).
- 1197 83. Yang, Z. *et al.* Correlation of an epigenetic mitotic clock with cancer risk. *Genome Biol* **17**, 205
1198 (2016).
- 1199 84. Zhou, W. *et al.* DNA methylation loss in late-replicating domains is linked to mitotic cell
1200 division. *Nat Genet* **50**, 591-602 (2018).
- 1201 85. Feil, R. & Fraga, M.F. Epigenetics and the environment: emerging patterns and implications.
1202 *Nat Rev Genet* **13**, 97-109 (2012).
- 1203 86. Bell, J.T. *et al.* Epigenome-wide scans identify differentially methylated regions for age and
1204 age-related phenotypes in a healthy ageing population. *PLoS Genet* **8**, e1002629 (2012).
- 1205 87. Grundberg, E. *et al.* Mapping cis- and trans-regulatory effects across multiple tissues in twins.
1206 *Nat Genet* **44**, 1084-1089 (2012).
- 1207 88. van Dongen, J. *et al.* Genetic and environmental influences interact with age and sex in shaping
1208 the human methylome. *Nat Commun* **7**, 11115 (2016).
- 1209 89. Teschendorff, A.E. & Relton, C.L. Statistical and integrative system-level analysis of DNA
1210 methylation data. *Nat Rev Genet* **19**, 129-147 (2018).
- 1211 90. Li, S. *et al.* Causes of blood methylomic variation for middle-aged women measured by the
1212 HumanMethylation450 array. *Epigenetics* **12**, 973-981 (2017).
- 1213 91. Boyce, W.T. & Kobor, M.S. Development and the epigenome: the 'synapse' of gene-
1214 environment interplay. *Dev Sci* **18**, 1-23 (2015).
- 1215 92. Fleiss, J.L. *Design and analysis of clinical experiments*, (New York: Wiley, 2011).
- 1216 93. Czamara, D. *et al.* Integrated analysis of environmental and genetic influences on cord blood
1217 DNA methylation in new-borns. *Nat Commun* **10**, 2548 (2019).
- 1218 94. Teh, A.L. *et al.* The effect of genotype and in utero environment on interindividual variation in
1219 neonate DNA methylomes. *Genome Res* **24**, 1064-74 (2014).

- 1220 95. Bergstedt, J., Urrutia, A., Albert, M.L., Quintana-Murci, L. & Patin, E. Accurate prediction of
1221 cell composition, age, smoking consumption and infection serostatus based on blood DNA
1222 methylation profiles. *bioRxiv* (2018).
- 1223 96. Furman, D. *et al.* Chronic inflammation in the etiology of disease across the life span. *Nat Med*
1224 **25**, 1822-1832 (2019).
- 1225 97. Mauvais-Jarvis, F. *et al.* Sex and gender: modifiers of health, disease, and medicine. *Lancet*
1226 **396**, 565-582 (2020).
- 1227 98. Niccoli, T. & Partridge, L. Ageing as a risk factor for disease. *Curr Biol* **22**, R741-52 (2012).
- 1228 99. Samet, J.M. Tobacco smoking: the leading cause of preventable disease worldwide. *Thorac*
1229 *Surg Clin* **23**, 103-12 (2013).
- 1230 100. Fortin, J.P., Triche, T.J., Jr. & Hansen, K.D. Preprocessing, normalization and integration of the
1231 Illumina HumanMethylationEPIC array with minfi. *Bioinformatics* **33**, 558-560 (2017).
- 1232 101. Triche, T.J., Jr., Weisenberger, D.J., Van Den Berg, D., Laird, P.W. & Siegmund, K.D. Low-
1233 level processing of Illumina Infinium DNA Methylation BeadArrays. *Nucleic Acids Res* **41**,
1234 e90 (2013).
- 1235 102. Bates, D., Mächler, M., Bolker, B. & Walker, S. Fitting Linear Mixed-Effects Models Using
1236 lme4. *J Stat Soft* **67**, 1–48 (2015).
- 1237 103. Johnson, W.E., Li, C. & Rabinovic, A. Adjusting batch effects in microarray expression data
1238 using empirical Bayes methods. *Biostatistics* **8**, 118-27 (2007).
- 1239 104. Price, M.E. *et al.* Additional annotation enhances potential for biologically-relevant analysis of
1240 the Illumina Infinium HumanMethylation450 BeadChip array. *Epigenetics Chromatin* **6**, 4
1241 (2013).
- 1242 105. Auton, A. *et al.* A global reference for human genetic variation. *Nature* **526**, 68-74 (2015).
- 1243 106. Stekhoven, D.J. & Bühlmann, P. MissForest--non-parametric missing value imputation for
1244 mixed-type data. *Bioinformatics* **28**, 112-8 (2012).
- 1245 107. Manichaikul, A. *et al.* Robust relationship inference in genome-wide association studies.
1246 *Bioinformatics* **26**, 2867-73 (2010).
- 1247 108. Pawlowsky-Glahn, V., José Egozcue, J. & Tolosana-Delgado, R. *Modeling and analysis of*
1248 *compositional data*, (Wiley, 2015).
- 1249 109. Arnold, K.F., Berrie, L., Tennant, P.W.G. & Gilthorpe, M.S. A causal inference perspective on
1250 the analysis of compositional data. *International Journal of Epidemiology* **49**, 1307-1313
1251 (2020).
- 1252 110. Templ, M., Hron, K. & Filzmoser, P. *robCompositions: an R-package for robust statistical*
1253 *analysis of compositional data*, (John Wiley and Sons, 2011).

- 1254 111. Halekoh, U. & Højsgaard, S. A Kenward-Roger Approximation and Parametric Bootstrap
1255 Methods for Tests in Linear Mixed Models – The R Package pbrtest. *2014* **59**, 32 (2014).
- 1256 112. Salas, L.A. *et al.* An optimized library for reference-based deconvolution of whole-blood
1257 biospecimens assayed using the Illumina HumanMethylationEPIC BeadArray. *Genome Biol*
1258 **19**, 64 (2018).
- 1259 113. Zou, H. & Hastie, T. Regularization and variable selection via the elastic net. *Journal of the*
1260 *Royal Statistical Society: Series B (Statistical Methodology)* **67**, 301-320 (2005).
- 1261 114. Friedman, J., Hastie, T. & Tibshirani, R. Regularization Paths for Generalized Linear Models
1262 via Coordinate Descent. *J Stat Soft* **33**, 1–22 (2010).
- 1263 115. Stasinopoulos, M.D., Rigby, R.A. & Bastiani, F.D. GAMLSS: A distributional regression
1264 approach. *Statistical Modelling* **18**, 248-273 (2018).
- 1265 116. Shabalin, A.A. Matrix eQTL: ultra fast eQTL analysis via large matrix operations.
1266 *Bioinformatics* **28**, 1353-8 (2012).
- 1267 117. Zeileis, A., Köll, S. & Graham, N. Various Versatile Variances: An Object-Oriented
1268 Implementation of Clustered Covariances in R. *2020* **95**, 36 (2020).
- 1269 118. Cheneby, J. *et al.* ReMap 2020: a database of regulatory regions from an integrative analysis of
1270 Human and Arabidopsis DNA-binding sequencing experiments. *Nucleic Acids Res* **48**, D180-
1271 D188 (2020).
- 1272 119. Phipson, B., Maksimovic, J. & Oshlack, A. missMethyl: an R package for analyzing data from
1273 Illumina's HumanMethylation450 platform. *Bioinformatics* **32**, 286-8 (2016).
- 1274 120. Chang, C.C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer
1275 datasets. *GigaScience* **4**, 7 (2015).
- 1276 121. Bergstedt, J. The Immune Factors Driving DNA Methylation Variation in Human Blood.
1277 *GitHub* <https://doi.org/10.5281/zenodo.7016878> (2022).
- 1278
- 1279
- 1280

1281 **Acknowledgments**

1282 We thank Sarah Merrill, Nicole Gladish, Violaine Saint-André, Lucas Husquin and the Milieu In-
1283 térieur scientific advisory board for comments and helpful discussions. We acknowledge the help of
1284 the HPC Core Facility of Institut Pasteur for this work. This research was enabled, in part, by the
1285 use of the FlowSorted.BloodExtended.EPIC R package developed at Dartmouth College, which
1286 software is subject to the licensing terms made available by Dartmouth Technology Transfer and
1287 which software is provided "as is" with no warranties whatsoever. This work benefited from support
1288 of the French government's program 'Investissement d'Avenir', managed by the Agence Nationale
1289 de la Recherche (reference 10-LABX-69-01).

1290

1291 **Author contributions**

1292 L.Q.-M. initiated the study. J.B., E.P. and L.Q.-M. conceived and developed the study. A.U. pre-
1293 pared DNA samples. D.T.S.L., J.L.M. and M.S.K. acquired Illumina Infinium MethylationEPIC ar-
1294 ray data. J.B. performed all analyses, with contributions from S.A.K.A., K.T. and E.P.. E.P. super-
1295 vised all analyses. A.J. developed the web browser. D.D. and M.L.A. advised on experiments. M.R.,
1296 M.S.K., D.D. and M.L.A. advised on data interpretation. J.B., E.P. and L.Q.-M. wrote the manu-
1297 script. All authors discussed the results and contributed to the final manuscript.

1298

1299 **Competing interests**

1300 The authors declare no competing interests.

1301

1302 **Milieu Intérieur Consortium members**

1303 Laurent Abel⁹, Andres Alcover¹⁰, Hugues Aschard¹¹, Philippe Bousso¹², Nollaig Bourke¹³, Petter
1304 Brodin^{14,15}, Pierre Bruhns¹⁶, Nadine Cerf-Bensussan¹⁷, Ana Cumano¹⁸, Christophe d'Enfert¹²,
1305 Ludovic Deriano¹⁹, Marie-Agnès Dillies²⁰, James Di Santo²¹, Françoise Dromer²², Gérard Eberl²³,
1306 Jost Enninga²⁴, Jacques Fellay^{25,26,27}, Ivo Gomperts-Boneca²⁸, Milena Hasan²⁹, Gunilla Karlsson
1307 Hedestam³⁰, Serge Hercberg³¹, Molly A. Ingersoll³², Olivier Lantz^{33,34}, Rose Anne Kenny^{35,36},
1308 Mickaël Ménager³⁷, Frédérique Michel³⁸, Hugo Mouquet³⁹, Cliona O'Farrelly^{40,41}, Etienne Patin¹,
1309 Sandra Pellegrini³⁸, Antonio Rausell⁴², Frédéric Rieux-Laucat⁴³, Lars Rogge⁴⁴, Magnus Fontes⁴⁵,
1310 Anavaj Sakuntabhai⁴⁶, Olivier Schwartz⁴⁷, Benno Schwikowski⁴⁸, Spencer Shorte⁴⁹, Frédéric
1311 Tangy⁵⁰, Antoine Toubert⁵¹, Mathilde Touvier³¹, Marie-Noëlle Ungeheuer⁵², Christophe Zimmer⁵³,
1312 Matthew L. Albert^{4,§}, Darragh Duffy^{6,§}, Lluís Quintana-Murci^{1,7§}

1313

- 1314 ¹Institut Pasteur, Université Paris Cité, CNRS UMR2000, Human Evolutionary Genetics Unit,
1315 Paris, France
- 1316 ⁹Imagine Institute, University Paris Cité, Necker Hospital for Sick Children, INSERM UMR 1163,
1317 Laboratory of Human Genetics of Infectious Diseases, Paris, France.
- 1318 ¹⁰Institut Pasteur, Université de Paris, INSERM-U1224, Unité Biologie Cellulaire des
1319 Lymphocytes, Ligue Nationale Contre le Cancer, Équipe Labellisée Ligue 2018, Paris, France
- 1320 ¹¹Institut Pasteur, Université de Paris, Department of Computational Biology, Paris, France
- 1321 ¹²Institut Pasteur, Université de Paris, INRAE USC2019, Unité Biologie et Pathogénicité
1322 Fongiques, Paris, France
- 1323 ¹³Department of Medical Gerontology, School of Medicine, Trinity College Dublin, Dublin, Ireland
- 1324 ¹⁴Department of Immunology and Inflammation, Imperial College London, London, UK
- 1325 ¹⁵Department of Women's and Children's Health, Karolinska Institutet, Stockholm, Sweden
- 1326 ¹⁶Institut Pasteur, Université Paris Cité, INSERM UMR1222, Unit of Antibodies in Therapy and
1327 Pathology, Paris, France
- 1328 ¹⁷Institut Imagine, Université Paris Cité, INSERM UMR1163, Laboratory Intestinal Immunity,
1329 Paris, France
- 1330 ¹⁸Institut Pasteur, Université de Paris, INSERM U1223, Unit Lymphocytes and Immunity, Paris,
1331 France
- 1332 ¹⁹Institut Pasteur, Université de Paris, INSERM U1223, Équipe Labellisée Ligue Contre Le Cancer,
1333 Genome Integrity, Immunity and Cancer Unit, Paris, France
- 1334 ²⁰Institut Pasteur, Université de Paris, Bioinformatics and Biostatistics Hub, Paris, France
- 1335 ²¹Institut Pasteur, Université de Paris, INSERM U1223, Innate Immunity Unit, Paris, France
- 1336 ²²Institut Pasteur, Université de Paris, CNRS UMR2000, Unité de Mycologie Moléculaire, Centre
1337 national de Référence Mycoses Invasives et Antifongiques, Paris, France
- 1338 ²³Institut Pasteur, Université de Paris, INSERM U1224, Microenvironment and Immunity Unit,
1339 Paris, France
- 1340 ²⁴Institut Pasteur, Université de Paris, CNRS UMR3691, Dynamics of Host-Pathogen Interactions
1341 Unit, Paris, France
- 1342 ²⁵School of Life Sciences, École Polytechnique Fédérale de Lausanne, Lausanne, Switzerland
- 1343 ²⁶Swiss Institute of Bioinformatics, Lausanne, Switzerland
- 1344 ²⁷Precision Medicine Unit, Lausanne University Hospital and University of Lausanne, Lausanne,
1345 Switzerland
- 1346 ²⁸Institut Pasteur, Université de Paris, CNRS UMR2001, Unité Biologie et Génétique de la Paroi
1347 Bactérienne, Paris, France

- 1348 ²⁹Institut Pasteur, Université de Paris, Cytometry and Biomarkers Unit of Technology and Service,
1349 Paris, France
- 1350 ³⁰Department of Microbiology, Tumor and Cell Biology, Karolinska Institutet, Stockholm, Sweden
- 1351 ³¹Université Sorbonne Paris Nord, Université de Paris, INSERM U1153, INRAE U1125, CNAM,
1352 Epidemiology and Statistics Research Center, Nutritional Epidemiology Research Team, Bobigny,
1353 France
- 1354 ³²Institut Pasteur, Université de Paris, Institut Cochin, INSERM U1016, CNRS UMR 8104,
1355 Mucosal Inflammation and Immunity Group, Paris, France
- 1356 ³³Institut Curie, Université Paris Sciences et Lettres, INSERM U932, Laboratoire d'Immunologie
1357 Clinique, Paris, France
- 1358 ³⁴Centre d'Investigation Clinique en Biothérapie Gustave-Roussy Institut Curie (CIC-BT1428),
1359 Paris, France
- 1360 ³⁵The Irish Longitudinal Study on Ageing (TILDA), Trinity College Dublin, Dublin, Ireland
- 1361 ³⁶Mercer's Institute for Successful Ageing, St James's Hospital, Dublin, Ireland
- 1362 ³⁷Imagine Institute, Université de Paris, INSERM UMR1163, Laboratory of Inflammatory
1363 Responses and Transcriptomic Networks in Diseases, Atip-Avenir Team, Paris, France
- 1364 ³⁸Institut Pasteur, Université de Paris, INSERM U1221, Cytokine Signaling Unit, Paris, France.
- 1365 ³⁹Institut Pasteur, Université de Paris, INSERM U1222, Laboratory of Humoral Immunology,
1366 75015, Paris, France
- 1367 ⁴⁰Comparative Immunology, School of Biochemistry and Immunology, Trinity Biomedical Sciences
1368 Institute, Dublin, Ireland
- 1369 ⁴¹School of Medicine, Trinity College Dublin, Dublin, Ireland
- 1370 ⁴²Imagine Institute, Université de Paris, INSERM UMR1163, Necker Hospital for Sick Children,
1371 Clinical Bioinformatics Laboratory, Paris, France
- 1372 ⁴³Imagine Institute, Université de Paris, INSERM UMR 1163, Laboratory of Immunogenetics of
1373 Autoimmune Diseases in Children, Paris, France
- 1374 ⁴⁴Institut Pasteur, Université de Paris, AP-HP Hôpital Cochin, Immunoregulation Unit, Paris, France
- 1375 ⁴⁵Institut Roche, Boulogne-Billancourt, France
- 1376 ⁴⁶Institut Pasteur, Université de Paris, CNRS UMR2000, Functional Genetics of Infectious Diseases
1377 Unit, Paris, France
- 1378 ⁴⁷Institut Pasteur, Université de Paris, CNRS UMR3569, Virus and Immunity Unit, Paris, France.
- 1379 ⁴⁸Institut Pasteur, Université de Paris, Computational Systems Biomedicine Lab, Paris, France
- 1380 ⁴⁹Institut Pasteur, Université de Paris, UTechS-PBI/Imagopole, Paris, France

1381 ⁵⁰Institut Pasteur, Université de Paris, CNRS UMR3965, Viral Genomics and Vaccination Unit,
1382 Paris, France

1383 ⁵¹AP-HP, Hôpital Saint-Louis, Université de Paris, INSERM UMR1160, Laboratoire
1384 d'Immunologie et d'Histocompatibilité, Paris, France

1385 ⁵²Institut Pasteur, Université de Paris, Investigation Clinique et Accès aux Ressources Biologiques
1386 (ICAReB), Paris, France

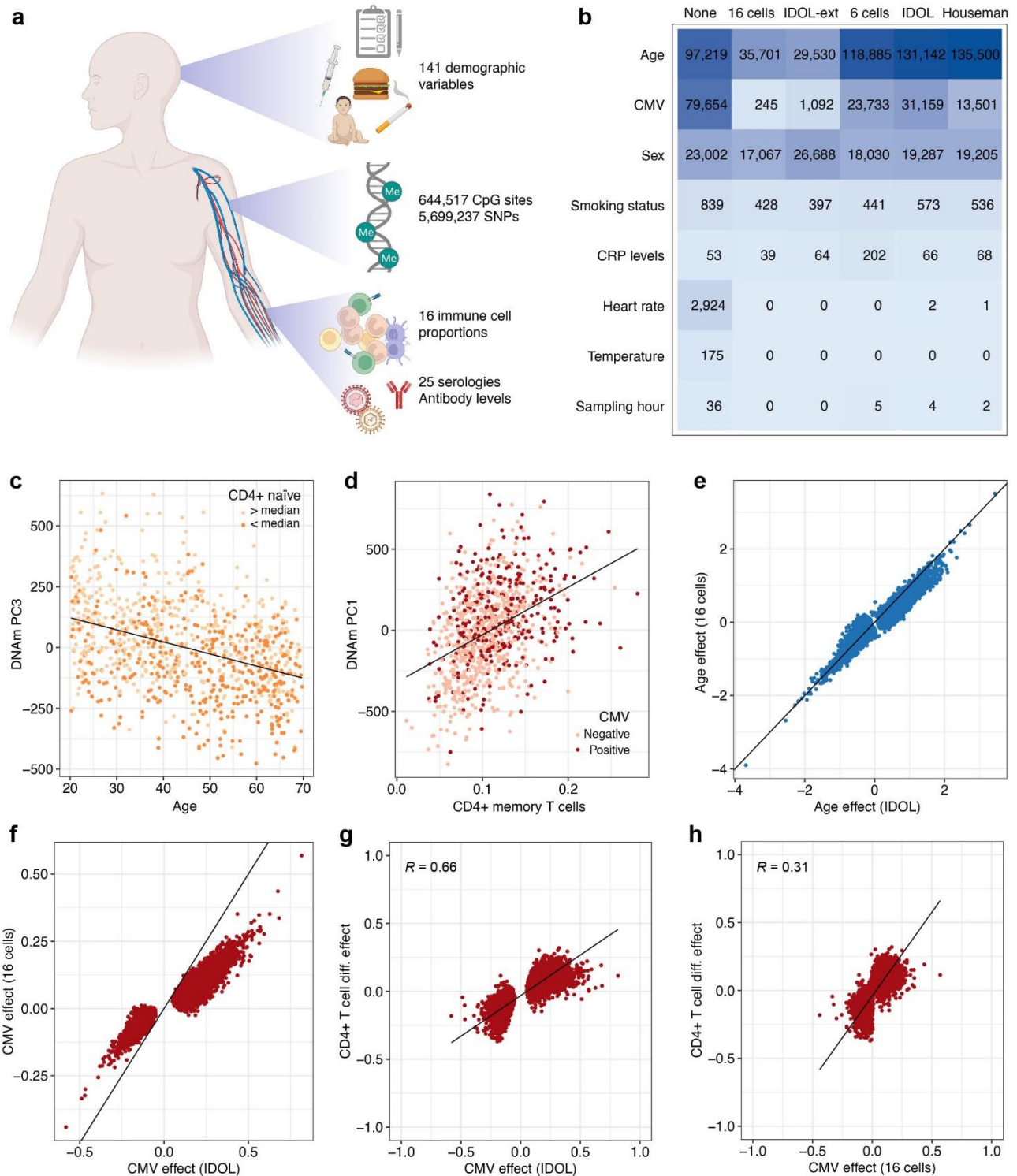
1387 ⁵³Institut Pasteur, Université de Paris, CNRS UMR3691, Imaging and Modeling Unit, Paris, France

1388 §Co-coordinators of the Milieu Intérieur Consortium (<https://milieuinterieur.fr/en/>)

1389

1390

1391 **Figure legends**



1392

1393

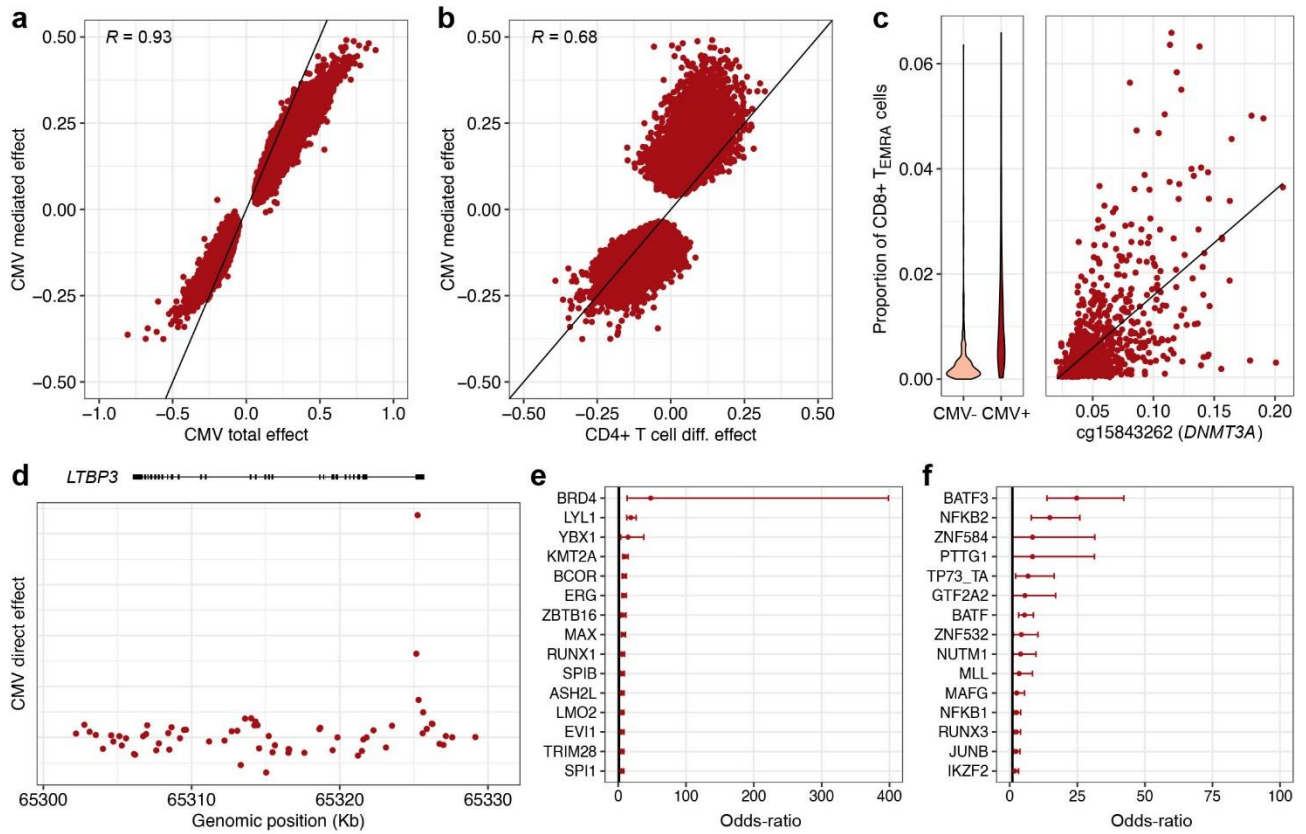
1394 **Fig. 1. Non-genetic effects on the blood DNA methylome according to different corrections for**

1395 **cellular heterogeneity. a** Study design. **b** Number of CpG sites associated with non-genetic factors,

1396 according to different corrections for cellular heterogeneity. Columns indicate adjustments for 16

1397 blood cell proportions measured by flow cytometry (“16 cells”), 12 blood cell proportions estimated

1398 by the EPIC IDOL-Ext deconvolution method²⁹ (“IDOL-ext”), 6 blood cell proportions measured
1399 by flow cytometry (“6 cells”), 6 cell proportions estimated by the IDOL deconvolution method³²
1400 (“IDOL”), 6 cell proportions estimated by Houseman *et al.*’s deconvolution method¹⁸ (“House-
1401 man”) and no adjustment for blood cell composition (“None”). **c** Age against the third Principal
1402 Component (PC) of DNA methylation levels. Colors indicate donors whose proportion of naïve
1403 CD8⁺ T cells in blood is below or above the cohort median. **d** Proportion of CD4⁺ memory T cells
1404 against the first PC of DNA methylation levels. Colors indicate the CMV serostatus of donors. **e** Di-
1405 rect effects of age on 5mC levels, adjusting on 6 cell proportions estimated by IDOL, against direct
1406 effects of age on 5mC levels, adjusting on 16 cell proportions measured by flow cytometry. **f** Direct
1407 effects of CMV serostatus on 5mC levels, adjusting on 6 cell proportions estimated by IDOL,
1408 against direct effects of CMV serostatus on 5mC levels, adjusting on 16 cell proportions measured
1409 by flow cytometry. **g** Effects of CD4⁺ T cell differentiation on 5mC levels against direct effects of
1410 CMV serostatus on 5mC levels, adjusting on 6 cell proportions estimated by IDOL. **h** Effects of
1411 CD4⁺ T cell differentiation on 5mC levels against direct effects of CMV serostatus on 5mC levels,
1412 adjusting on 16 cell proportions measured by flow cytometry. **e-h** Effect sizes are given in the M
1413 value scale. Only associations significant either with the model adjusting for IDOL-estimated cell
1414 proportions or the model adjusted for 16 measured cell proportions are shown ($P_{adj} < 0.05$). **e-f** The
1415 black line indicates the identity line. **c-d, g-h** The black line indicates the linear regression line. Sta-
1416 tistics were computed based on a sample size of $n = 884$ and for 644,517 CpG sites.
1417



1418

1419

1420

1421

1422

1423

1424

1425

1426

1427

1428

1429

1430

1431

1432

1433

1434

Fig. 2. Effects of cytomegalovirus infection on the blood DNA methylome. a Total effects

against cell-composition-mediated effects of CMV infection on 5mC levels. **b** Effects of CD4⁺ T cell differentiation on 5mC levels against cell-composition-mediated effects of CMV infection on

5mC levels. **c** Proportion of CD8⁺ T_{EMRA} cells in CMV⁻ and CMV⁺ donors (left panel). 5mC levels

at the *DNMT3A* locus against the proportion of CD8⁺ T_{EMRA} cells (right panel). 5mC levels are

given in the β value scale. The black line indicates the linear regression line. **d** Genomic distribution

of direct effects of CMV infection at the *LTBP3* locus. **e** Enrichment of CpG sites with a significant

direct, positive effect of CMV infection in binding sites for TFs. **f** Enrichment of CpG sites with a

significant direct, negative effect of CMV infection in binding sites for TFs. **a, b** Only CpG sites

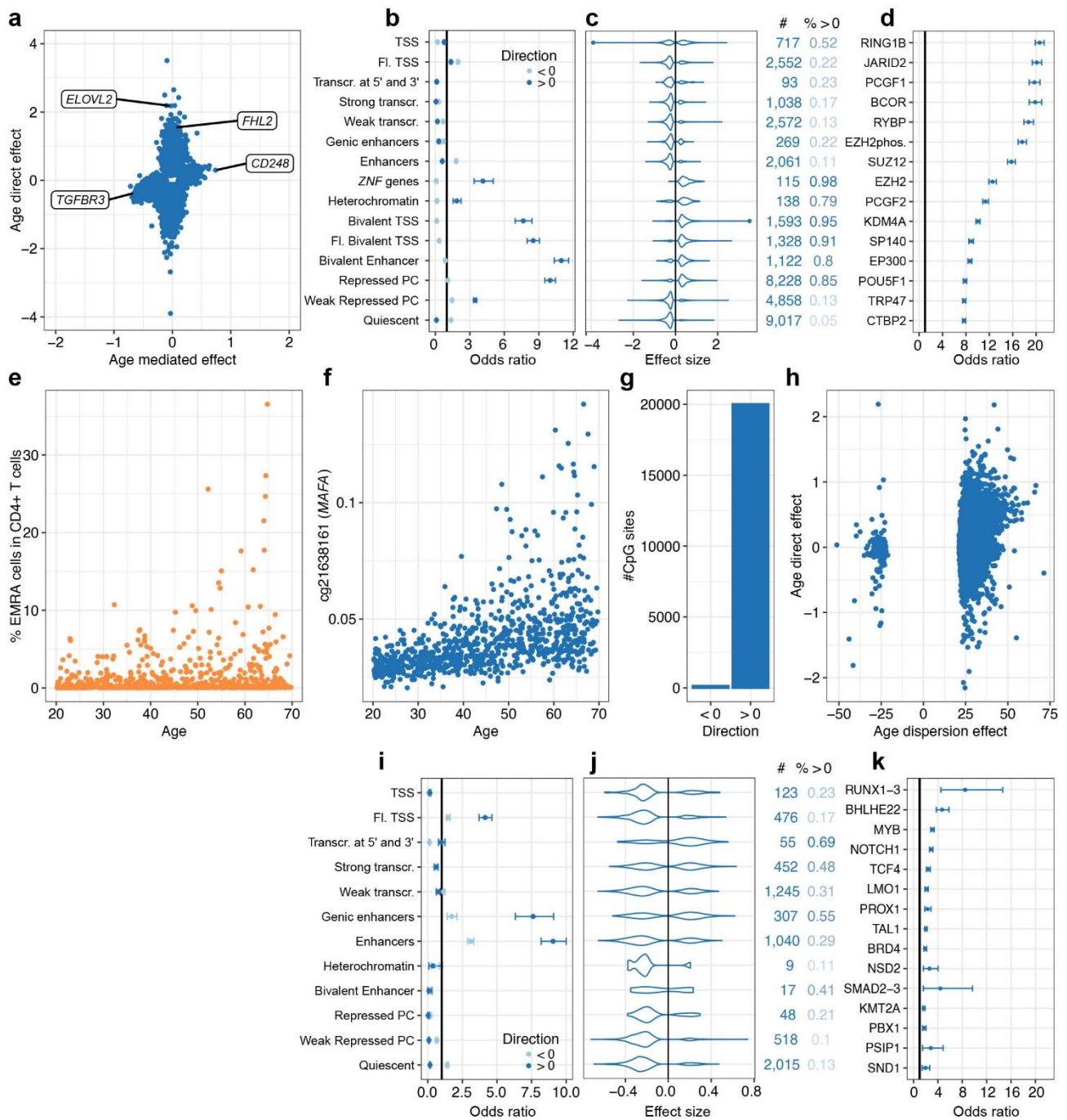
with a significant cell-composition-mediated effect are shown. The black line indicates the identity

line. **a, b, d** Effect sizes are given in the M value scale. **e, f** The 15 most enriched TFs are shown,

out of 1,165 tested TFs. The point and error bars indicate the odds-ratio and 95% CI. CIs were esti-

mated by the Fisher's exact method. Statistics were computed based on a sample size of $n = 884$

and for 644,517 CpG sites.



1435

1436

1437

1438

1439

1440

1441

1442

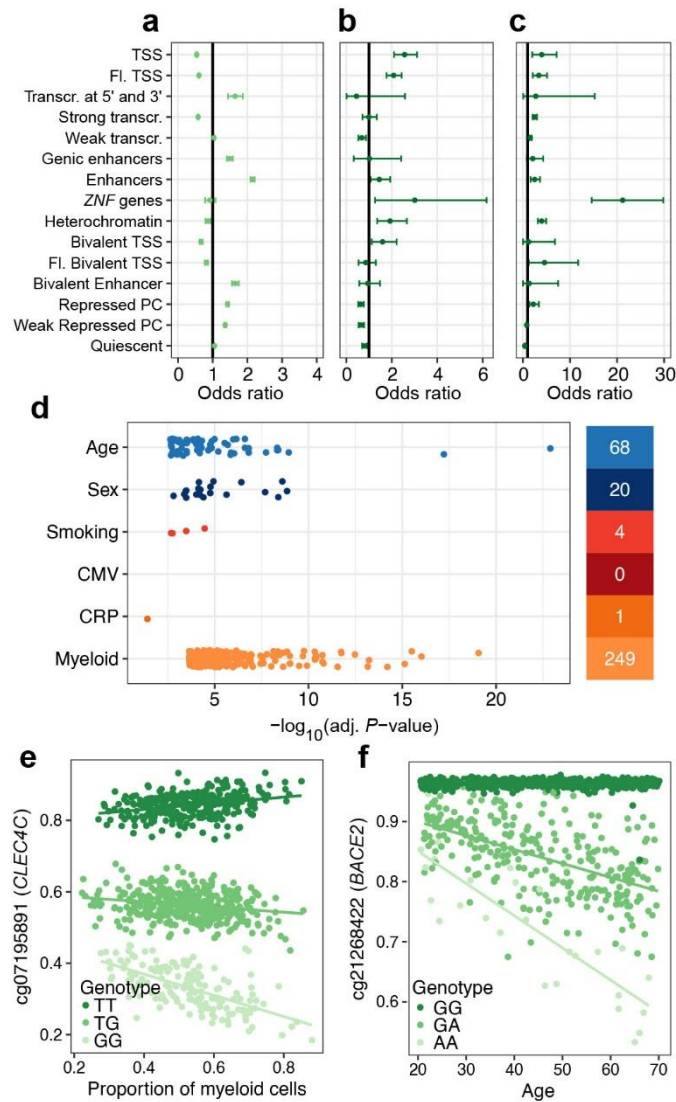
1443

1444

Fig. 3. Direct and cell-composition-mediated effects of aging on the blood DNA methylome. a Direct effects against cell-composition-mediated effects of age on 5mC levels (50-year effect). Only CpG sites with a significant direct or cell-composition-mediated effect are shown. Labels denote genes with strong direct or cell-composition-mediated effects of age. **b** Enrichment in CpG sites with significant direct effects of age, across 15 chromatin states. **c** Distributions of significant direct effects of age, across 15 chromatin states. **d** Enrichment of CpG sites with a significant positive, direct effect of age in binding sites for TFs. **e** Increased variance of the proportion of CD4⁺ TEMRA cells with age. **f** Increased variance of 5mC levels with age at the *MAFA* locus. 5mC levels are given

1445 in the β value scale. **g** Number of CpG sites with a significant increase or decrease in variance with
1446 age. **h** Direct effects against dispersion effects of age on 5mC levels. **i** Enrichment of CpG sites with
1447 significant cell-composition-mediated effects of age, across 12 chromatin states. **j** Distributions of
1448 significant cell-composition-mediated effects of age, across 12 chromatin states. **k** Enrichment of
1449 CpG sites with significant cell-composition-mediated, positive effects of age in binding sites for
1450 TFs. **a,c,h,j** Effect sizes are given in the M value scale. **c,j** Numbers on the right indicate the num-
1451 ber of associated CpG sites and proportion of positive effects. **b,d,i,k** The point and error bars indi-
1452 cate the odds-ratio and 95% CI. CIs were estimated by the Fisher's exact method. Statistics were
1453 computed based on a sample size of $n = 884$ and for 644,517 CpG sites. **d,k** The 15 most enriched
1454 TFs are shown, out of 1,165 tested TFs. **b,c,i,j** Chromatin states were defined in PBMCs¹⁵. Chroma-
1455 tin states were not shown when < 5 associated CpG sites were observed. TSS, Fl. and PC denote
1456 transcription start site, flanking and Polycomb, respectively. Statistics were computed based on a
1457 sample size of $n = 884$ and for 644,517 CpG sites.

1458



1459

1460

1461 **Fig. 4. Effects of genetics and gene \times environment interactions on the blood DNA methylome.**

1462 **a** Enrichment in CpG sites associated with local meQTL variants, across 15 chromatin states. **b** En-

1463 richment in CpG sites associated with remote meQTL variants, across 15 chromatin states. **c** En-

1464 richment in remote meQTL variants, across 15 chromatin states. **d** P -value distributions for signifi-

1465 cant effects of genotype \times age, genotype \times sex, genotype \times smoking, genotype \times CMV serostatus,

1466 genotype \times CRP levels and genotype \times cell-type interactions. The number of significant associa-

1467 tions is indicated on the right. Associations were tested by two-sided Wald tests with heteroscedas-

1468 ticity-consistent standard errors estimated by the sandwich R package¹¹⁷. Multiple testing was done

1469 by the Bonferroni correction separately for each term. **e** Myeloid lineage-dependent effect of the

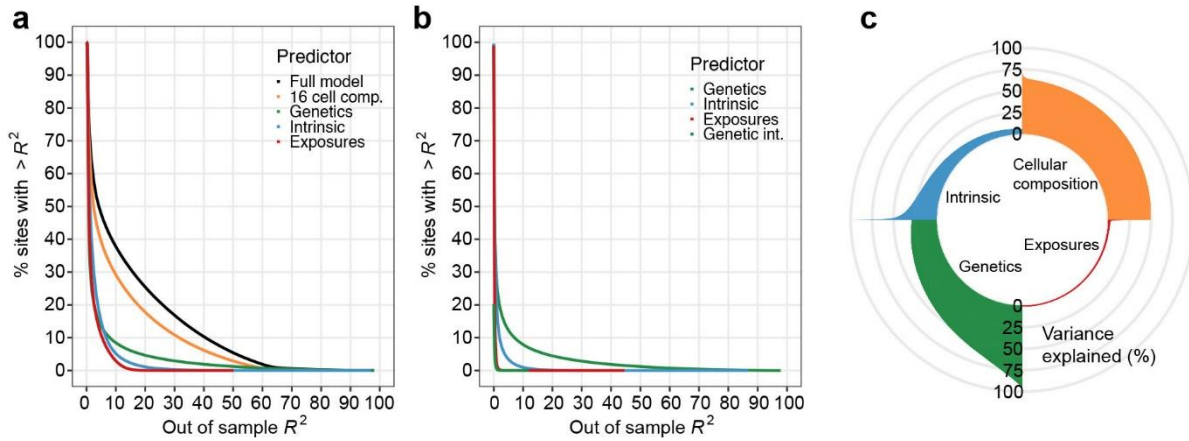
1470 rs11055602 variant on 5mC levels at the *CLEC4C* locus. **f** Age-dependent effect of the rs2837990

1471 variant on 5mC levels at the *BACE2* locus. **a-c** The point and error bars indicate the odds-ratio and

1472 95% CI. CIs were estimated by the Fisher's exact method. Chromatin states were defined in

1473 PBMCs¹⁵. TSS, Fl. and PC denote transcription start site, flanking and Polycomb, respectively. **e-f**
1474 5mC levels are given in the β value scale. Solid lines indicate linear regression lines. Statistics were
1475 computed based on a sample size of $n = 884$ and for 644,517 CpG sites.
1476

1477



1478

1479

1480

1481

1482

1483

1484

1485

1486

1487

1488

1489

Fig. 5. Best predictors of the blood DNA methylome of adults. **a** Complementary cumulative distribution function of the out-of-sample variance explained by the full model, blood cell composition, genetic factors, intrinsic factors (i.e., age and sex) and environmental exposures (i.e., smoking, CMV infection and CRP levels), for 644,517 CpG sites. **b** Complementary cumulative distribution function of the out-of-sample variance explained by genetic factors, intrinsic factors, environmental exposures and gene \times environment ($G \times E$) interactions, when conditioning on blood cell composition, for 644,517 CpG sites. **c** Proportion of the explained out-of-sample variance of 5mC levels for the 20,000 CpG sites with the variance most explained by blood cell composition, genetic factors, intrinsic factors and environmental exposures, respectively.