1  **A flexible, efficient binomial mixed model for identifying differential DNA**

2  **methylation in bisulfite sequencing data**

3  AJ Lea[1], J Tung[1,2,3,4*†], X Zhou[5,6*†]

4  *These authors contributed equally to this work.

5

6  1. Department of Biology, Duke University, Box 90338, Durham, NC 27708, USA

7  2. Institute of Primate Research, National Museums of Kenya, P. O. Box 24481,

8  Karen 00502, Nairobi, Kenya

9  3. Department of Evolutionary Anthropology, Duke University, Box 90383, Durham,

10  NC 27708, USA

11  4. Duke University Population Research Institute, Duke University, Box 90420,

12  Durham, NC 27708, USA

13  5. Department of Biostatistics, University of Michigan, 1415 Washington Heights,

14  Ann Arbor, MI 48109

15  6. Center for Statistical Genetics, University of Michigan, 1415 Washington Heights,

16  Ann Arbor, MI 48109

17

18  [†]Corresponding author email: xzhousph@umich.edu, jt5@duke.edu

19  Other author emails: amanda.lea@duke.edu

20

21

## Abstract

Identifying sources of variation in DNA methylation levels is important for understanding gene regulation. Recently, bisulfite sequencing has become a popular tool for investigating DNA methylation levels. However, modeling bisulfite sequencing data is complicated by dramatic variation in coverage across sites and individual samples, and because of the computational challenges of controlling for genetic covariance in count data. To address these challenges, we present a binomial mixed model and an efficient, sampling-based algorithm (MACAU: Mixed model association for count data via data augmentation) for approximate parameter estimation and $p$-value computation. This framework allows us to simultaneously account for both the over-dispersed, count-based nature of bisulfite sequencing data, as well as genetic relatedness among individuals. Using simulations and two real data sets (whole genome bisulfite sequencing (WGBS) data from *Arabidopsis thaliana* and reduced representation bisulfite sequencing (RRBS) data from baboons), we show that our method provides well-calibrated test statistics in the presence of population structure. Further, it improves power to detect differentially methylated sites: in the RRBS data set, MACAU detected 1.6-fold more age-associated CpG sites than a beta-binomial model (the next best approach). Changes in these sites are consistent with known age-related shifts in DNA methylation levels, and are enriched near genes that are differentially expressed with age in the same population. Taken together, our results indicate that MACAU is an efficient, effective tool for analyzing bisulfite sequencing data, with particular salience to analyses of structured populations. MACAU is freely available at www.xzlab.org/software.html.

## Author Summary

45

46      DNA methylation is an important epigenetic modification involved in regulating

47   gene expression. It can be measured at base-pair resolution, on a genome-wide

48   scale, by coupling sodium bisulfite conversion with high-throughput sequencing (a

49   technique known as 'bisulfite sequencing'). However, the data generated by such

50   methods present several challenges for statistical analysis. In particular, while the

51   raw data generated from bisulfite sequencing experiments are read counts, they are

52   often converted to proportions for ease of modeling, resulting in loss of information.

53   Furthermore, although DNA methylation levels are known to be heritable—and are

54   thus affected by kinship and population structure—existing approaches for modeling

55   bisulfite sequencing data fail to account for this covariance. Such failure can lead to

56   spurious associations and reduced power. Here, we present a new approach that

57   models bisulfite sequencing data using raw read counts, while also taking into

58   account population structure and other sources of data over-dispersion. Using

59   simulations and two real data sets (publicly available data from *Arabidopsis thaliana*

60   and newly generated data from *Papio cynocephalus*), we demonstrate that our

61   model provides well-calibrated *p*-values and improves power compared with

62   previous methods. In addition, the DNA methylation patterns identified by our

63   method agree with those reported in previous studies.

64

## Introduction

66    DNA methylation — the covalent addition of methyl groups to cytosine bases

67    — is a major epigenetic gene regulatory mechanism observed in a wide variety of

68    species. DNA methylation influences genome-wide gene expression patterns, is

69    involved in genomic imprinting and X-inactivation, and functions to suppress the

70    activity of transposable elements [1–3]. In addition, DNA methylation is essential for

71    normal development. For example, mutant *Arabidopsis* plants with reduced levels of

72    DNA methylation display a range of abnormalities including reduced overall size,

73    altered leaf size and shape, and reduced fertility [4–6]. In humans, DNA methylation

74    levels are strongly linked to disease, including major public health burdens such as

75    diabetes [7,8], Alzheimer's disease [9,10], and many forms of cancer [7,11–15].

76    Together, these observations point to a central role for DNA methylation in shaping

77    genome architecture, influencing development, and driving trait variation.

78    Consequently, there is substantial interest in identifying the genetic [16–19] and

79    environmental [20–23] factors that shape DNA methylation levels. Progress toward

80    this goal requires statistical approaches that can handle the complexities of real

81    world, population-based datasets. Here, we present one such approach, designed

82    specifically for analyses of differential methylation levels in bisulfite sequencing

83    datasets.

84    High-throughput bisulfite sequencing approaches, which include whole

85    genome bisulfite sequencing (WGBS or BS-seq) [24], reduced representation

86    bisulfite sequencing (RRBS) [25,26], and sequence capture followed by bisulfite

87    conversion [27,28], are used to estimate genome-wide DNA methylation levels at

4

88  base-pair resolution. All such methods rely on the differential sensitivity of

89  methylated versus unmethylated cytosines to the chemical sodium bisulfite.

90  Specifically, sodium bisulfite converts unmethylated cytosines to uracil (and

91  ultimately thymine following PCR), while methylated cytosines are protected from

92  conversion. Estimates of DNA methylation levels for each cytosine base can thus be

93  obtained directly from high-throughput sequencing data by comparing the number of

94  C's (reflecting an originally methylated version of the base) versus T's (reflecting an

95  originally unmethylated version of the base) at that position in the mapped reads.

96      The raw data produced by bisulfite sequencing methods are therefore count

97  data, in which both the number of methylated reads and the total coverage at a site

98  contain useful information. Higher total coverage corresponds to a more reliable

99  estimate of the true DNA methylation level, which, in a typical experiment, can vary

100  dramatically across individuals and sites (e.g., by several orders of magnitude: S1

101  Figure). Many commonly used methods for testing for differential methylation

102  (whether by genotype, environmental predictor, or experimental perturbation) ignore

103  this variability by converting counts to percentages or proportions (e.g., t-tests,

104  Mann-Whitney U tests, linear models, and all tools initially designed for array-based

105  data [29,30]; Table 1). Thus, a site at which 5 of 10 reads are designated as

106  methylated (i.e., read as a cytosine) is treated identically to a site at which 50 of 100

107  reads are designated as methylated. This assumption reduces the power to uncover

108  true predictors of variation in DNA methylation levels, because it treats noisy

109  measurements the same way as accurate ones.

110

5

111 **Table 1.** Approaches for identifying differentially methylated loci in bisulfite
112 sequencing data sets.

| Statistical method | Directly models counts? | Controls for biological covariates? | Controls for genetic covariance? | Programs that implement the method |
|---|---|---|---|---|
| t-test or Wilcoxon rank-sum test | No | No | No | R and many others |
| Fisher's exact test | Yes | No | No | R and many others |
| Binomial regression | Yes | Yes | No | R and many others |
| Linear regression | No | Yes | No | R and many others |
| Beta-binomial model | Yes | Some[1] | No | DSS [31], MOABS [32], RadMeth [33] |
| Linear mixed model | No | Yes | Yes | GEMMA [34], EMMA [35], EMMAX [36], FaST-LMM [37] |
| Binomial mixed model | Yes | Yes | Yes | MACAU |

113 [1]Only RadMeth; the implementations of the beta-binomial model in MOABS and DSS do not allow the
114 user to control for covariates.
115

116

117      To address this problem, several recently introduced methods for differential

118 DNA methylation analysis implement a beta-binomial model (e.g., 'DSS: Dispersion

119 Shrinkage for Sequencing data' [31], 'RADMeth: Regression Analysis of Differential

120 Methylation' [33], and 'MOABS: Model Based Analysis of Bisulfite Sequencing data'

121 [32]). These methods model the binomial nature of bisulfite sequencing data, while

122 taking into account the well-known problem of over-dispersion in sequencing reads.

123 Because these methods work directly on count data, they can reliably account for

124 variation in read coverage across sites and individuals. Consequently, beta-binomial

125 methods consistently provide increased power to detect true associations between

126 genetic or environmental sources of variance and DNA methylation levels [31–33].

127      However, methods based on beta-binomial models only account for over-

128 dispersion due to independent variation, making them unsuited for data sets

6

129    containing population structure or related individuals. Accounting for genetic

130    relatedness is important because genetic variation can exert strong and pervasive

131    effects on DNA methylation levels [17,19,38,39]. In humans, methylation levels at

132    more than ten thousand CpG sites are influenced by local genetic variation [18], and

133    DNA methylation levels in whole blood are 18%-20% heritable on average, with the

134    heritability estimates for the most heritable loci (top 10%) averaging around 68%

135    [38,39]. As a result, DNA methylation levels will frequently covary with kinship or

136    population structure, and failure to account for this covariance could lead to spurious

137    associations or reduced power to detect true effects. This phenomenon has been

138    extensively documented for genotype-phenotype association studies [35,36,40–42],

139    and controlling for genetic covariance between samples is now a basic requirement

140    for genome-wide association studies. Similar logic applies to analyses of gene

141    regulatory phenotypes and studies of gene expression variation often do take

142    genetic structure into account by using mixed model approaches [43–45]. However,

143    despite growing interest in environmental epigenetics and epigenome-wide

144    association studies (EWAS), none of the currently available count-based methods

145    appropriately control for genetic effects on DNA methylation levels in bisulfite

146    sequencing data (Table 1). Consequently, even though count-based methods have

147    been shown to be more powerful, recent bisulfite sequencing studies have turned to

148    linear mixed models to deal with the confounding effects of population structure

149    [19,46].

150         To address this gap, we present a binomial mixed model (BMM) for

151    identifying differentially methylated sites that directly models raw read counts while

7

152    accounting for both covariance between samples and extra over-dispersion caused

153    by independent noise. We also present an efficient, sampling-based inference

154    algorithm to accompany this model, called MACAU (Mixed model association for

155    count data via data augmentation). MACAU works directly on binomially distributed

156    count data from any high-throughput bisulfite sequencing method (e.g., WGBS,

157    RRBS, targeted sequence capture) and uses random effects to not only model over-

158    dispersion (as in the standard beta-binomial approach [47]), but also to model

159    relatedness/population structure. Hence, MACAU enables users to identify

160    differentially methylated sites in a wide variety of settings, with little cost to power

161    even when genetic effects on DNA methylation levels are negligible.

162          We compared MACAU's performance with currently available methods under

163    two realistic scenarios, using both real bisulfite sequencing data sets (WGBS and

164    RRBS) and simulations parameterized based on properties of real data. In the first

165    scenario, we analyzed publicly available data from *Arabidopsis thaliana* [48] to show

166    that, when a predictor variable of interest is correlated with population structure,

167    MACAU provides better control of type I error than existing methods. This setting is

168    particularly relevant to understanding geographic variation in DNA methylation levels

169    (e.g., [19,48–50]) and for identifying genetic or environmental predictors of DNA

170    methylation in structured samples (e.g., [50,51]). In the second scenario, we used

171    newly generated RRBS data from wild baboons (*Papio cynocephalus*) to

172    demonstrate that MACAU also provides increased power to detect truly differentially

173    methylated sites in the presence of kinship—a condition that often holds in analyses

174    of natural populations (e.g., [48,52,53]) and in tests for epigenetic discordance

8

175    between siblings [22,53–55]. As interest in epigenome-wide association studies

176    (EWAS), environmental epigenetics, and the epigenetic correlates of disease grows,

177    these types of complex data sets will become increasingly common.

178

## Results

180    **The binomial mixed model and the MACAU algorithm**

181    Here, we briefly describe the model and the algorithm. Additional information

182    is provided in the Supplementary Information Text File, which includes details on the

183    model, inference method, and algorithm (including descriptions of the data

184    augmentation approach and efficient MCMC sampling steps).

185    To detect differentially methylated sites, we model each potential target of

186    DNA methylation individually (i.e., we model each CpG site one at a time) as a

187    function of $x$, a predictor variable of interest. Here, $x$ could be a genotype value, as

188    in methylation QTL mapping analyses; an environmental predictor of interest, such

189    as temperature, chemical exposure, or social environment; an individual

190    characteristic, such as age or sex; or an experimental perturbation, as in a

191    treatment-control design. For each site, we consider the following binomial mixed

192    model (BMM):

$$y_i = Bin(r_i, \pi_i), \tag{1}$$

193    where $r_i$ is the total read count for $i$th individual; $y_i$ is the methylated read count for

194    that individual, constrained to be an integer value less than or equal to $r_i$; and $\pi_i$ is

195    an unknown parameter that represents the underlying proportion of methylated

9

196     reads for the individual at the site. We use a logit link to model $\pi_i$ as a linear function

197     of several parameters:

$$log\left(\frac{\pi_i}{1-\pi_i}\right) = \boldsymbol{w}_i^T\boldsymbol{\alpha} + x_i\beta + g_i + e_i, \tag{2}$$

$$\boldsymbol{g} = (g_1, \cdots, g_n)^T \sim MVN(0, \sigma^2 h^2 \boldsymbol{K}), \tag{3}$$

$$\boldsymbol{e} = (e_1, \cdots, e_n)^T \sim MVN(0, \sigma^2(1-h^2)\boldsymbol{I}), \tag{4}$$

198     where, for a data set including *c* covariates and *n* individuals, $\boldsymbol{w}_i$ is a *c*-vector of

199     covariates including an intercept; $\boldsymbol{\alpha}$ is a *c*-vector of corresponding coefficients; $x_i$ is

200     the predictor of interest for individual *i* and $\beta$ is its coefficient; $\boldsymbol{g}$ is an *n*-vector of

201     genetic random effects that model correlation due to population structure or kinship;

202     MVN denotes the multivariate normal distribution; *e* is an *n*-vector of environmental

203     residual errors that model independent variation; $\boldsymbol{K}$ is a known *n* by *n* relatedness

204     matrix that can be calculated based on pedigree or genotype data; $\boldsymbol{I}$ is an *n* by *n*

205     identity matrix; $\sigma^2 h^2$ is the genetic variance component; $\sigma^2(1-h^2)$ is the

206     environmental variance component; and $h^2$ is the heritability of the logit transformed

207     methylation proportion (i.e. $logit(\boldsymbol{\pi})$). Note that $\boldsymbol{K}$ has been standardized to ensure

208     $tr(\boldsymbol{K})/n = 1$, so that $h^2$ lies between 0 and 1 and can be interpreted as heritability

209     (see [56]; *tr* denotes the trace norm).

210        Both $\boldsymbol{g}$ and *e* model over-dispersion (i.e., the increased variance in the data

211     that is not explained by the binomial model). However, they model different aspects

212     of over-dispersion: *e* models the variation that is due to independent environmental

213     noise (a known problem in data sets based on sequencing reads [57–60], including

214     analyses of read proportions [61]), while $\boldsymbol{g}$ models the variation that is explained by

215     kinship or population structure. Effectively, our model improves and generalizes the

216   beta-binomial model by introducing this extra $g$ term to model individual relatedness

217   due to population structure or stratification. In the absence of $g$, our model becomes

218   similar to other beta-binomial models previously developed for modeling count data

219   [31,33,47,62].

220        We are interested in testing the null hypothesis that the predictor of interest

221   has no effect on DNA methylation levels: $H_0: \beta = 0$. This test requires obtaining the

222   maximum likelihood estimate $\hat{\beta}$ from the model. Unlike its linear counterpart,

223   estimating $\hat{\beta}$ from the binomial mixed model is notoriously difficult, as the joint

224   likelihood consists of an *n*-dimensional integral that cannot be solved analytically

225   [63,64]. Standard approaches rely on numerical integration [65] or Laplace

226   approximation [66,67], but neither strategy scales well with the increasing dimension

227   of the integral, which in our case is equal to the sample size. Because of this

228   problem, standard implementations of binomial mixed models often produce biased

229   estimates and overly narrow (i.e., anti-conservative) confidence intervals [68–72]. To

230   overcome this problem, we instead use a Markov chain Monte Carlo (MCMC)

231   algorithm-based approach for inference, using un-informative priors for the hyper-

232   parameters $h^2$ and $\sigma^2$. After drawing accurate posterior samples of $\beta$, we rely on the

233   asymptotic normality of both the likelihood and the posterior distributions [73] to

234   obtain the approximate maximum likelihood estimate $\hat{\beta}$ and its standard error se($\hat{\beta}$).

235   This procedure allows us to construct approximate Wald test statistics and *p*-values

236   for hypothesis testing. Despite the stochastic nature of the procedure, the MCMC

237   errors are small enough to ensure stable *p*-value computation across multiple

238   MCMC runs (S2 Figure). We note that with reasonably large sample sizes (n=50 or

11

239    more), the resulting p-values are also robust to prior perturbation on hyper-

240    parameters (S3 Figure); however, all results reported here are based on calculations

241    with un-informative priors.

242        In addition to the approximate inference procedure described above, we also

243    developed a novel MCMC algorithm based on an auxiliary variable representation of

244    the binomial distribution for efficient, approximate $p$-value computation [74–76] (see

245    SI Text File Section 2: Inference Method Overview and SI Text File Section 3.1:

246    Data Augmentation for more details). We did so to reduce the heavy computational

247    burden of standard MCMC algorithms, which would otherwise be prohibitive in terms

248    of run time for large datasets. Building on the auxiliary variable representation, our

249    main technical contribution is a new framework that approximates the distribution of

250    the auxiliary variables (S4 Figure, S1-S2 Tables) while simultaneously taking

251    advantage of recent innovations for fitting mixed effects models [34,35,37,77] (see

252    SI Text File Sections 3.2 and 3.3). This framework reduces per-MCMC iteration

253    computational complexity from cubic to quadratic with respect to the sample size,

254    and results in an approximate $n$-fold speed up in practice compared with the popular

255    Bayesian software MCMCglmm [78], where $n$ is the sample size (S5 Figure, S3

256    Table; we note that this speed-up is generalizable to other GLMM problems as well).

257    Our implementation of the BMM is therefore efficient for data sets ranging up to

258    hundreds of samples and millions of sites, as computational complexity scales only

259    linearly with respect to the number of analyzed sites (S5 Figure).

260        Because our model effectively includes the beta-binomial model as a special

261    case, we expect it to perform similarly to the beta-binomial model in settings in which

12

262    population structure is absent (we say "effectively" because the beta-binomial model

263    uses a beta distribution to model independent noise while we use a log-normal

264    distribution). However, we expect our model to outperform the beta binomial in

265    settings in which population structure is present. In addition, in the presence of

266    population stratification, we expect the beta-binomial model to produce inflated test

267    statistics (thus increasing the false positive rate) while our model should provide

268    calibrated ones. Below, we test these predictions using two different bisulfite

269    sequencing data sets. We begin with simulations in which the true value of $\beta$ is

270    known, and the over-dispersion parameter and genetic covariance between samples

271    are motivated by the real data sets. We also motivate our choice of simulated

272    sample sizes based on real bisulfite sequencing data sets, which currently range

273    from ~20 – 150 samples [19,26,46,53,79–82]. However, because sample sizes are

274    only likely to grow in the future, for the data set types of most direct interest (i.e.,

275    those that contain population structure and heritable DNA methylation levels) we

276    further consider sample sizes that are much larger than currently represented in the

277    literature (n = 500 and n = 1000). Finally, we apply our model directly to the real

278    data.

279

280    **Count-based models perform well in the absence of genetic effects on DNA**

281    **methylation levels**

282        We first compared the performance of the BMM implemented in MACAU with

283    the performance of other currently available methods for analyzing bisulfite

284    sequencing data in the absence of genetic effects. Intuitively, we expected MACAU

13

285   and the beta-binomial model to perform similarly, and we expected both methods to

286   outperform those that first transform the raw count data. To test our prediction, we

287   simulated the effect of a predictor variable on DNA methylation levels across 5000

288   CpG sites (4500 true negatives and 500 true positives). Motivated by our analysis of

289   age effects on DNA methylation levels in the baboon RRBS data set (below), we

290   conducted this simulation by sampling from a distribution of known age values from

291   the same baboon population. For all simulations, we set the effect of genetic

292   variation on DNA methylation levels equal to zero, which is equivalent to setting

293   either (i) the heritability of DNA methylation levels to zero (unlikely based on prior

294   findings [38,39]), or (ii) studying completely unrelated individuals in the absence of

295   population structure. To explore MACAU's performance across a range of

296   conditions, we simulated age effects on DNA methylation levels across three effect

297   sizes (percent of variance in DNA methylation explained (PVE) = 5%, 10%, or 15%)

298   and three sample sizes (n = 20, 50, and 80). These values capture the majority of

299   effect sizes and sample sizes documented in recent genome-wide bisulfite

300   sequencing studies (e.g., [45,52,53,83]).

301       Because age is naturally modeled as a continuous variable, we focused our

302   comparisons only on approaches that could accommodate continuous predictor

303   variables (comparisons in which we artificially binarized age, which allowed us to

304   include a larger set of approaches, are shown in S6 Figure and S7 Figure for cases

305   excluding and including genetic effects on DNA methylation, respectively; however,

306   binarizing a truly continuous variable consistently results in poorer performance: see

307   S6 Figure versus S9 Figure). Specifically, in addition to the BMM implemented in

14

308    MACAU, we considered the performance of a beta-binomial model, a binomial

309    model, a linear model, and a linear mixed model (implemented in the software

310    GEMMA [34]). For the linear and linear mixed model case, methylation proportions

311    were quantile normalized to a standard normal prior to modeling (see Methods and

312    S8 Figure for parallel results using logit, M-value, and arcsin(sqrt) transformations

313    prior to linear mixed modeling as alternatives to quantile normalization). As

314    expected, we found that MACAU performed similarly to the beta-binomial model,

315    and that these two approaches consistently detected more true positive age effects

316    on DNA methylation levels (at a 10% empirical FDR) than all other methods (S9

317    Figure). For example, in the "easiest" case we simulated (PVE = 15%, n = 80), we

318    found that the beta-binomial model detected 30% of simulated true positives, while

319    the BMM implemented in MACAU detected 27.8%. The slight loss of power in the

320    BMM is a consequence of the smaller degrees of freedom caused by the additional

321    genetic variance component. In comparison, the linear model detected 21.2% of true

322    positives; the linear mixed effects model, 14%; and the binomial model, 8.4% (S9

323    Figure). Although it is often used to test for differential methylation [53,84,85], the

324    binomial model exhibits low power when an empirical FDR is used to control for

325    multiple hypothesis testing due to poor type I error calibration, as has been

326    previously reported [33]. Area under a receiver operating characteristic curve (AUC)

327    was also consistently very similar between the beta-binomial and MACAU (S9

328    Figure), although the advantage of the count-based methods was less clear by this

329    measure. This reduced contrast is because AUC is based on true positive-false

330    positive trade-offs across the entire range of $p$-value thresholds: methods can

331   consequently yield high AUCs even when they harbor little power to detect true

332   positives at FDR thresholds that are frequently used in practice. Taken together, our

333   simulations suggest a general advantage to count-based models for samples that

334   contain no genetic structure. Further, the differences in performance between the

335   beta-binomial model and the BMM implemented in MACAU were consistently small

336   in this setting (S9 Figure).

337

338   **Binomial mixed models control for false positive associations that arise from**

339   **population structure: simulations and a real data example from *Arabidopsis***

340        We next evaluated each model's performance in a more realistic setting, in

341   which genetic covariance between samples could potentially confound tests for

342   environmental or genetic effects on DNA methylation levels. As a case study

343   example, we drew from publicly available phenotype data and SNP genotype data

344   for 24 *Arabidopsis thaliana* accessions [86,87] in which leaf tissue samples had

345   been recently subjected to whole genome bisulfite sequencing [48]. Among these

346   accessions, a secondary dormancy phenotype (measured as the slope of the

347   relationship between length of cold treatment and seed germination percentages

348   [88]) is correlated with population structure ($R^2 = 0.38$ against the first principal

349   component of the genotype matrix for these accessions; $p = 7.84 \times 10^{-4}$; S10

350   Figure). Because secondary dormancy is associated with environmental conditions

351   that are experienced after the seed has already dispersed, we have no expectation

352   that secondary dormancy should be associated with DNA methylation levels in leaf

353   tissue. Consequently, this data set provided the opportunity to evaluate calibration of

354 Type I error (false positives) using MACAU, which controls for population structure,

355 versus other available approaches.

356 To do so, we first used the true distribution of secondary dormancy

357 characteristics and the true genetic structure among these 24 accessions to simulate

358 a dataset that consisted entirely of null associations. Specifically, we simulated data

359 sets (containing 4000 sites each) in which the secondary dormancy had no effect on

360 DNA methylation levels, but the effect of genetic variation on DNA methylation levels

361 was either moderate ($h^2 = 0.3$) or large ($h^2 = 0.6$). Thus, in these data sets,

362 population structure could confound the relationship between the predictor variable

363 (the capacity for secondary dormancy) and DNA methylation levels if not taken into

364 account.

365 As predicted, we found that the BMM implemented in MACAU appropriately

366 controlled for genetic effects on DNA methylation levels: whether DNA methylation

367 levels were moderately ($h^2 = 0.3$) or strongly ($h^2 = 0.6$) heritable, MACAU did not

368 detect any sites associated with secondary dormancy at a relatively liberal false

369 discovery rate threshold of 20% (whether calculated against empirical permutations

370 or calculated using the R package *qvalue* [32]). In addition, the *p*-value distributions

371 for secondary dormancy effects on DNA methylation levels, in both simulations, did

372 not differ from the expected uniform distribution (Fig. 1; Kolmogorov-Smirnov (KS)

373 test when $h^2 = 0.3$: D = 0.015, p = 0.909; when $h^2 = 0.6$: D = 0.016, p = 0.874;

374 genomic control factors: 0.90 when $h^2 = 0.3$, 0.93 when $h^2 = 0.6$). In contrast, when

375 we analyzed the same simulated data sets with a beta-binomial model, we

376 erroneously detected 2 CpG sites associated with secondary dormancy when

17

377    heritability was set to 0.3, and 4 CpG sites when heritability was set to 0.6 (at a 20%

378    FDR in both cases). More concerningly, the distributions of $p$-values produced by the

379    beta-binomial model were significantly different from the expected uniform

380    distribution and skewed towards low (significant) values (KS test when $h^2$ = 0.3: D =

381    0.084, p = 1.75 x $10^{-8}$; when $h^2$ = 0.6: D = 0.096, p = 2.80 x $10^{-11}$; genomic control

382    factors: 1.18 when $h^2$ = 0.3, 1.32 when $h^2$ = 0.6). These results suggest an

383    increasing problem with false positives as the heritability of DNA methylation levels

384    increases (see S11 Figure for similar results when comparing a linear model to a

385    linear mixed model).

386

387    **Fig. 1. MACAU appropriately controls for genetic covariance in**

388    **simulated and real WGBS data and eliminates false positive**

389    **identification of differentially methylated sites.** (A, B) The distribution of $p$-

390    values for 4000 simulated true negative sites (n = 24 accessions; effect of

391    secondary dormancy on DNA methylation levels = 0). For each simulation, $h^2$

392    was set to 0.3 (A) or 0.6 (B). Simulated data were analyzed with a beta-

393    binomial model or MACAU, and compared against the expected uniform

394    distribution. (C) QQ-plots comparing the $p$-value distributions for (i) a model

395    testing for effects of secondary dormancy on DNA methylation levels in real

396    WGBS data, with quantiles plotted on the y-axis; and (ii) the same model

397    when the secondary dormancy values were permuted across individuals, with

398    quantiles plotted on the x-axis. The genomic control factor, λ, is shown for

399    each set of results.

400

18

401

402    Notably, this problem should become more acute with increasing sample size,

403    which provides greater power to detect false positives generated by this type of

404    confounding [89]. Indeed, both increasing the simulated sample size and increasing

405    the simulated correlation between the predictor variable and genetic structure

406    produces increasingly poorly calibrated results. For example, when sample sizes

407    were simulated from 25 up to 1000 individuals (and the heritability of DNA

408    methylation levels was set to 0.6), we observed genomic inflation factors ranging

409    from 1.03 – 3.49 for data sets analyzed with a beta-binomial (Fig. 2a). Not

410    surprisingly, for a dataset of a fixed size, the beta-binomial genomic control factor

411    increased as the confounding between population structure and the predictor

412    variable of interest became more extreme (see S12a Figure for comparable results

413    for a linear model). In contrast, when we analyzed the same simulated datasets with

414    the BMM implemented in MACAU, the genomic control factors consistently ranged

415    from 0.82 – 1.08, even when sample sizes were large and/or the correlation between

416    population structure and the predictor variable was substantial (Fig. 2b; see S12b

417    Figure for comparable results from a linear mixed model). Importantly, these

418    differences in genomic control factors can translate into substantial differences in the

419    results suggested by a given method. For example, when n = 1000 and the predictor

420    variable is highly confounded with population structure ($R^2 = 0.5$), a beta-binomial

421    falsely identified 32% of sites in the data set as differentially methylated (10% FDR),

422    while MACAU correctly identified no differentially methylated sites (10% FDR; S13

423    Figure).

19

424

425     **Fig 2. MACAU controls for genetic covariance in data sets that span a**

426     **range of sample sizes and levels of correlation between population**

427     **structure and a predictor variable of interest.** Genomic control factor when

428     simulated datasets (n=5000 sites per dataset; $h^2 = 0.6$) were analyzed with

429     either (A) a beta-binomial model or (B) a BMM implemented in MACAU. The

430     correlation between the simulated predictor variable and the first principal

431     component of genome-wide genotype data is plotted on the x-axis. Genotype

432     data are for *Arabidopsis* accessions, as reported in [87].

433

434     To investigate the calibration of test statistics in the real data set, we then

435     analyzed the relationship between the secondary dormancy phenotype and WGBS

436     data for the 24 *Arabidopsis* accessions in which both phenotype and WGBS data

437     were available (n = 830,676 CpG sites tested [32,33,34]). We again compared the

438     performance of a simple linear model, a binomial model, a beta-binomial model, the

439     BMM implemented in MACAU, and an LMM implemented in GEMMA. Further

440     illustrating its poor handling of Type I error, the binomial model detected more than

441     100,000 secondary dormancy-associated sites at a 10% empirical FDR threshold,

442     respectively, with a genomic control factor of 3.81. A beta-binomial model

443     substantially improved over the binomial model, but still detected 39 secondary

444     dormancy-associated sites at a 20% empirical FDR threshold, and 150 sites and 690

445     sites at a 10% or 20% FDR *qvalue* threshold, respectively (genomic control factor =

446     1.16). Given the clear confounding of population structure and secondary dormancy

20

447     in this sample, as well as the results of our simulations, these associations are

448     probably largely, if not completely, spurious. In contrast, MACAU, the linear mixed

449     model (GEMMA), and the simple linear model did not identify any CpG sites

450     associated with secondary dormancy, either at a 10% or a 20% false discovery rate

451     threshold (Fig. 1 and S11 Figure; genomic control factors: MACAU – 0.89, GEMMA

452     – 0.97, Linear model – 0.99). Based on our earlier simulations, the similarity of

453     performance among the three approaches likely stems from different reasons: the

454     linear model is poorly powered to detect positive hits with this sample size (either

455     true positives or false positives); the linear mixed model controls for population

456     structure but has low power to detect true associations; while MACAU combines

457     both the increased power conferred by modeling the raw count data with appropriate

458     controls for population structure (see Fig. 1 and results below).

459

460     **MACAU provides increased power to detect true positives in the presence of**

461     **kinship: simulations based on data from baboons**

462           In other data sets, a predictor variable of interest may not be confounded with

463     genetic structure, but modeling genetic similarity between samples could reduce

464     residual error variance and improve power. To investigate this scenario, we focused

465     on the relationship between age and DNA methylation levels in a wild baboon

466     population. Female baboons remain in their natal groups throughout their lives,

467     producing relatedness values that are primarily due to matrilineal descent. The

468     resulting genetic structure is one in which females tend to be more closely related to

469     each other, on average, than males or male-female dyads [90], but in which not all

21

470     females are related (because multiple matrilines co-reside in a single group). Data

471     sets drawn from baboon populations therefore include a substantial number of

472     unrelated individuals, but also some dyads that are genetically non-independent

473     (i.e., relatives: S14 Figure).

474        To test the relative performance of different modeling approaches in this

475     setting, we first simulated moderate to large genetic effects on DNA methylation

476     levels ($h^2$ = 0.3 and 0.6 respectively, as in the *Arabidopsis* simulation above) and

477     relatedness values based on the observed distribution of relatedness values within

478     baboon social groups (n = 80, 500, or 1000 baboons). We again simulated a range

479     of non-zero effect sizes (percent variance explained by age = 5%, 10%, or 15%) for

480     500 true positive sites, and an effect size of zero for 4500 true negative sites.

481        In simulations in which age had a moderate effect on DNA methylation levels

482     (PVE = 10%), MACAU detected 11.4% (when $h^2$ = 0.3) and 20.6% (when $h^2$ = 0.6)

483     of simulated true positives at a 10% empirical FDR, and produced well calibrated p-

484     values for sites with no simulated age effect (S15 Figure). In comparison, the beta-

485     binomial model (the next best model) detected 8.2% and 10.4% of true positives,

486     respectively (Fig. 3). As in the simulations, we again observed that a simple binomial

487     model was prone to type I error, which resulted in failure to detect true age-

488     associated sites when empirical FDRs were calculated against permuted data. Our

489     additional simulations at PVE = 5% or PVE = 15%, and n = 500 or n = 1000,

490     confirmed MACAU's advantage over other methods across a range of conditions

491     (S16-S17 Figure). As expected, the magnitude of this advantage was positively

492     correlated with the heritability of DNA methylation levels.

493

494     **Fig. 3. MACAU exhibits increased power to detect differential**

495     **methylation when DNA methylation levels are heritable**. Receiver

496     operating characteristic (ROC) curves and true positive rates at a 10% false

497     discovery rate threshold for simulated age effects on DNA methylation levels

498     at (A-C) simulated sites with moderately heritable DNA methylation levels ($h^2$

499     = 0.3) and (D-F) simulated sites with highly heritable DNA methylation levels

500     ($h^2$ = 0.6). Panels B and E are enlarged versions of panels A and D,

501     respectively. They focus on false positive rates below 0.1, because the

502     performance of alternative methods at low false positive rates tends to be

503     most important to researchers in practice; that is, it is unlikely to matter if

504     method performance is identical when accepting a 50% false positive rate,

505     which would yield very poor inferential power. Each simulated dataset

506     contained n=80 individuals and 5000 simulated CpG sites, with 500 true

507     positives and 4500 true negatives. Here, we show results where the

508     simulated percent variance explained by age = 10%. A binomial model could

509     not detect true positives at a false positive rate below 0.10 (when $h^2$ = 0.3) or

510     below 0.9 (when $h^2$ = 0.6); the binomial is therefore removed from panel B,

511     and only shown for large false positive rates in panel E.

512

513     **Age-associated DNA methylation levels in wild baboons**

514     Finally, we analyzed the new baboon RRBS data set for differential

515     methylation patterns by age (n = 50, age range = 1.76 – 18.01 years in our sample,

23

516    S4 Table). Because age-related effects on DNA methylation levels are well

517    described, this approach allowed us to not only evaluate MACAU's ability to detect

518    differentially methylated sites, but also to identify known age-related signatures in

519    DNA methylation data [38,39,91–93]. This data set included 433,871 CpG sites,

520    enriched for putatively functional regions of the genome (e.g., genes, gene

521    promoters, CpG islands, as expected in RRBS data sets [25,26]: S18 Figure; see

522    also S19 Figure and S4 Table for additional information on data quality, including

523    bisulfite conversion rates, *MspI* digest efficiency, correlation with gene expression

524    levels, and methylation level distributions by genomic regions). As in our simulations,

525    we found that MACAU provided increased power to detect age effects in the

526    presence of familial relatedness. We detected 1.6-fold more age-associated CpG

527    sites at a 10% empirical FDR using MACAU compared to the results of a beta-

528    binomial model, the next best approach (1.4-fold more sites at a 20% empirical FDR;

529    Fig. 4 and S20 Figure). This advantage was consistently observed across all FDR

530    thresholds we considered, except for relatively low (<7.5%) empirical FDR

531    thresholds, when all of the methods were very low powered as a result of the modest

532    sample size.

533

534    **Fig. 4. Age-associated CpG sites identified by MACAU in the baboon**

535    **RRBS data.** (A) The number of age-associated CpG sites detected at a given

536    empirical FDR. The binomial model cannot detect age-associated sites at a

537    false discovery rate below 0.20 and is consequently removed from the panel.

538    (B) For age-associated sites detected by MACAU (at a 10% FDR), the

24

539        proportion of sites that gain or lose methylation with age is shown by genomic

540        region. Positive = DNA methylation levels increase with age; Negative = DNA

541        methylation levels decrease with age. (C) Age-associated CpG sites detected

542        using MACAU (10% FDR) are more likely to fall near genes that are

543        expressed in whole blood, compared to the background set of CpG sites near

544        genes (**p < $10^{-10}$). Further, age-associated CpG sites are more likely to

545        occur near genes that are differentially expressed (DE) with age, compared to

546        CpG sites near genes that are not DE with age (*p = 0.032).

547

548        We performed several analyses to investigate the likely validity and functional

549   importance of the age-associated CpG sites we identified. Based on the results of

550   previous studies, we expected that age-associated sites in CpG islands would tend

551   to gain methylation with age [92,93], while sites in other regions of the genome (e.g.,

552   CpG island shores, gene bodies) would tend to lose methylation with age [92,93]. In

553   addition, we expected that, in whole blood, bivalent/poised promoters should gain

554   DNA methylation with age, while enhancers should lose methylation with age (as

555   discussed in [91,92,94]). Finally, we expected that stretches of differentially

556   methylated sites (i.e., differentially methylated regions, or DMRs) would tend to

557   occur in or near CpG islands and CpG shores, potentially altering how steeply

558   methylation levels change between islands and their surrounding shelves (e.g., [95]).

559        Our results conformed to these patterns: sites in CpG islands tended to gain

560   methylation with age (71.4% of sites were positively correlated with age); and sites

561   in promoters, CpG island shores, and gene bodies tended to lose methylation with

562    age (72.7%, 75.4%, and 75.2% of sites were negatively correlated with age,

563    respectively; Fig. 4). In addition, we found that positively correlated, age-associated

564    sites were highly enriched in chromatin states associated with bivalent/poised

565    promoters (as defined by the Roadmap Epigenomics Project [96]). Specifically, age-

566    associated CpG sites in bivalent/poised promoters were 3.4 times more likely to

567    show increases in DNA methylation with age, compared to age-associated CpG

568    sites in other regions ($p < 10^{-10}$, Fisher's exact test). Negatively correlated age-

569    associated sites (i.e., sites where DNA methylation levels decreased with age) were

570    strongly enriched in enhancers (defined as sites either marked by H3K4me1 in

571    human PBMCs [97] or sites within chromatin states annotated as 'enhancers' by the

572    Roadmap Epigenomics Project [96], $p = 2 \times 10^{-4}$, Fisher's exact test). Finally, we

573    detected 142 age-related DMRs, the majority of which were found in CpG islands,

574    shores, and bridging islands and shores (S21 Figure and S5 Table).

575        We also reasoned that true positive age-associated CpG sites should contain

576    information about age-associated gene expression levels. To test this hypothesis,

577    we turned to previously generated whole blood RNA-seq data [43] from the same

578    baboon population (n = 63; only four baboons in the RNA-seq data set were also

579    included in the DNA methylation data set). Overall, we observed a strong enrichment

580    of differentially methylated CpG sites in or near (within 10 kb) blood-expressed

581    genes (n = 12,018 genes), compared to the background set of all CpG sites near

582    genes (Fisher's exact test, $p < 10^{-10}$). Further, CpG sites near age-associated genes

583    (n = 1396 genes, 10% FDR) were 30.5% more likely to be differentially methylated

584    with age compared to the background set of all CpG sites near genes (Fisher's exact

26

585    test, p = 0.032; Figure 4). Notably, this enrichment was almost always stronger for

586    the set of differentially methylated sites identified by MACAU than for the same

587    number of top sites identified when running the linear model, linear mixed model,

588    binomial, or beta-binomial approaches, across different FDR thresholds (S22

589    Figure).

590

591

## Discussion

593        DNA methylation levels can have potent effects on downstream gene

594    regulation, and, in doing so, can shape key behavioral, physiological, and disease-

595    related phenotypes [7,20,98–100]. These observations have motivated an increasing

596    number of DNA methylation studies in humans and other organisms, highlighting the

597    need for sophisticated statistical methods that can accommodate the complexities of

598    a broad array of data sets [19,46]. Here, we demonstrate that the binomial mixed

599    model implemented in our software MACAU can (i) effectively control for

600    confounding relationships between genetic background and a predictor variable of

601    interest and (ii) provide increased power to detect true sources of variance in DNA

602    methylation levels in data sets that contain kinship or population structure. In

603    addition, MACAU provides increased flexibility over current count-based methods

604    that cannot accommodate biological replicates (e.g., Fisher's exact test), continuous

605    predictor variables (e.g., DSS, MOABS, RadMeth), or biological or technical

606    covariates (e.g., MOABS, DSS; see also Table 1). Given the increasing interest in

607    both the environmental [21,101,102] and genetic [16,17,19,103] architecture of DNA

608    methylation levels, we believe MACAU will be a useful tool for generalizing

609    epigenomic studies to a larger range of populations. MACAU is particularly well

610    suited to data sets that contain related individuals or population structure; notably,

611    several major population genomic resources contain structure of these kinds (e.g.,

612    the HapMap population samples [104], the Human Genome Diversity Panel [105],

613    and the 1000 Genomes Project in humans [106]; the Hybrid Mouse Diversity Panel

614    in mice [107]; and the 1001 Genomes Project in *Arabidopsis* [108]).

615         Indeed, our results suggest MACAU is a useful tool even in data sets that are

616    less affected by genetic structure, or when the heritability of DNA methylation levels

617    is unclear. Because the beta-binomial model is effectively incorporated as a special

618    case, MACAU exhibits only a slight loss of power relative to a beta-binomial model

619    without genetic random effects when $h^2 = 0$, while conferring better power and better

620    test statistic calibration when $h^2 > 0$ (S9, S16-S17 Figures, Fig. 1). Previous studies

621    in humans have shown that, while the heritability of DNA methylation levels varies

622    across loci, an appreciable proportion of loci are either modestly ($h^2 \geq 0.3$: 21.06% of

623    all CpG sites) or highly ($h^2 \geq 0.6$: 6.95% of all CpG sites) heritable [39,109]. Further,

624    DNA methylation QTLs are widespread across the genome [18,38,103]. Thus,

625    because investigators will rarely have *a priori* knowledge of the heritability of DNA

626    methylation levels at a given locus, and because the advantage of a beta-binomial

627    model is small even when heritability is zero, we recommend applying MACAU in

628    cases in which genetic effects on DNA methylation levels are poorly understood. In

629    addition, our model provides a natural framework for incorporating the spatial

630    dependency of DNA methylation levels across neighboring sites [110,111], which we

28

631    expect to increase power even further [110,111]. However, we do note that, even

632    with the efficient algorithm implemented here, fitting the binomial mixed model (or its

633    extensions) remains more computationally expensive than other approaches for

634    moderately sized datasets (S3 Table). While it remains appropriate for the sample

635    sizes used in current studies (e.g., dozens to hundreds of individuals), or even larger

636    with the support of a moderate-sized computing cluster (because MACAU is easily

637    parallelizable with respect to sites), rapid increases in sample size—especially in the

638    context of EWAS—strongly motivate additional algorithm development to scale up

639    the binomial mixed model for data sets that include thousands or tens of thousands

640    of individuals. This is particularly important given that methods tailored for other

641    types of studies (e.g., quantile normalization followed by linear mixed modeling or

642    *voom + limma*, both commonly used for RNA-seq) do not appear to translate well to

643    bisulfite sequencing data sets (Figure S8; see Methods for additional information on

644    the *voom + limma* comparison).

645        Although we developed MACAU with the analysis of bisulfite sequencing data

646    in mind, we note that a count-based binomial mixed model may be an appropriate

647    tool in other settings as well. For example, allele-specific gene expression (ASE) can

648    be measured in RNA-seq data by comparing the number of reads originating from a

649    given variant to the total number of mapped reads for that site [77,112–114].

650    Similarly, alternative isoform usage can be represented as a proportion of reads

651    containing a non-constitutive exon versus the total reads for the same gene [47].

652    The structure of these data are highly similar to the structure of bisulfite sequencing

653    data, which focus on counts of methylated versus total reads. Unsurprisingly, beta-

29

654   binomial models have also emerged as one of the most popular methods for

655   estimating both ASE values [114–116] and alternative isoform usage [47].

656   Researchers interested in the predictors of variation in either of these measures —

657   which could include *trans*-acting genetic effects, environmental conditions, or

658   properties of the individual (e.g., sex or disease status) — might also benefit from

659   using MACAU. Recent work from the TwinsUK study motivates the need for such a

660   model: Grundberg et al. demonstrated a strong heritable component to ASE levels

661   [117], which could be effectively taken into account using the random effects

662   approach implemented here.

663        Finally, linear mixed models have been recently proposed to account for cell

664   type heterogeneity in epigenome-wide association studies focused on array data

665   [118]. In this framework, the random effect covariance structure is based on overall

666   covariance in DNA methylation levels between samples, which is assumed to be

667   largely attributable to variation in tissue composition. MACAU provides a potential

668   avenue for extending these ideas to sequencing-based data sets.

669

670   **Materials and Methods**

671   ***Arabidopsis thaliana* whole genome bisulfite sequencing (WGBS) data set**

672        We downloaded publicly available WGBS data generated by Schmitz et al.

673   [48], as well as previously published SNP genotype data [87] and secondary

674   dormancy data [86] for 24 *Arabidopsis* accessions. We used the SNP genotype data

675   (specifically, 188,093 sites with minor allele frequency >5%) to construct a pairwise

676   genetic relatedness matrix, $K$, as the product of a standardized genotype matrix $X$,

677    or $K=XX^T/p$ [56], where genotypes were expressed as 0, 1, or 2 depending on the

678    number of reference alleles for that site-sample combination. We used this estimate

679    of $K$ for both the simulations and our analyses of the real WGBS data.

680         In these analyses, we focused on CpG sites measured in ≥50% of

681    accessions, and excluded sites that were constitutively hypermethylated (average

682    DNA methylation level >0.90) or hypomethylated (average DNA methylation level

683    <0.10, following [101,118]). We also excluded highly invariable sites (i.e., sites

684    where the standard deviation of DNA methylation levels fell in the lowest 5% of the

685    overall data set) and sites with very low coverage (i.e., sites where the mean

686    coverage fell in the lowest quartile for the overall data set, below a mean of 3.34

687    reads). After filtering, the final data set consisted of 830,676 sites.

688         For the analysis of test statistic calibration as a function of sample size (Fig.

689    2), we also used *Arabidopsis* data, but simulated the phenotype data as a function of

690    genetic covariance between the accessions. Genotype data were obtained from [87].

691

692    **Baboon reduced representation bisulfite sequencing (RRBS) data set**

693         **Study subjects and sample collection.** To investigate age effects on DNA

694    methylation levels, in both real and simulated data sets, we drew on data and

695    samples from a wild population of yellow baboons in the Amboseli ecosystem of

696    southern Kenya. This population has been monitored for over four decades by the

697    Amboseli Baboon Research Project (ABRP) [119], and the ages of animals born in

698    the study population (n = 37; 74% of the data set) were therefore known to within a

699    few days' error. For animals that immigrated into the study population, ages were

31

700    estimated from morphological features by trained observers (n = 13; 26% of the data

701    set) [120]. Pairwise relatedness values were calculated based on previously

702    collected microsatellite data (14 highly variable loci) [121,122], using the likelihood-

703    based estimator of Lynch and Ritland [123] implemented in the program

704    COANCESTRY [124]. Using the age and relatedness data sets, we simulated age

705    effects on DNA methylation levels for either n = 20, 50, or 80 baboons from a single

706    social group. For simulations with larger sample sizes, we extrapolated both age

707    values and pairwise relatedness values from the n = 80 dataset to maintain the

708    same level of age variation and genetic structure; notably, our results are highly

709    stable in the face of realistic levels of noise in the estimate of $K$ (S23 Figure). In

710    addition, we used previously collected blood samples from the Amboseli population,

711    paired with age and microsatellite genotype records, to investigate age effects on

712    DNA methylation levels in a newly generated RRBS data set.

713        To generate the new RRBS data, we used whole blood samples collected

714    from 50 animals (35 males and 15 females) by the ABRP between 1989 and 2011

715    following well-established procedures [43,125,126]. Briefly, animals were

716    immobilized by an anesthetic-bearing dart delivered through a hand-held blow gun.

717    They were then quickly transferred to a processing site for blood sample collection.

718    Following sample collection, study subjects were allowed to regain consciousness in

719    a covered holding cage until they were fully recovered from the effects of the

720    anesthetic. Upon recovery, study subjects were released near their social group and

721    closely monitored. Blood samples were stored at the field site or at an ABRP-

32

722    affiliated lab at the University of Nairobi until they were transported to the United

723    States.

724        Importantly, given the large range in sample collection dates, we observed no

725    correlation between the age of our study subjects at sample collection and sample

726    age (i.e., time since the collection date; Spearman rank correlation, p = 0.779).

727    Further, to ensure that variation in sample collection dates did not influence our

728    results, we also controlled for sample age as a covariate in our final analyses of the

729    RRBS dataset (see *Analysis of age-related changes in DNA methylation levels*).

730

731        **RRBS data generation and low-level processing.** Genomic DNA was

732    extracted from whole blood samples using the DNeasy Blood and Tissue Kit

733    (QIAGEN) according to the manufacturer's instructions. RRBS libraries were created

734    from 180 ng of genomic DNA per individual, following the protocol by Boyle et al.

735    [25]. In addition, 1 ng of unmethylated lambda phage DNA (Sigma Aldrich) was

736    incorporated into each library to assess the efficiency of the bisulfite conversion

737    (>98% in all case: S4 Table). All RRBS libraries were sequenced using 100 bp

738    single end sequencing on an Illumina HiSeq 2000 platform, yielding a mean of 28.97

739    ±8.97 million reads per analyzed sample (range: 9.59 – 79.78 million reads; Table

740    S4).

741        We removed adaptor contamination and low-quality bases from all reads

742    using the program TRIMMOMATIC [127]. We then mapped the trimmed reads to the

743    olive baboon genome (*Panu* 2.0) using BSMAP, a tool designed for high-throughput

744    DNA methylation data [128]. We used a Python script packaged with BSMAP to

33

745    extract the number of reads as cytosine (reflecting an originally methylated base)

746    and the total read count for each individual and CpG site. We performed the same

747    set of filtering steps described for the *Arabidopsis* WGBS data set to produce our

748    final data set for the baboons. Specifically, we excluded sites that were constitutively

749    hypermethylated or hypomethylated, sites that were highly invariable, and sites that

750    had low average coverage across individuals (in this case, the lowest quartile for

751    mean coverage levels was 4.74 reads). The final filtered data set consisted of

752    433,871 CpG sites.

753

754    **Simulations**

755         To simulate the methylated read counts and total read counts that result from

756    WGBS and RRBS, we performed the following procedure:

757         First, we simulated the proportion of methylated reads for each site. To do so,

758    we drew secondary dormancy values or age values, $x$, as the predictor of interest,

759    from the actual values for the *Arabidopsis* accessions or from the baboon

760    population, respectively. For simulations that focused on *Arabidopsis* data sets of

761    various sizes (e.g., Figure 2), we simulated $x$ and varied the degree to which it was

762    confounded with population structure. Specifically, for each dataset (ranging from

763    n=25 to n=1000 accessions) we performed principal components analysis on the

764    SNP genotype data, and extracted the first principal component to capture the major

765    axis of population structure (PC1). We then added environmental noise from a zero-

766    centered normal distribution to achieve a correlation ($R^2$) between the simulated

767    phenotype and PC1 that reached the desired value (ranging from $R^2 = 0.1$ to 0.5).

768    For each simulated data set, we simulated the DNA methylation level at each

769    CpG site, $\pi$, as a linear function of $x$ and its effect size, $\beta$. In addition, we included

770    the effects of genetic variation ($g$) and random environmental variation ($e$), passed

771    through a logit link (based on the model described in the Results section).

772    For the baboon RRBS and the *Arabidopsis* WGBS simulations, we

773    determined *K* from 14 highly variable microsatellite loci or from the publicly available

774    SNP data, as described above. For each simulation, we set $h^2$ to 0, 0.3, or 0.6 to

775    simulate non-heritable, modestly heritable, or highly heritable DNA methylation

776    levels. We also estimated the variance term $\sigma^2$ from the real data sets. Specifically,

777    we took the mean estimate of $\sigma^2$ across all sites (calculated in MACAU) for each real

778    data set, and used this value as the fixed value of $\sigma^2$ in the corresponding

779    simulations.

780    Next, for each site, we simulated total read counts $r_i$ for each individual *i* from

781    a negative binomial distribution that models the extra variation observed in the real

782    data:

$$r_i \sim NB(t, p), \tag{5}$$

783    where *t* and *p* are site specific parameters estimated from the real data. Specifically,

784    we generated 10,000 sets of *t* and *p* parameters by fitting a negative binomial

785    distribution to the total read count data from 10,000 randomly selected CpG sites in

786    the real baboon RRBS data set or the real *Arabidopsis* data set, using the function

787    'fitdistr' in the R package *MASS* [129]. To simulate counts for a given CpG site, we

788    randomly selected one of these parameter sets to produce the total number of

789    reads. Finally, we simulated the number of methylated reads for each individual at

790    that locus ($y$) by drawing from a binomial distribution parameterized by the number

791    of total reads ($r$) and the DNA methylation level ($\pi$).

792

793    **Comparison of MACAU to existing methods**

794         For all simulated and real data sets, we used raw methylated and total read

795    counts to compare the results of a beta-binomial model (using a custom R script), a

796    binomial model (implemented via 'glm' in R), and the binomial mixed model

797    implemented in MACAU. For computation time comparison, we used the

798    MCMCglmm software, which also provides an implementation of a binomial mixed

799    model [78]. In addition, we used the same count data to run a Fisher's exact test

800    (implemented in R), DSS [31], and RadMeth [33] in the subset of analyses that

801    utilized these programs. To analyze simulated and real data sets using a linear

802    model (implemented using '*lm*' in R) or the linear mixed model implemented in

803    GEMMA [34], we estimated DNA methylation levels by dividing the number of

804    methylated reads by the total read count for each individual and CpG site. We then

805    quantile normalized the resulting proportions for each CpG site to a standard normal

806    distribution, and imputed any missing data using the K-nearest neighbors algorithm

807    in the R package *impute* [130].

808         In addition to the quantile normalization approach, we also evaluated three

809    other methods for transforming methylation proportions: a logit transformation,

810    following [110]; the "M" value transformation ($\log_2$((methylated counts +

811    $\alpha$)/(unmethylated counts + $\alpha$)), where $\alpha = 0.01$, following [30]; and an arcsin square

812    root transformation, following [131]. All four approaches produced qualitatively

36

813    identical results (S8 Figure), so we elected to concentrate on the results from

814    quantile normalization in the main text. Finally, we also tested the performance of a

815    powerful, commonly used method for modeling RNA-seq data: the combination of

816    the *voom* function for data weighting with *limma*, a linear model approach [132]. Our

817    results indicated that *voom* + *limma* performs more poorly than even a simple linear

818    model (S24 Figure), probably because read depth variation is much more

819    complicated in bisulfite sequencing studies than in RNA-seq studies (Figure S1).

820    Because *voom* + *limma* also cannot account for population structure, we report

821    these results in the SI but focus on results from the simple linear model in the main

822    text.

823        To compute empirical false discovery rates in simulated data, we divided the

824    number of false positives detected at a given *p*-value threshold by the total number

825    of sites called by the model as significant at that threshold (i.e., the sum of false

826    positives and true positives). To compute empirical false discovery rates in the real

827    data, in which the false positives and true positives were unknown, we used

828    permutations. Specifically, we permuted the predictor variable for each data set four

829    times, reran our analyses, and then calculated the false discovery rate as the

830    average number of sites detected at a given *p*-value threshold in the permuted data

831    divided by the total number of sites detected at that threshold in the real data. For

832    simulated data sets only, we also calculated the area under the receiver operating

833    characteristic curve (AUC) to produce a measure of the overall tradeoff between

834    detecting true positives and calling false positives.

835

**Analysis of age-related changes in DNA methylation levels**

836

837    Our initial analyses of the baboon RRBS dataset focused only on the relative

838    ability of each method to detect age-associated sites. For these analyses, we

839    therefore did not control for other biological covariates that may contribute to

840    variance in DNA methylation levels (note that biological covariates cannot be

841    incorporated into several implementations of the beta-binomial model [31,32]: see

842    Table 1). However, to investigate patterns of age-related changes in DNA

843    methylation levels, and to compare them to previously described patterns in the

844    literature, we wished to control for such covariates. To do so, we reran the

845    differential methylation analysis in MACAU, this time controlling for sex, sample age,

846    and efficiency of the bisulfite conversion rate estimated from the lambda phage

847    spike-in.

848    First, we investigated whether age-associated sites were enriched in

849    functionally coherent regions of the genome, many of which have previously been

850    identified as age-related [38,92,93]. To do so, we defined gene bodies as the

851    regions between the 5'-most transcription start site (TSS) and 3'-most transcription

852    end site (TES) of each gene using *Panu* 2.0 annotations from Ensembl [133]. We

853    defined promoter regions as the 2 kb upstream of the TSS. CpG were annotated

854    based on the UCSC Genome Browser track for baboon [134], with CpG island

855    shores defined as the 2 kb regions flanking either side of the CpG island boundary

856    (following [26,135,136]). Finally, because no enhancer annotations are available that

857    are specific to baboons, we used H3K4me1 ChIP-seq data generated by ENCODE

858    (from human peripheral blood mononuclear cells) to define enhancer regions [97]. In

38

859   addition, we used chromatin state annotations from the Roadmap Epigenomics

860   Project (also generated from human peripheral blood mononuclear cells) to further

861   investigate biases in the locations of age-associated sites [96]. Using these

862   annotation sets, we performed Fisher's Exact Tests to ask whether age-associated

863   sites were enriched or underrepresented in specific genomic regions. To identify

864   differentially methylated regions (DMRs), we used the criteria proposed by [137].

865   Specifically, DMRs contained at least 3 differentially methylated sites with an inter-

866   CpG distance ≤1 kb, with only 3 non-differentially methylated sites permitted in the

867   DMR as a whole.

868        Second, we asked whether differentially methylated sites were more likely to

869   fall close to blood-expressed genes. For this comparison, we drew on previously

870   published RNA-seq data, generated from whole blood samples collected in the

871   Amboseli baboon population [43]. We defined blood-expressed genes as those

872   genes that had non-zero counts in more than 10% of individuals in the RNA-seq

873   data sets, and that had mean read counts greater than or equal to 10. We then

874   compared the number of differentially methylated CpG sites near blood-expressed

875   genes (i.e., within the gene body or within 10 kb of the gene TSS or TES) to the

876   number of differentially methylated CpG sites near genes that were not expressed in

877   blood, using a Fisher's Exact Test.

878        Finally, we investigated whether CpG sites that occur near genes that are

879   differentially expressed with age were also more likely to be differentially methylated

880   with age. For this comparison, we defined 'age-associated genes' as genes

881   differentially expressed with age (at a 10% FDR) in the RNA-seq data set [43]. We

39

882    compared the number of differentially methylated CpG sites near blood-expressed,

883    age-associated genes to the number of differentially methylated CpG sites near

884    genes that were not within this set of genes, again using a Fisher's Exact Test.

885

886    **Ethics statement**

887    The baboon data used in this study was generated from samples collected

888    from wild baboons living in the Amboseli ecosystem of southern Kenya. This

889    research is conducted under the authority of the Kenya Wildlife Service (KWS), the

890    Kenyan governmental body that oversees wildlife (permit number

891    NCST/RCD/12B/012/57 to Jenny Tung). As the animals are members of a wild

892    population, KWS requires that we do not interfere with injuries to study subjects

893    inflicted by predators, conspecifics, or through other naturally occurring events.

894    Permission to perform temporary immobilizations (for blood sample collection) was

895    granted by KWS; further, these immobilizations were supervised by a KWS-

896    approved Kenyan veterinarian, who monitored anesthetized animals for

897    hypothermia, hyperthermia, and trauma (no such events occurred during our sample

898    collection efforts). Observational and sample collection protocols were approved

899    though IACUC committees at Duke University (current protocol is A020-15-01 to

900    Jenny Tung and Susan C. Alberts).

901

902    **Software and data availability**

903    The MACAU software and a custom script for implementing a beta-binomial

904    model in R is available at: www.xzlab.org/software.html. Previously published data

40

905 sets are available at http://bergelson.uchicago.edu/regmap-data/regmap.html/

906 (*Arabidopsis* SNP genotype data); http://www.ncbi.nlm.nih.gov/geo/ (*Arabidopsis*

907 WGBS data: GSE43857);

908 http://www.nature.com/nature/journal/v465/n7298/full/nature08800.html#supplement

909 ary-information (*Arabidopsis* phenotype data); and http://www.ncbi.nlm.nih.gov/sra

910 (Baboon RNA-seq data: GSE63788). Baboon RRBS data generated in this study are

911 deposited in NCBI (project accession SRP058411).

912

913 **Acknowledgments**

41

928

## **Supporting Information**

929

930     Text S1: Supplementary text

931     Figures S1-S24: Supplementary figures

932     Tables S1-S5: Supplementary tables

933

## References

934

935  1.  Mohandas T, Sparkes R, Shapiro L. Reactivation of an inactive human X
936      chromosome: evidence for X inactivation by DNA methylation. Science.
937      1981;211: 393–396.

938  2.  Li E, Beard C, Jaenisch R. Role for DNA methylation in genomic imprinting.
939      Nature. 1993;366: 362–365. doi:10.1038/366362a0

940  3.  Jones P. Functions of DNA methylation: islands, start sites, gene bodies and
941      beyond. Nat Rev Genet. 2012;13: 484–92. doi:10.1038/nrg3230

942  4.  Kakutani T, Jeddeloh J, Richards EJ. Characterization of an Arabidopsis
943      thaliana DNA hypomethylation mutant. Nucleic Acids Res. 1995;23: 130–137.

944  5.  Ronemus MJ, Galbiati M, Ticknor C, Chen J, Dellaporta SL. Demethylation-
945      induced developmental pleiotropy in Arabidopsis. Science. 1996;273: 654–
946      657. doi:10.1126/science.273.5275.654

947  6.  Finnegan EJ, Peacock WJ, Dennis ES. Reduced DNA methylation in
948      Arabidopsis thaliana results in abnormal plant development. Proc Natl Acad
949      Sci. 1996;93: 8449–8454. doi:10.1073/pnas.93.16.8449

950  7.  Rakyan VK, Beyan H, Down T, Hawa MI, Maslau S, Aden D, et al.
951      Identification of type 1 Diabetes-associated DNA methylation variable
952      positions that precede disease diagnosis. PLoS Genet. 2011;7: 1–9.
953      doi:10.1371/journal.pgen.1002300

954  8.  Dayeh T, Volkov P, Salö S, Hall E, Nilsson E, Olsson AH, et al. Genome-wide
955      Dna methylation analysis of human pancreatic islets from type 2 diabetic and
956      non-diabetic donors identifies candidate genes that influence insulin secretion.
957      PLoS Genet. 2014;10. doi:10.1371/journal.pgen.1004160

958  9.  De Jager PL, Srivastava G, Lunnon K, Burgess J, Schalkwyk LC, Yu L, et al.
959      Alzheimer's disease: early alterations in brain DNA methylation at ANK1,
960      BIN1, RHBDF2 and other loci. Nat Neurosci. 2014;17: 1156–1163.
961      doi:10.1038/nn.3786

962  10. Bakulskia K, Dolinoya D, Sartorb M, Paulsond H, Konend J, Liebermane A, et
963      al. Genome-wide DNA methylation differences between late-onset Alzheimer's
964      disease and cognitively normal controls in human frontal cortex. J Alzheimers
965      Dis. 2012;29: 1–28. doi:10.3233/JAD-2012-111223.Genome-Wide

966  11. Liu Y, Aryee MJ, Padyukov L, Fallin MD, Hesselberg E, Runarsson A, et al.
967      Epigenome-wide association data implicate DNA methylation as an

968          intermediary of genetic risk in rheumatoid arthritis. Nat Biotechnol. 2013;31:
969          142–147. doi:10.1038/nbt.2487

970   12.   Irizarry R, Ladd-Acosta C, Wen B, Wu Z, Montano C, Onyango P, et al. The
971          human colon cancer methylome shows similar hypo- and hypermethylation at
972          conserved tissue-specific CpG island shores. Nat Genet. 2009;41: 178–86.
973          doi:10.1038/ng.298

974   13.   Gluckman PD, Hanson M, Buklijas T, Low FM, Beedle AS. Epigenetic
975          mechanisms that underpin metabolic and cardiovascular diseases. Nat Rev
976          Endocrinol. 2009;5: 401–8. doi:10.1038/nrendo.2009.102

977   14.   Suarez-Alvarez B, Rodriguez RM, Fraga MF, López-Larrea C. DNA
978          methylation: a promising landscape for immune system-related diseases.
979          Trends Genet. 2012;28: 506–14. doi:10.1016/j.tig.2012.06.005

980   15.   Aran D, Sabato S, Hellman A. DNA methylation of distal regulatory sites
981          characterizes dysregulation of cancer genes. Genome Biol. 2013;14: R21.
982          doi:10.1186/gb-2013-14-3-r21

983   16.   Shah S, McRae AF, Marioni RE, Harris SE, Gibson J, Henders AK, et al.
984          Genetic and environmental exposures constrain epigenetic drift over the
985          human life course. Genome Res. 2014; doi:10.1101/gr.176933.114

986   17.   Bell JT, Pai A, Pickrell JK, Gaffney DJ, Pique-Regi R, Degner JF, et al. DNA
987          methylation patterns associate with genetic and gene expression variation in
988          HapMap cell lines. Genome Biol. 2011;12: R10. doi:10.1186/gb-2011-12-1-r10

989   18.   Banovich NE, Lan X, Mcvicker G, Degner JF, Blischak JD, Roux J, et al.
990          Methylation QTLs are associated with coordinated changes in transcription
991          factor binding, histone modifications, and gene expression levels. PLoS
992          Genet. 2014;10: 1–12. doi:10.1371/journal.pgen.1004663

993   19.   Dubin MJ, Zhang P, Meng D, Remigereau M, Osborne EJ, Casale FP, et al.
994          DNA methylation variation in Arabidopsis has a genetic basis and appears to
995          be involved in local adaptation. eLife. 2015;4: e05255.
996          doi:10.7554/eLife.05255

997   20.   Weaver ICG, Cervoni N, Champagne F a, D'Alessio AC, Sharma S, Seckl JR,
998          et al. Epigenetic programming by maternal behavior. Nat Neurosci. 2004;7:
999          847–54. doi:10.1038/nn1276

1000   21.   Waterland R a, Kellermayer R, Laritsky E, Rayco-Solon P, Harris RA,
1001          Travisano M, et al. Season of conception in rural gambia affects DNA
1002          methylation at putative human metastable epialleles. PLoS Genet. 2010;6:
1003          e1001252. doi:10.1371/journal.pgen.1001252

1004 22. Heijmans BT, Tobi EW, Stein AD, Putter H, Blauw GJ, Susser ES, et al.
1005      Persistent epigenetic differences associated with prenatal exposure to famine
1006      in humans. Proc Natl Acad Sci. 2008;105: 17046–9.
1007      doi:10.1073/pnas.0806560105

1008 23. Wolff GL, Kodell RL, Moore SR, Cooney C. Maternal epigenetics and methyl
1009      supplements affect agouti gene expression in Avy/a mice. Am Soc Exp Biol.
1010      1998;12: 949–57.

1011 24. Cokus SJ, Feng S, Zhang X, Chen Z, Merriman B, Haudenschild CD, et al.
1012      Shotgun bisulphite sequencing of the Arabidopsis genome reveals DNA
1013      methylation patterning. Nature. 2008;452: 215–219. doi:10.1038/nature06745

1014 25. Boyle P, Clement K, Gu H, Smith Z. Gel-free multiplexed reduced
1015      representation bisulfite sequencing for large-scale DNA methylation profiling.
1016      Genome Biol. 2012;13: R92. doi:10.1186/gb-2012-13-10-R92

1017 26. Gu H, Smith ZD, Bock C, Boyle P, Gnirke A, Meissner A. Preparation of
1018      reduced representation bisulfite sequencing libraries for genome-scale DNA
1019      methylation profiling. Nat Protoc. 2011;6: 468–81. doi:10.1038/nprot.2010.190

1020 27. Ivanov M, Kals M, Kacevska M, Metspalu A, Ingelman-Sundberg M, Milani L.
1021      In-solution hybrid capture of bisulfite-converted DNA for targeted bisulfite
1022      sequencing of 174 ADME genes. Nucleic Acids Res. 2013;41.
1023      doi:10.1093/nar/gks1467

1024 28. Deng J, Shoemaker R, Xie B, Gore A, LeProust EM, Antosiewicz-Bourget J, et
1025      al. Targeted bisulfite sequencing reveals changes in DNA methylation
1026      associated with nuclear reprogramming. na. 2009;27: 353–60.
1027      doi:10.1038/nbt.1530

1028 29. Aryee MJ, Jaffe AE, Corrada-Bravo H, Ladd-Acosta C, Feinberg AP, Hansen
1029      KD, et al. Minfi: A flexible and comprehensive Bioconductor package for the
1030      analysis of Infinium DNA methylation microarrays. Bioinformatics. 2014;30:
1031      1363–1369. doi:10.1093/bioinformatics/btu049

1032 30. Du P, Zhang X, Huang C-C, Jafari N, Kibbe W, Hou L, et al. Comparison of
1033      Beta-value and M-value methods for quantifying methylation levels by
1034      microarray analysis. BMC Bioinformatics. 2010;11: 587. doi:10.1186/1471-
1035      2105-11-587

1036 31. Feng H, Conneely KN, Wu H. A Bayesian hierarchical model to detect
1037      differentially methylated loci from single nucleotide resolution sequencing data.
1038      Nucleic Acids Res. 2014;42: 1–11. doi:10.1093/nar/gku154

1039  32.  Sun D, Xi Y, Rodriguez B, Park HJ, Tong P, Meong M, et al. MOABS: model
1040       based analysis of bisulfite sequencing data. Genome Biol. 2014;15: R38.
1041       doi:10.1186/gb-2014-15-2-r38

1042  33.  Dolzhenko E, Smith AD. Using beta-binomial regression for high-precision
1043       differential methylation analysis in multifactor whole-genome bisulfite
1044       sequencing experiments. BMC Bioinformatics. 2014;15: 215.
1045       doi:10.1186/1471-2105-15-215

1046  34.  Zhou X, Stephens M. Genome-wide efficient mixed-model analysis for
1047       association studies. Nat Genet. 2012;44: 821–4. doi:10.1038/ng.2310

1048  35.  Kang HM, Zaitlen N, Wade CM, Kirby A, Heckerman D, Daly MJ, et al.
1049       Efficient control of population structure in model organism association
1050       mapping. Genetics. 2008;178: 1709–23. doi:10.1534/genetics.107.080101

1051  36.  Kang H, Sul J, Zaitlen N, Kong S, Freimer NB, Sabatti C, et al. Variance
1052       component model to account for sample structure in genome-wide association
1053       studies. Nat Genet. 2010;42: 348–354. doi:10.1038/ng.548.Variance

1054  37.  Lippert C, Listgarten J, Liu Y, Kadie CM, Davidson RI, Heckerman D. FaST
1055       linear mixed models for genome-wide association studies. Nat Methods.
1056       2011;8. doi:10.1038/nmeth.1681

1057  38.  Bell JT, Tsai PC, Yang TP, Pidsley R, Nisbet J, Glass D, et al. Epigenome-
1058       wide scans identify differentially methylated regions for age and age-related
1059       phenotypes in a healthy ageing population. PLoS Genet. 2012;8.
1060       doi:10.1371/journal.pgen.1002629

1061  39.  McRae AF, Powell JE, Henders AK, Bowdler L, Hemani G, Shah S, et al.
1062       Contribution of genetic variation to transgenerational inheritance of DNA
1063       methylation. Genome Biol. 2014;15: R73. doi:10.1186/gb-2014-15-5-r73

1064  40.  Novembre J, Johnson T, Bryc K, Kutalik Z, Boyko AR, Auton A, et al. Genes
1065       mirror geography within Europe. Nature. 2008;456: 98–101.
1066       doi:10.1038/nature07566

1067  41.  Price AL, Patterson NJ, Plenge RM, Weinblatt ME, Shadick N, Reich D.
1068       Principal components analysis corrects for stratification in genome-wide
1069       association studies. Nat Genet. 2006;38: 904–909. doi:10.1038/ng1847

1070  42.  Yu J, Pressoir G, Briggs WH, Vroh Bi I, Yamasaki M, Doebley JF, et al. A
1071       unified mixed-model method for association mapping that accounts for multiple
1072       levels of relatedness. Nat Genet. 2006;38: 203–8. doi:10.1038/ng1702

1073    43.    Tung J, Zhou X, Alberts SC, Stephens M, Gilad Y. The genetic architecture of
1074           gene expression levels in wild baboons. eLife. 2015;4: 1–22.
1075           doi:10.7554/eLife.04729

1076    44.    Turner L, Harr B. Genome-wide mapping in a house mouse hybrid zone
1077           reveals hybrid sterility loci and Dobzhansky-Muller interactions. eLife. 2014;3:
1078           e02504. doi:10.7554/eLife.02504

1079    45.    Tung J, Barreiro LB, Johnson ZP, Hansen KD, Michopoulos V, Toufexis D, et
1080           al. Social environment is associated with gene regulatory variation in the
1081           rhesus macaque immune system. Proc Natl Acad Sci. 2012;109: 6490–5.
1082           doi:10.1073/pnas.1202734109

1083    46.    Orozco LD, Morselli M, Rubbi L, Guo W, Go J, Shi H, et al. Epigenome-Wide
1084           Association of Liver Methylation Patterns and Complex Metabolic Traits in
1085           Mice. Cell Metab. Elsevier Inc.; 2015;21: 905–917.
1086           doi:10.1016/j.cmet.2015.04.025

1087    47.    Zhao K, Lu Z-X, Park JW, Zhou Q, Xing Y. GLiMMPS: Robust statistical model
1088           for regulatory variation of alternative splicing using RNA-seq data. Genome
1089           Biol. BioMed Central Ltd; 2013;14: R74. doi:10.1186/gb-2013-14-7-r74

1090    48.    Schmitz RJ, Schultz MD, Urich M, Nery JR, Pelizzola M, Libiger O, et al.
1091           Patterns of population epigenomic diversity. Nature. 2013;
1092           doi:10.1038/nature11968

1093    49.    Platt A, Gugger PF, Pellegrini M, Sork VL. Genome-wide signature of local
1094           adaptation linked to variable CpG methylation in oak populations. Mol Ecol.
1095           2015;1: n/a–n/a. doi:10.1111/mec.13230

1096    50.    Heyn H, Moran S, Hernando-herraez I, Res G, Sayols S, Gomez A, et al. DNA
1097           methylation contributes to natural human variation DNA methylation
1098           contributes to natural human variation. 2013; 1363–1372.
1099           doi:10.1101/gr.154187.112

1100    51.    Eichten SR, Briskine R, Song J, Li Q, Swanson-Wagner R, Hermanson PJ, et
1101           al. Epigenetic and genetic influences on DNA methylation variation in maize
1102           populations. Plant Cell. 2013;25: 2783–97. doi:10.1105/tpc.113.114793

1103    52.    Gertz J, Varley KE, Reddy TE, Bowling KM, Pauli F, Parker SL, et al. Analysis
1104           of DNA methylation in a three-generation family reveals widespread genetic
1105           influence on epigenetic regulation. PLoS Genet. 2011;7: e1002228.
1106           doi:10.1371/journal.pgen.1002228

53.  Tobi EW, Goeman JJ, Monajemi R, Gu H, Putter H, Zhang Y, et al. DNA methylation signatures link prenatal famine exposure to growth and metabolism. Nat Commun. 2014;5: 1–13. doi:10.1038/ncomms6592

54.  Wong CCY, Caspi A, Williams B, Craig IW, Houts R, Ambler A, et al. A longitudinal study of epigenetic variation in twins. Epigenetics. 2010;5: 516–526. doi:10.4161/epi.5.6.12226

55.  Gordon L, Joo JE, Powell JE, Ollikainen M, Novakovic B, Li X, et al. Neonatal DNA methylation profile in human twins is specified by a complex interplay between intrauterine environmental and genetic factors, subject to tissue-specific influence. Genome Res. 2012; doi:10.1101/gr.136598.111

56.  Zhou X, Carbonetto P, Stephens M. Polygenic Modeling with Bayesian Sparse Linear Mixed Models. PLoS Genet. 2013;9. doi:10.1371/journal.pgen.1003264

57.  Robinson MD, McCarthy DJ, Smyth GK. edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. Bioinformatics. 2010;26: 139–40. doi:10.1093/bioinformatics/btp616

58.  Anders S, Huber W. Differential expression analysis for sequence count data. Genome Biol. BioMed Central Ltd; 2010;11: R106. doi:10.1186/gb-2010-11-10-r106

59.  Rapaport F, Khanin R, Liang Y, Pirun M, Krek A, Zumbo P, et al. Comprehensive evaluation of differential gene expression analysis methods for RNA-seq data. Genome Biol. 2013;14: R95. doi:10.1186/gb-2013-14-9-r95

60.  Robinson MD, Smyth GK. Moderated statistical tests for assessing differences in tag abundance. Bioinformatics. 2007;23: 2881–2887. doi:10.1093/bioinformatics/btm453

61.  Pickrell JK, Marioni JC, Pai AA, Degner JF, Engelhardt BE, Nkadori E, et al. Understanding mechanisms underlying human gene expression variation with {RNA} sequencing. Nature. 2010;464: 768–772.

62.  Knowles DA, Davis JR, Raj A, Zhu X, Potash JB, Myrna M, et al. Allele-specific expression reveals interactions between genetic variation and environment. bioRxiv. 2015; doi:http://dx.doi.org/10.1101/025874

63.  McCulloch CE, Searle SR, Neuhaus JM. Generalized, Linear, and Mixed Models. New York, NY, USA: Wiley-Interscience; 2008.

64.  Bolker BMB, Brooks MEM, Clark CJ, Geange SW, Poulsen JR, Stevens MHH, et al. Generalized linear mixed models: a practical guide for ecology and

1141    evolution. Trends Ecol Evol. 2009;24: 127–135.
1142    doi:10.1016/j.tree.2008.10.008

1143    65.    Pinheiro JC, Chao EC. Efficient Laplacian and Adaptive Gaussian Quadrature
1144           Algorithms for Multilevel Generalized Linear Mixed Models. J Comput Graph
1145           Stat. 2006;15: 58–81. doi:10.1198/106186006X96962

1146    66.    Breslow NE, Clayton DG. Approximate inference in generalized linear mixed
1147           models. J Am Stat Assoc. 1993;88: 9–25.

1148    67.    Goldstein H. Nonlinear multilevel models, with an application to discrete
1149           response data. Biometrika. 1991;78: 45–51.

1150    68.    Goldstein H, Rasbash J. Improved approximations for multilevel models with
1151           binary responses. J R Stat Soc Ser A. 1996;159: 505–513.

1152    69.    Rodriguez G, Goldman N. Improved estimation procedures for multilevel
1153           models with binary response: {A} case-study. J R Stat Soc Ser A. 2001;164:
1154           339–355.

1155    70.    Browne WJ, Draper D. A comparison of {B}ayesian and likelihood-based
1156           methods for fitting multilevel models. Bayesian Anal. 2006;3: 473–514.

1157    71.    Jang W, Lim J. A numerical study of {PQL} estimation biases in generalized
1158           linear mixed models under heterogeneity of random effects. Commun Stat -
1159           Simul Comput. 2009;38: 692–702.

1160    72.    Fong Y, Rue H, Wakefield J. Bayesian inference for generalized linear mixed
1161           models. Biostatistics. 2010;11: 397–412.

1162    73.    Schwartz L. On Bayes procedures. Z Wahrscheinlichkeitstheorie. 1965;4: 10–
1163           26. doi:10.1007/BF00535479

1164    74.    Frühwirth-Schnatter S, Frühwirth R, Held L, Rue H. Improved auxiliary mixture
1165           sampling for hierarchical models of non-Gaussian data. Stat Comput. 2009;19:
1166           479–492. doi:10.1007/s11222-008-9109-4

1167    75.    Scott SL. Data augmentation, frequentist estimation, and the Bayesian
1168           analysis of multinomial logit models. Stat Pap. 2011;52: 87–109.
1169           doi:10.1007/s00362-009-0205-0

1170    76.    Fruhwirth-Schnatter S, Fruhwirth R. Data augmentation and MCMC for binary
1171           and multinomial logit models. In: Kneib T, Tutz G, editors. Statistical Modelling
1172           and Regression Structures: Festschrift in Honour of Ludwig Fahrmeir. New
1173           York: Springer; 2010. pp. 111–132. doi:10.1007/978-3-7908-2413-1

1174    77.    Pirinen M, Donnelly P, Spencer CC. Efficient computation with a linear mixed
1175           model on large-scale data sets with applications to genetic studies. Ann Appl
1176           Stat. 2013;7: 369–390. doi:10.1214/12-AOAS586

1177    78.    Hadfield JD. MCMC methods for multi-response generalized linear mixed
1178           models: the MCMCglmm R package. J Stat Softw. 2010;33: 1–22.

1179    79.    Landau DA, Clement K, Ziller MJ, Boyle P, Fan J, Gu H, et al. Locally
1180           Disordered Methylation Forms the Basis of Intratumor Methylome Variation in
1181           Chronic Lymphocytic Leukemia. Cancer Cell. Elsevier Inc.; 2014;26: 813–825.
1182           doi:10.1016/j.ccell.2014.10.012

1183    80.    Plongthongkum N, van Eijk KR, de Jong S, Wang T, Sul JH, Boks MPM, et al.
1184           Characterization of Genome-Methylome Interactions in 22 Nuclear Pedigrees.
1185           PLoS One. 2014;9: e99313. doi:10.1371/journal.pone.0099313

1186    81.    Ziller MJ, Müller F, Liao J, Zhang Y, Gu H, Bock C, et al. Genomic distribution
1187           and Inter-Sample variation of Non-CpG methylation across human cell types.
1188           PLoS Genet. 2011;7. doi:10.1371/journal.pgen.1002389

1189    82.    Becker C, Hagmann J, Müller J, Koenig D, Stegle O, Borgwardt K, et al.
1190           Spontaneous epigenetic variation in the Arabidopsis thaliana methylome.
1191           Nature. 2011;480: 245–9. doi:10.1038/nature10555

1192    83.    Carone BR, Fauquier L, Habib N, Shea JM, Hart CE, Li R, et al. Paternally
1193           induced transgenerational environmental reprogramming of metabolic gene
1194           expression in mammals. Cell. Elsevier Inc.; 2010;143: 1084–96.
1195           doi:10.1016/j.cell.2010.12.008

1196    84.    Murria R, Palanca S, Juan I De, Egoavil C, Alenda C, García-casado Z, et al.
1197           Methylation of tumor suppressor genes is related with copy number
1198           aberrations in breast cancer. Am J Cancer Res. 2015;5: 375–385.

1199    85.    Lockett G a., Kucharski R, Maleszka R. DNA methylation changes elicited by
1200           social stimuli in the brains of worker honey bees. Genes, Brain Behav.
1201           2012;11: 235–242. doi:10.1111/j.1601-183X.2011.00751.x

1202    86.    Atwell S, Huang YS, Vilhjálmsson BJ, Willems G, Horton M, Li Y, et al.
1203           Genome-wide association study of 107 phenotypes in Arabidopsis thaliana
1204           inbred lines. Nature. 2010;465: 627–631. doi:10.1038/nature08800

1205    87.    Horton MW, Hancock AM, Huang YS, Toomajian C, Atwell S, Auton A, et al.
1206           Genome-wide patterns of genetic variation in worldwide Arabidopsis thaliana
1207           accessions from the RegMap panel. Nat Genet. 2012;44: 212–216.
1208           doi:10.1038/ng.1042

1209 88. Cadman CSC, Toorop PE, Hilhorst HWM, Finch-Savage WE. Gene
1210      expression profiles of Arabidopsis Cvi seeds during dormancy cycling indicate
1211      a common underlying dormancy control mechanism. Plant J. 2006;46: 805–
1212      822. doi:10.1111/j.1365-313X.2006.02738.x

1213 89. Price AL, Zaitlen N a, Reich D, Patterson N. New approaches to population
1214      stratification in genome-wide association studies. Nat Rev Genet. 2010;11:
1215      459–463. doi:10.1038/nrg2813

1216 90. Altmann J, Alberts S, Haines S, Dubach J, Muruthi PM, Coote T, et al.
1217      Behavior predicts genetic structure in a wild primate group. Proc Natl Acad
1218      Sci. 1996;93: 5797–5801.

1219 91. Winnefeld M, Lyko F. The aging epigenome: DNA methylation from the cradle
1220      to the grave. Genome Biol. 2012;13: 165. doi:10.1186/gb4033

1221 92. Day K, Waite LL, Thalacker-Mercer A, West A, Bamman MM, Brooks JD, et al.
1222      Differential DNA methylation with age displays both common and dynamic
1223      features across human tissues that are influenced by CpG landscape.
1224      Genome Biol. 2013;14: R102. doi:10.1186/gb-2013-14-9-r102

1225 93. Christensen BC, Houseman EA, Marsit CJ, Zheng S, Wrensch MR, Wiemels
1226      JL, et al. Aging and environmental exposures alter tissue-specific DNA
1227      methylation dependent upon CPG island context. PLoS Genet. 2009;5.
1228      doi:10.1371/journal.pgen.1000602

1229 94. Rakyan VK, Down TA, Maslau S, Andrew T, Yang T, Beyan H, et al. Human
1230      aging-associated DNA hypermethylation occurs preferentially at bivalent
1231      chromatin domains. Genome Res. 2010;4: 434–439.

1232 95. Hansen KD, Timp W, Bravo HC, Sabunciyan S, Langmead B, McDonald OG,
1233      et al. Increased methylation variation in epigenetic domains across cancer
1234      types. Nat Genet. Nature Publishing Group; 2011;43: 768–75.
1235      doi:10.1038/ng.865

1236 96. Roadmap Epigenomics Consortium, Kundaje A, Meuleman W, Ernst J, Bilenky
1237      M, Yen A, et al. Integrative analysis of 111 reference human epigenomes.
1238      Nature. 2015;518: 317–330. doi:10.1038/nature14248

1239 97. Dunham I, Kundaje A, Aldred SF, Collins PJ, Davis C, Doyle F, et al. An
1240      integrated encyclopedia of DNA elements in the human genome. Nature.
1241      2012;489: 57–74. doi:10.1038/nature11247

1242 98. Murgatroyd C, Patchev A V, Wu Y, Micale V, Bockmühl Y, Fischer D, et al.
1243      Dynamic DNA methylation programs persistent adverse effects of early-life
1244      stress. Nat Neurosci. 2009;12: 1559–66. doi:10.1038/nn.2436

1245  99.  Ikegame T, Bundo M, Murata Y, Kasai K, Kato T, Iwamoto K. DNA methylation
1246       of the BDNF gene and its relevance to psychiatric disorders. J Hum Genet.
1247       2013;58: 434–8. doi:10.1038/jhg.2013.65

1248  100. Elliott E, Ezra-Nevo G, Regev L, Neufeld-Cohen A, Chen A. Resilience to
1249       social stress coincides with functional DNA methylation of the CRF gene in
1250       adult mice. Nat Neurosci. 2010;13: 1351–3. doi:10.1038/nn.2642

1251  101. Lam LL, Emberly E, Fraser HB, Neumann SM, Chen E, Miller GE, et al.
1252       Factors underlying variable DNA methylation in a human community cohort.
1253       Proc Natl Acad Sci. 2012;109: 17253–60. doi:10.1073/pnas.1121249109

1254  102. Feil R, Fraga MF. Epigenetics and the environment: emerging patterns and
1255       implications. Nat Rev Genet. 2011;13: 97–109. doi:10.1038/nrg3142

1256  103. Shi J, Marconett CN, Duan J, Hyland PL, Li P, Wang Z, et al. Characterizing
1257       the genetic basis of methylome diversity in histologically normal human lung
1258       tissue. Nat Commun. 2014;5: 3365. doi:10.1038/ncomms4365

1259  104. The International HapMap Consortium. The International HapMap Project.
1260       Nature. 2003;426: 789–796. doi:10.1038/nature02168

1261  105. Cann H, Toma D, Cazes L, Legrand M, Morel V, Piouffre L, et al. A human
1262       genome diversity cell line panel. Science. 2002;296: 261–2.
1263       doi:http://dx.doi.org/10.1108/17506200710779521

1264  106. The 1000 Genomes Project Consortium. An integrated map of genetic
1265       variation from 1,092 human genomes. Nature. 2012;135: 0–9.
1266       doi:10.1038/nature11632

1267  107. Bennett BJ, Farber CR, Orozco L, Kang HM, Ghazalpour A, Siemers N, et al.
1268       A high-resolution association mapping panel for the dissection of complex
1269       traits in mice. Genome Res. 2010; 281–290. doi:10.1101/gr.099234.109

1270  108. Weigel D, Mott R. The 1001 genomes project for Arabidopsis thaliana.
1271       Genome Biol. 2009;10: 107. doi:10.1186/gb-2009-10-5-107

1272  109. Quon G, Lippert C, Heckerman D, Listgarten J. Patterns of methylation
1273       heritability in a genome-wide analysis of four brain regions. Nucleic Acids Res.
1274       2013;41: 2095–2104. doi:10.1093/nar/gks1449

1275  110. Akalin A, Kormaksson M. methylKit: a comprehensive R package for the
1276       analysis of genome-wide DNA methylation profiles. Genome Biol. BioMed
1277       Central Ltd; 2012;13: R87. doi:10.1186/gb-2012-13-10-R87

1278   111.   Hansen K, Langmead B, Irizarry R. BSmooth□: from whole genome bisulfite
1279          sequencing reads to differentially methylated regions. Genome Biol. BioMed
1280          Central Ltd; 2012;13: R83. doi:10.1186/gb-2012-13-10-R83

1281   112.   Battle A, Mostafavi S, Zhu X, Potash JB, Weissman MM, McCormick C, et al.
1282          Characterizing the genetic basis of transcriptome diversity through RNA-
1283          sequencing of 922 individuals. Genome Res. 2014;24: 14–24.
1284          doi:10.1101/gr.155192.113

1285   113.   Crowley JJ, Zhabotynsky V, Sun W, Huang S, Pakatci IK, Kim Y, et al.
1286          Analyses of allele-specific gene expression in highly divergent mouse crosses
1287          identifies pervasive allelic imbalance. Nat Genet. 2015;47: 353–360.
1288          doi:10.1038/ng.3222

1289   114.   Pickrell JJK, Marioni JJC, Pai AA, Degner JF, Engelhardt BE, Nkadori E, et al.
1290          Understanding mechanisms underlying human gene expression variation with
1291          RNA sequencing. Nature. 2010;464: 768–772.
1292          doi:10.1038/nature08872.Understanding

1293   115.   Skelly D, Johansson M, Madeoy J, Wakefield J, Akey JM. A powerful and
1294          flexible statistical framework for testing hypotheses of allele-specific gene
1295          expression from RNA-seq data. Genome Res. 2011;21: 1728–1737.
1296          doi:10.1101/gr.119784.110

1297   116.   Harvey C, Moyebrailean G, Davis O, Wen X, Luca F, Pique-Regi R. QuASAR:
1298          Quantitative allele specific analysis of reads. Bioinformatics. 2014; 1–7.
1299          doi:10.1101/007492

1300   117.   Grundberg E, Small KS, Hedman ÅK, Nica AC, Buil A, Keildson S, et al.
1301          Mapping cis- and trans-regulatory effects across multiple tissues in twins. Nat
1302          Genet. 2012;44: 1084–1089. doi:10.1038/ng.2394

1303   118.   Zou J, Lippert C, Heckerman D, Aryee M, Listgarten J. Epigenome-wide
1304          association studies without the need for cell-type composition. Nat Methods.
1305          2014;11: 309–11. doi:10.1038/nmeth.2815

1306   119.   Alberts SC, Altmann J. The Amboseli Baboon Research Project: 40 years of
1307          continuity and change. In: Kappeler P, Watts DP, editors. Long-Term Field
1308          Studies of Primates. New York: Springer; 2012. pp. 261–288.

1309   120.   Altmann J, Altmann S, Hausfater G. Physical maturation and age estimates of
1310          yellow baboons, Papio cynocephalus, in Amboseli National Park, Kenya. Am J
1311          Primatol. 1981;1: 389–399. doi:10.1002/ajp.1350010404

1312   121.   Buchan JC, Alberts SC, Silk JB, Altmann J. True paternal care in a multi-male
1313          primate society. Nature. 2003;425: 179–81. doi:10.1038/nature01866

1314     122.   Alberts SC, Buchan JC, Altmann J. Sexual selection in wild baboons: from
1315            mating opportunities to paternity success. Anim Behav. 2006;72: 1177–1196.
1316            doi:10.1016/j.anbehav.2006.05.001

1317     123.   Lynch M, Ritland K. Estimation of pairwise relatedness with molecular
1318            markers. Genetics. 1999;152: 1753–1766.

1319     124.   Wang J. COANCESTRY: a program for simulating, estimating and analysing
1320            relatedness and inbreeding coefficients. Mol Ecol Resour. 2011;11: 141–5.
1321            doi:10.1111/j.1755-0998.2010.02885.x

1322     125.   Tung J, Primus A, Bouley AJ, Severson TF, Alberts SC, Wray G. Evolution of
1323            a malaria resistance gene in wild primates. Nature. 2009;460: 388–91.
1324            doi:10.1038/nature08149

1325     126.   Tung J, Akinyi MY, Mutura S, Altmann J, Wray G, Alberts SC. Allele-specific
1326            gene expression in a wild nonhuman primate population. Mol Ecol. 2011;20:
1327            725–39. doi:10.1111/j.1365-294X.2010.04970.x

1328     127.   Bolger AM, Lohse M, Usadel B. Trimmomatic: A flexible trimmer for Illumina
1329            sequence data. Bioinformatics. 2014;30: 2114–2120.
1330            doi:10.1093/bioinformatics/btu170

1331     128.   Xi Y, Li W. BSMAP: whole genome bisulfite sequence MAPping program.
1332            BMC Bioinformatics. 2009;10: 232. doi:10.1186/1471-2105-10-232

1333     129.   Venables WN, Ripley BD. Modern Applied Statistics with S. Fourth. New York,
1334            NY: Springer; 2002.

1335     130.   Hastie T, Tibshirani R, Narasimhan B, Chu G. Impute: imputation for
1336            microarray data. R package version 1.42.0. 2015.

1337     131.   Johnson KC, Koestler DC, Cheng C, Christensen BC. Age-related DNA
1338            methylation in normal breast tissue and its relationship with invasive breast
1339            tumor methylation. Epigenetics. 2014;9: 268–275. doi:10.4161/epi.27015

1340     132.   Law C, Chen Y, Shi W, Smyth G. Voom! Precision weights unlock linear model
1341            analysis tools for RNA-seq read counts. Melbourne, Australia; 2013.

1342     133.   Cunningham F, Amode MR, Barrell D, Beal K, Billis K, Brent S, et al. Ensembl
1343            2015. Nucleic Acids Res. 2014;43: D662–D669. doi:10.1093/nar/gku1010

1344     134.   Karolchik D, Barber GP, Casper J, Clawson H, Cline MS, Diekhans M, et al.
1345            The UCSC Genome Browser database: 2014 update. Nucleic Acids Res.
1346            2014;42: 764–770. doi:10.1093/nar/gkt1168

1347  135.  Hernando-Herraez I, Prado-Martinez J, Garg P, Fernandez-Callejo M, Heyn H,
1348        Hvilsom C, et al. Dynamics of DNA methylation in recent human and great ape
1349        evolution. PLoS Genet. 2013;9: e1003763. doi:10.1371/journal.pgen.1003763

1350  136.  Rönn T, Volkov P, Davegårdh C, Dayeh T, Hall E, Olsson AH, et al. A six
1351        months exercise intervention influences the genome-wide DNA methylation
1352        pattern in human adipose tissue. PLoS Genet. 2013;9: e1003572.
1353        doi:10.1371/journal.pgen.1003572

1354  137.  Slieker RC, Bos SD, Goeman JJ, Bovée JV, Talens RP, van der Breggen R, et
1355        al. Identification and systematic annotation of tissue-specific differentially
1356        methylated regions using the Illumina 450k array. Epigenetics Chromatin.
1357        2013;6: 26. doi:10.1186/1756-8935-6-26

**Supplementary Figure 1. In a real WGBS dataset (from *Arabidopsis*) and a real RRBS dataset (from yellow baboons), coverage varies widely across CpG sites and individuals.** For each CpG site represented in each data set (n=433,871 for baboon and n=830,676 for *Arabidopsis*), we calculated the mean site-specific coverage across individuals, as well as the standard deviation of coverage values for those sites. The distribution of these values are are shown for the baboon RRBS dataset (A-B, in blue) and the *Arabidopis* WGBS dataset (C-D, in green). Average coverage values are depicted in A and C, and coverage standard deviation values are depicted in B and D.

**Supplementary Figure 2. MACAU p-values are consistent across runs**. QQ-plots comparing the p-value distributions for 3 independent runs of MACAU on the same data sets, with different simulated heritability values (Panels A, D - $h^2 = 0$; Panels B, E - $h^2 = 0.3$; Panels C, F - $h^2 = 0.6$). Pairwise correlations between each independent run were $R > 0.95$ for $h^2 = 0$:,$R > 0.97$ for $h^2 = 0.3$; and $R > 0.98$ for $h^2 = 0.6$. Distributions shown are for analyses of simulated secondary dormancy effects on DNA methylation levels in the *Arabidopsis* data set (4000 sites, n=24 accessions).

**Supplementary Figure 3. MACAU results are robust to prior perturbation.** QQ-plots comparing the results from MACAU implemented with an uninformative prior ($\sigma^2 \sim U(0,1)$, as in the main text, x-axis) versus an alternative prior ($\log(\sigma^2) \sim U(0,1)$, y-axis). All analyses tested for age effects on DNA methylation levels in a simulated baboon data (based on properties of the real baboon RRBS data and age information). Sample sizes and heritabilities are shown on each plot, as are the results from a Kolmogorov-Smirnov test comparing the two distributions represented in each plot. In all cases, the simulated percent variance explained by age was set to 10%. The number of age-associated sites detected in each analysis were identical for all simulations where n=80 (10% empirical FDR), and very similar when n=50 (0.4-0.8% more age-associated sites were detected with the alternative prior than with the uninformative prior).

**Supplementary Figure 4. The normal mixture provides an accurate approximation to the negative log gamma distribution. (**A) Density plot and (B) quantile-quantile plots demonstrating that the normal mixture approximation approximates $-\log(Ga(r, 1))$ well even in the most difficult case when r=1.

**Supplementary Figure 5. A binomial mixed model (BMM) implemented in MACAU is more efficent than a BMM implemented in the software MCMCglmm.** (A) Computation time (in hours) is plotted for datasets containing varying numbers of individuals, but each containing 100 sites. Computation time is plotted on a log10 scale in the main plot, and on a traditional scale in the inset. (B) Computation time (in hours) is plotted for a dataset containing 150 individuals, but varying numbers of sites (in thousands) as noted on the x-axis. All computation was performed on a single core of an Intel Xeon L5420 2.50 GHz processor.

1402

1403 **Supplementary Figure 6. Comparisons between methods when DNA**
1404 **methylation levels are not heritable, and the predictor variable is binarized.** To
1405 include methods that can only analyze categorical differences in DNA methylation
1406 levels between two groups, we binarized age values in our simulated RRBS
1407 datasets (individuals below median age = young versus individuals above median
1408 age = old). We compared the AUC of each method (open circles), as well as their
1409 ability to detect true positives at a 10% FDR (closed circles). For these comparisons,
1410 we used simulated datasets with a fixed $h^2$ of 0 (n = 5000 sites including 500 true
1411 positives and 4500 true negatives; percent variance explained by age varies as
1412 noted in the panel headings). Results for simulations with (A) n = 50 or (B) n = 80
1413 individuals are plotted below.  Note that the right-hand y axis for the proportion of
1414 true positives detected varies depending on sample size.

1415 **Supplementary Figure 7. Comparisons between methods when DNA**
1416 **methylation levels are heritable, and the predictor variable is binarized.** To
1417 include methods that can only analyze categorical differences in DNA methylation
1418 levels between two groups, we binarized age values in our simulated RRBS
1419 datasets (individuals below median age = young versus individuals above median
1420 age = old). We compared the AUC of each method (open circles), as well as their
1421 ability to detect true positives at a 10% FDR (closed circles). For these comparisons,
1422 we used simulated datasets with a fixed sample size of 80 (n = 5000 sites including
1423 500 true positives and 4500 true negatives; percent variance explained by age
1424 varies as noted in the panel headings). Results for simulations with (A) $h^2$ = 0.3 or
1425 (B) $h^2$ = 0.6 are plotted below.

1426

1427 **Supplementary Figure 8. MACAU outperforms linear mixed models**
1428 **implemented after a variety of standard data transformation approaches.** We
1429 performed four different transformations on simulated baboon bisulfite sequencing
1430 count data (n = 5000 sites including 500 true positives and 4500 true negatives;
1431 percent variance explained by age = 10%; sample size = 80, $h^2$ = 0.6). Below, we
1432 use QQ-plots to compare the distribution of p-values produced by GEMMA
1433 (operating on the transformed data) versus MACAU (analyzing the raw count data).
1434 In all panels, the observed p-values are plotted against quantiles for the distribution
1435 of p-values obtained from running each method (MACAU or GEMMA, respectively)
1436 on permuted data. We also note the proportion of simulated true positives detected
1437 by each approach (for comparison, MACAU detects 20.6% of simulated true
1438 positives in the same dataset).

1439 **Supplementary Figure 9. Comparison across methods when DNA methylation**
1440 **levels are not heritable.** We compared the AUC of each method (open circles) and
1441 their ability to detect true positives at a 10% FDR (closed circles). We did so using
1442 simulated data sets (n = 5000 sites including 500 true positives and 4500 true
1443 negatives; percent variance explained by age varies as noted in the panel
1444 headings). For all simulations shown below, $h^2$ was set to 0. (A) Results for

1445  simulations with n=20 individuals; (B) with n=50 individuals; and (C) with n=80
1446  individuals. Note that the right-hand y axis for the proportion of true positives
1447  detected varies depending on sample size.

1448  **Supplementary Figure 10. Secondary dormancy is correlated with population**
1449  **structure in the *Arabidopsis* WGBS dataset.** Principal components analysis on
1450  188,093 genotyped sites with minor allele frequency >5% reveals that genetic
1451  background is correlated with secondary dormancy values. The correlation between
1452  the secondary dormancy phenotype values and the first principal component of the
1453  genetic relatedness matrix is $R^2 = 0.38$, $p = 7.84 \times 10^{-4}$ (n = 24). The first principal
1454  component (PC1) explains 8.5% of the genetic variance in the data set.
1455
1456  **Supplementary Figure 11. A mixed modeling approach (implemented in**
1457  **GEMMA) appropriately controls for genetic covariance in simulated and real**
1458  **WGBS data.** (A, B) The distribution of *p*-values for 4000 simulated true negative
1459  sites (n = 24 accessions; effect of secondary dormancy on DNA methylation levels =
1460  0). For each simulation, $h^2$ was set to 0.3 (A) or 0.6 (B). Simulated data were
1461  analyzed with a linear model or GEMMA, and compared against the expected
1462  uniform distribution. (C) QQ-plots comparing the *p*-value distributions for (i) a model
1463  testing for effects of secondary dormancy on DNA methylation levels in real WGBS
1464  data, plotted on the y-axis; and (ii) the same model when the secondary dormancy
1465  values were permuted across individuals, plotted on the x-axis. Here, the lack of
1466  inflated test statistics in the case of the linear model is likely due to the model's low
1467  power (see Figure S12b, for n=25). The genomic control factor, λ, is shown for each
1468  set of results.
1469
1470  **Supplementary Figure 12. A mixed modeling approach (implemented in**
1471  **GEMMA) controls for genetic covariance in data sets that span a range of**
1472  **sample sizes and levels of correlation between population structure and a**
1473  **predictor variable of interest.** Genomic control factor when simulated datasets
1474  (n=5000 sites per dataset; $h^2 = 0.6$) were analyzed with either (A) a linear model or
1475  (B) a linear mixed model implemented in GEMMA. The correlation between the
1476  simulated predictor variable and the first principal component of genome-wide
1477  genotype data is plotted on the x-axis.
1478
1479  **Supplementary Figure 13. MACAU controls for genetic covariance in data sets**
1480  **that span a range of sample sizes and levels of correlation between population**
1481  **structure and a predictor variable of interest.** Percent of dataset associated with
1482  the predictor variable (at a 10% FDR) when simulated datasets (n=5000 sites per
1483  dataset; $h^2 = 0.6$) were analyzed with either (A) a beta-binomial model or (B) a
1484  binomial mixed model implemented in MACAU. The correlation between the
1485  simulated predictor variable and the first principal component of genome-wide
1486  genotype data is plotted on the x-axis.
1487
1488  **Supplementary Figure 14. Distribution of pairwise relatedness values for**
1489  **baboons (n=80) from a single social group, used in simulations.** Approximately

1490     half of the individuals are unrelated, while a small proportion (~10%) are highly
1491     related (i.e., related at the level of half siblings or higher, r = 0.25).

1492

1493     **Supplementary Figure 15. MACAU produces well-calibrated *p*-values when the**
1494     **simulated effect of age is set to 0.** Results from 4500 simulated sites, where we
1495     set the effect of age on DNA methylation levels equal to 0 and the heritability of DNA
1496     methylation levels equal to (A) 0, (B) 0.3, or (C) 0.6. All QQ-plots compare the
1497     distribution of p-values produced by MACAU to the expected uniform distribution.

1498     **Supplementary Figure 16. MACAU provides increased power to detect age-**
1499     **associated sites when DNA methylation levels are heritable.** We simulated age
1500     effects on DNA methylation levels, in presence of genetic effects (panel A, $h^2$ = 0.3;
1501     panel B, $h^2$ = 0.6) across a range of effect sizes. The proportion of true positives
1502     detected at a 10% empirical FDR is plotted for each method (closed circles) as is the
1503     AUC (open circles). For all simulations shown here, the sample size was set to 80
1504     individuals.

1505     **Supplementary Figure 17. MACAU provides increased power to detect age-**
1506     **associated sites when DNA methylation levels are heritable.** We simulated age
1507     effects on DNA methylation levels in datasets of 500 (A-B) and 1000 individuals (C-
1508     D). For all simulations, we included genetic effects on DNA methylation levels
1509     (panels A and C: $h^2$ = 0.3; panels B and D: $h^2$ = 0.6). Below, we show the proportion
1510     of true positives detected at a 1% empirical FDR (closed circles) as well as the AUC
1511     (open circles) for each method.

1512

1513     **Supplementary Figure 18. Distribution of sites covered in the baboon RRBS**
1514     **dataset (n = 433,871 CpG sites).** (A) Absolute number of sites analyzed for a given
1515     genomic region. See *Materials and Methods* for information on how we defined each
1516     genomic region. (B) Proportion of total annotated features in the baboon genome for
1517     which a least one CpG site was analyzed in this data set.

1518     **Supplementary Figure 19. DNA methylation patterns in the baboon RRBS data**
1519     **.** (A) The distributions of bisulfite conversion rates (estimated from a spike-in sample
1520     of unmethylated lambda phage DNA) and proportions of reads starting or ending
1521     with an *Msp1* digest site, for each sample. (B) Barplots showing the distribution of
1522     DNA methylation levels by genomic compartment. As expected, CpG islands,
1523     H3K3me1-marked enhancers and promoters tend to be lowly methylated, while
1524     gene bodies and the background set of all sites analyzed tend to be
1525     hypermethylated. See [1] for similar results from a human RRBS dataset. (C) For
1526     each CpG site within 5000 bp of an annotated Ensembl TSS, we calculated the
1527     mean DNA methylation level at that site across all 50 baboons. These mean levels
1528     are plotted as a smoothed function of distance from the TSS, stratified by gene
1529     expression level quartiles obtained from baboon whole blood RNA-seq [2]. As
1530     expected, more highly methylated regions are associated with more lowly expressed
1531     genes. Only expressed genes were included.

1532 **Supplementary Figure 20. Distribution of p-values from four different methods**
1533 **for the real RRBS data.** QQ-plots comparing the p-value distributions for (i) a model
1534 testing for effects of age on DNA methylation levels in real RRBS data, plotted on
1535 the y-axis; and (ii) the same model when the age values were permuted across
1536 individuals, plotted on the x-axis. For each method, the number of sites detected at a
1537 10% FDR was as follows: Beta-binomial = 747, GEMMA = 205, Linear 324, MACAU
1538 = 1018.
1539
1540 **Supplementary Figure 21. MACAU detects differentially methylated regions in**
1541 **the baboon genome.** Using the criteria of [3], we detected 142 age-related DMRs.
1542 Two representative DMRs are plotted in panels A and B (location of DMR in panel A:
1543 Chr14, 908111-908168; and panel B: Chr 20: 996106-996139; see Table S5 for the
1544 locations of additional DMRs). To detect DMRs, baboon ages were binarized into
1545 two categories, based on whether an individual's age fell above or below the median
1546 age in our sample. Smoothed estimates of DNA methylation levels are shown for
1547 each age group, and the location of measured CpG sites are noted along the x-axis
1548 by black dots. Panel C shows the proportion of all identified DMRs that fell in a CpG
1549 island, CpG island shore, or both.

1550 **Supplementary Figure 22. Sites identified by MACAU are consistently enriched**
1551 **near genes identified as age-associated in the same population**. For each
1552 method below, we asked whether CpG sites that occur near age-associated genes
1553 (identified using RNA-seq data from [2]) were more likely to be differentially
1554 methylated with age compared to the background set of all CpG sites near genes
1555 (using a Fisher's exact test). We report the enrichment observed and show whether
1556 the p-value associated with the Fisher's exact test (FET) was below 0.05 (triangles).
1557 We repeated this analysis using a varying number of top CpG sites from each
1558 method, with the number for each analysis shown on the x-axis. Dotted vertical lines
1559 correspond to the number of sites detected by MACAU at a 10% empirical FDR (a
1560 more conservative approach), or at a 10% FDR calculated in the R package *qvalue*
1561 [4] (a less conservative approach).

1562 **Supplementary Figure 23. MACAU is robust to error in the estimation of**
1563 **pairwise genetic relatedness.** To understand how the performance of MACAU
1564 varies when there is error in the estimation of pairwise genetic relatedness, we
1565 added random error drawn from a normal distribution with mean 0 and standard
1566 deviation as shown on the x-axis. We then reran our analyses of simulated data sets
1567 with varying heritabilities (as shown in the figure legend inset) where n=80 and
1568 percent variance explained by age=10%. For each analysis, we show the number of
1569 simulated true positives detected by MACAU at a 10% empirical FDR (note that the
1570 results from our original analyses, with no error in the estimation of pairwise genetic
1571 relatedness, corresponds to the results for SD = 0 on the x-axis).

1572 **Supplementary Figure 24. MACAU outperforms the linear modeling approach**
1573 **implemented in 'voom + limma'.** We tested the performance of a commonly used
1574 method for modeling RNA-seq data: the combination of the *voom* function for data

60

1575 weighting with *limma*, a linear model approach [5]. We used simulated baboon
1576 bisulfite sequencing count data where the percent variance explained by age = 10%,
1577 sample size = 80, and $h^2$ = 0.6 (n = 5000 sites including 500 true positives and 4500
1578 true negatives). QQ-plots show the results for both the *voom* + *limma* approach
1579 (purple), as well as results from the same dataset using MACAU (orange). QQ-plots
1580 compare the p-value distributions for (i) a model testing for the effect of age on DNA
1581 methylation levels, plotted on the y-axis; and (ii) the same model when the age
1582 values were permuted across individuals, plotted on the x-axis (i.e., the null
1583 distribution of p-values). MACAU detects 20.6% of simulated true positives at a 10%
1584 FDR, while the *voom* + *limma* approach detects less than 1% of simulated true
1585 positives.

1587 **Table S1. Normal Mixture Approximations to -log(Ga(r, 1)) for r in [1, 5].** Normal
1588 mixture approximations to -log(Ga(r, 1)) for r in [1, 5]. A separate normal mixture
1589 distribution is used to approximate each negative log gamma distribution. The
1590 estimated parameters in the normal mixture distribution ensure that the Kullback-
1591 Leibler (KL) divergence between the two distributions is below $5x10^{-4}$. The
1592 parameters in the normal mixture distribution include the number of normal
1593 components (k), their weights (w), means (m) and variances ($\sigma^2$). Means and
1594 variances are shown in their standardized version, where $\Psi(r)$ denotes the
1595 diagamma function and $\Psi'(r)$ denotes the trigamma function.

1597 **Table S2. Normal Mixture Approximations to -log(Ga(r, 1)) for r in [6, 170].**
1598 Normal mixture approximations to -log(Ga(r, 1)) for r in [6, 170]. A separate normal
1599 mixture distribution is used to approximate each negative log gamma distribution.
1600 The estimated parameters in the normal mixture distribution ensure that the
1601 Kullback-Leibler (KL) divergence between the two distributions is below $5x10^{-4}$. The
1602 parameters in the normal mixture distribution include the number of normal
1603 components (k), their weights (w), means (m) and variances ($\sigma^2$), all of which are
1604 functions of r. Means and variances are shown in their standardized version, where
1605 $\Psi(r)$ denotes the diagamma function and $\Psi'(r)$ denotes the trigamma function.

1607 **Table S3. Computation times for each method on the two real datasets.**
1608 Computation was performed on a single core of an Intel Xeon L5420 2.50 GHz
1609 processor. *n* = number of individuals; *m* = number of sites.

1611 **Table S4. Baboon RRBS dataset sample characteristics and read mapping**
1612 **summary.**

1614 **Table S5. Locations of identified age-DMRs in the baboon genome.**

1615 **Supplementary References**

1616 1.    Wang J, Xia Y, Li L, Gong D, Yao Y, Luo H, et al. Double restriction-enzyme
1617       digestion improves the coverage and accuracy of genome-wide CpG

1618    methylation profiling by reduced representation bisulfite sequencing. BMC
1619    Genomics. 2013;14: 11. doi:10.1186/1471-2164-14-11

1620  2.   Tung J, Zhou X, Alberts SC, Stephens M, Gilad Y. The genetic architecture of
1621    gene expression levels in wild baboons. eLife. 2015;4: 1–22.
1622    doi:10.7554/eLife.04729

1623  3.   Slieker RC, Bos SD, Goeman JJ, Bovée JV, Talens RP, van der Breggen R, et
1624    al. Identification and systematic annotation of tissue-specific differentially
1625    methylated regions using the Illumina 450k array. Epigenetics Chromatin.
1626    2013;6: 26. doi:10.1186/1756-8935-6-26

1627  4.   Dabney A, Storey J. qvalue: Q-value estimation for false discovery rate
1628    control. R package version 1.43.0. 2015.

1629  5.   Law C, Chen Y, Shi W, Smyth G. Voom! Precision weights unlock linear model
1630    analysis tools for RNA-seq read counts. Melbourne, Australia; 2013.

1631

1632

1633