

The Multi-State Epigenetic Pacemaker enables the identification of combinations of factors that influence DNA methylation

Colin Farrell^{1,4}, Keshiv Tandon¹, Roberto Ferrari², Kalsuda Lapborisuth¹, Rahil Modi¹, Sagi Snir³, and Matteo Pellegrini^{1,4}

¹Dept. of Molecular, Cell and Developmental Biology;
University of California, Los Angeles, CA 90095, USA;;

²Dept. of Chemistry, Life Sciences and Environmental Sustainability, Laboratory of Molecular Cell Biology of the Epigenome (MCBE), University of Parma, Italy;

³Dept. of Evolutionary Biology, University of Haifa, Israel;

⁴Corresponding Authors; colinpfarrell@gmail.com, matteop@mcdb.ucla.edu

Abstract

Epigenetic clocks, DNA methylation based predictive models of chronological age, are often utilized to study aging associated biology. Despite their widespread use, these methods do not account for other factors that also contribute to the variability of DNA methylation data. For example, many CpG sites show strong sex-specific or cell type specific patterns that likely impact the predictions of epigenetic age. To overcome these limitations, we developed a multidimensional extension of the Epigenetic Pacemaker, the Multi-State Epigenetic Pacemaker (MSEPM). We show that the MSEPM is capable of accurately modeling multiple methylation associated factors simultaneously, while also providing site specific models that describe the per site relationship between methylation and these factors. We utilized the MSEPM with a large aggregate cohort of blood methylation data to construct models of the effects of age, sex and cell type heterogeneity on DNA methylation. We found that these models capture a large fraction of the variability at thousands of DNA methylation sites. Moreover, this approach allows us to identify sites that are primarily affected by aging and no other factors. An analysis of these sites reveals that those that lose methylation over time are enriched for CTCF transcription factor chip peaks, while those that gain methylation over time are associated with bivalent promoters of genes that are not expressed in blood. These observations suggest mechanisms that underlie age associated methylation changes and suggest that age associated increases in methylation may not have strong functional consequences on cell states. In conclusion, the MSEPM is capable of accurately modeling multiple methylation associated factors and the models produced can illuminate site specific combinations of factors that affect methylation dynamics.

1 Introduction

DNA methylation, the addition of a methyl group to the fifth carbon of the cytosine pyrimidine ring, is associated with the topological organization of the cellular genome, gene expression and the state of a cell. Within a population of cells the methylation pattern at certain sites can change

predictably with the age of the individual from which the cells are drawn. This predictable nature of DNA methylation has led to the development of accurate DNA methylation based predictive models for age and health, termed epigenetic clocks. The difference between the predicted and the expected epigenetic age given an individual's chronological age has been interpreted as a measure of age acceleration[**Horvath2018-ia**], and has been associated with mortality[**Perna2016-pi**, **Marioni2015-sn**] and other adverse health outcomes[**Dugue2018-ad**, **Huang2019-hf**, **Armstrong2017-vg**, **Chuang2017-nk**, **Horvath2015-af**].

However, epigenetic clocks suffer from several limitations that limit the interpretability of their predictions and the underlying mechanisms. Epigenetic clocks are generally trained by using penalized regression based methods that attempt to minimize the difference between the predicted and observed value of age. As a result, as the error between predicted and observed age is decreased, the associations between age acceleration and mortality disappears[**Zhang2019-br**]. Second generation epigenetic clocks attempt to resolve this issue by fitting a measure of human health, rather than age, and as a result these clocks are generally more sensitive to individual health status[**Lu2019-lg**, **Levine2018-en**, **Belsky2020-ha**]. However, while the response variable is modified in these clocks the method used to fit the clock is largely the same. Epigenetic clocks are generally trained using regularized regression models, where the likelihood is maximized by minimizing the difference between the observed and predicted response variable subject to the elastic net penalty, λ_1 and λ_2 . Methylation sites that increase model error and are influenced by other relevant factors such as smoking or obesity, may be discarded during model fitting, thus limiting the ability of this approach to account for the effects of these extraneous factors on epigenetic aging.

As an alternative to penalized regression based methods we previously developed an evolutionary based model for epigenetic dynamics, the Epigenetic Pacemaker (EPM)[**Farrell2020-bn**, **Snir2016-dv**]. The EPM attempts to minimize the difference between observed and predicted methylation values amongst a collection of sites through the implementation of a conditional expectation maximization algorithm[**Snir2020-tc**]. Under the EPM the observed methylation status of a collection of sites is modeled linearly with respect to an input factor of interest, such as age. A hidden epigenetic state, that is related to the initial factor, but not necessarily linearly, is learned through the course of model fitting. The EPM can capture the non-linear relationship between methylation and age[**Snir2019-ii**] and outputs an interpretable model for each site. However, both the EPM and regression based methods suffer from the same limitation, which is that they are limited to a single trait predicted by, or used to model, observed methylation patterns. In reality, the observed methylation landscape is likely impacted by a variety of factors that act simultaneously to produce the observed methylome of an individual.

To overcome this limitation, we have developed a multidimensional extension of the EPM, the Multi-State Epigenetic Pacemaker (MSEPM). We show that the (MSEPM) can accurately model site specific methylation variation driven by several factors, and given a trained model, accurately predict the values of the factors associated with an individual's

84 observed methylation profile in both simulated methylation datasets and
 85 a large aggregate blood tissue methylation dataset. Importantly, as factors
 86 that explain the observed methylation profile of an individual are added
 87 to the model the ability to model the factors and methylation
 88 values improves. Additionally, we show that sites with similar associa-
 89 tions to modeled factors cluster together and are enriched for specific
 90 transcription factors. Therefore, unlike traditional epigenetic clocks, the
 91 MSEPM allows us to study mechanisms that may underlie age associated
 92 methylation changes. In our large dataset of blood samples, we find that
 93 sites that increase methylation with age are enriched for bivalent promoters,
 94 and are proximal to genes that are lowly expressed in blood. These
 95 results suggest that positively age associated sites may not have a sig-
 96 nificant functional impact on aging traits. The MSEPM is available as
 97 a Python package with scikit-learn style syntax under a MIT license at
 98 <https://github.com/NuttyLogic/MultistateEpigeneticPacemaker>.

99 2 Methods

100 2.1 Multi-State Epigenetic Pacemaker Model

The MSEPM model describes the observed methylation at site i and for individual j , $\hat{m}_{i,j}$, as a weighted linear combination of k individual epigenetic factors $p_{k,j}$.

$$\hat{m}_{i,j} = r_i^0 + \sum_{k=1}^n p_{j,k} r_{i,k}$$

101 Where k epigenetic factors are weighted by k site specific epigenetic rates
 102 of change, $r_{i,k}$, and offset by a sites specific intercept term, r_i^0 . Site
 103 parameters, $r_{i,k}$ and r_i^0 , are characteristic of the site and shared amongst all
 104 individuals while epigenetic factors, $p_{j,k}$, are characteristic of an individ-
 105 ual and are the same across all sites for that individual. In practice, the
 106 observed methylation value is also dependent on a normally distributed
 107 error term $\epsilon_{i,j}$.

$$\hat{m}_{i,j} = r_i^0 + \sum_{k=1}^n p_{j,k} r_{i,k} + \epsilon_{i,j}$$

108 Under this model epigenetic factors are related to observable individ-
 109 ual factors $p_{k,j}^0$, such as chronological age, sex and cell types, but may
 110 be transformed relative to observable factors. The epigenetic age factor,
 111 for example, often has a non-linear relationship with the observed
 112 age[Snir2019-ii]. The MSEPM learns the appropriate transformation
 113 during model fitting to describe the observed methylation status linearly
 114 in terms of the epigenetic age factor, but not linearly with age.

115 Given an input matrix $\hat{M} = [\hat{m}_{i,j}]$ of methylation values for i sites
 116 and matched observable epigenetic factors $\hat{P}^0 = [p_{j,k}^0]$ for j individuals
 117 the objective of the MSEPM is to find the optimal values of $r_{i,k}$ and $p_{j,k}$
 118 that minimize the residual sum of square (RSS) error,

$$\epsilon_{i,j}^2 = (m_{i,j} - r_i^0 - \sum_{k=1}^n p_{j,k} r_{i,k})^2$$

This is accomplished through the implementation of a conditional expectation maximization algorithm. The maximum likelihood (ML) values of $r_{i,k}$ and r_i^0 can be solved using ordinary least squares (OLS) regression. Provided the ML estimates for $r_{i,k}$, the site coefficients are fixed and epigenetic factors, $p_{j,k}$, are updated by minimizing the RSS across all i sites using gradient descent,

$$p_{j,k}^{n+1} = p_{j,k}^n - \lambda \nabla F(p_{j,k})$$

where λ is a specified learning rate. The optimization is accomplished by alternating between optimizing $r_{i,K}$ and $p_{j,k}$ until the reduction in sum of the site RSS is below a specified threshold or a set number of iterations is reached. Importantly, while the ML values of $p_{j,k}$ are by definition linear with the methylation status at any site, the original input factors for $p_{j,k}^0$ may not be.

Provided a trained MSEPM model and an unobserved methylation matrix, epigenetic factors are estimated by calculating each independent OLS for solution all i sites given the $r_{i,k}$ coefficients set for the respective input factor. These epigenetic factors can then be used to find the expected methylation value using the trained individual site models where

$$E[m_{i,j}] = r_{i,0} + P_j \dot{R}_i$$

where $P_j \dot{R}_i$ is a matrix of point values p and r.

2.2 MSEPM Simulation Framework

We implemented a simulation framework using the MSEPM formulation to evaluate the performance of the MSEPM model under various conditions. To simulate the association between methylation status and an observable factor we modeled the epigenetic factor $p_{k,j}$ as function of time, or age, and magnitude, $p_{k,j}^0$ with a non-linear transformation γ_k , where $p_{k,j} = Age_j p_{k,j}^{0,\gamma_k}$. In practice the value of the $p_{k,j}$ is often unknown and the association between methylation status and $p_{k,j}$ is inferred through the observable factor $p_{k,j}^0$.

Methylation sites were simulated by first randomly setting the range of the methylation site, $-1 < \delta < 1$ a site intercept, r_i^0 , and the site error, $\sigma_i \sim \mathcal{U}(0.025, 0.05)$. The possible range of the methylation site is described by the initial methylation value, $m_0 \beta(.2, .2)$, and the target methylation value, m_t , where the range is $\delta_i = m_t - m_0$. $\beta(.2, .2)$ is the beta distribution with its parameters for randomly setting the sample. The value of m_t is set conditionally to ensure site variability is always larger than some specified threshold, θ , where $\theta \leq |\delta| \geq r_i^0 \beta(.2, .2)$.

Simulated methylation sites are then randomly associated with a combination of zero, one, or multiple epigenetic factors. Rates for sites associated with multiple factors were set by sampling from a uniform distribution. The weighted factor rates are normalized so the input combination

147 of traits describes the range of the simulated site, δ . If a site is associated
148 with no factors the observed methylation status of a site is described by
149 a random normal with a characteristic offset, $\hat{m}_i = r_{0,i} + N(\mu, \sigma)$.

150 2.3 Blood MSEPM Model Training

151 MSEPM models were trained using a large aggregate dataset of blood de-
152 rived methylation data from 17 publicly available datasets[**Ventham2016-qj**,
153 **Demetriou2013-wb**, **Polidoro2013-xv**, **Johansson2013-of**, **Arloth2020-lo**,
154 **Liu2013-dg**, **Soriano-Tarraga2016-uq**, **Chuang2017-nk**, **Zannas2019-me**,
155 **Kurushima2019-pe**, **Voisin2015-lh**, **Tan2014-sg**, **Tserel2015-ro**, **Butcher2017-oz**,
156 **Dabin2020-iw**, **Marabita2013-cj**, **ValleG2019-xi157 lation 450K Beadchip methylation array IDAT files were processed using
158 minfi[**Aryee2014-ky**] (v1.34.0). Sample IDAT files were processed in
159 batches according to GEO series and Beadchip identification. Methyla-
160 tion values within each batch were normal-exponential normalized using
161 out-of-band probes[**Triche2013-pp**]. Blood cell types counts were es-
162 timated using a regression calibration approach[**Houseman2012-rr**] and
163 sex predictions were made using the median intensity measurements of the
164 X and Y chromosomes as implemented in minfi[**Aryee2014-ky**]. Sam-
165 ples were filtered for quality control using the relative intensity of
166 the methylated and unmethylated probes. Samples were used for down-
167 stream analysis if the sample median methylation probe intensity was
168 greater than 10.5 and the difference between the observed and expected
169 median unmethylation probe intensity is less than 0.4, where the expected
170 median unmethylated intensity is described by $E[intensity_{unmethylated}] =$
171 $0.66intensity_{methylated} + 3.718$. This resulted in a total of 5687 samples.**

172 We trained MSEPM models using data assembled from four GEO
173 series[**Johansson2013-of**, **Liu2013-dg**, **Butcher2017-oz**, **Damaso2020-gd**]
174 ($n = 1605$). The samples were randomly split into training ($n = 1203$)
175 and validation ($n = 402$) sets stratified by age. Methylation values for
176 all samples were quantile normalized by probe type[**Horvath2013-sk**]
177 using the median site methylation values across all training samples for
178 each methylation site. Training set blood cell type abundance estimates
179 were used to train a principal component analysis (PCA) model which
180 was then used to calculate cell type PCA estimates for the validation and
181 testing sets. Methylation sites were selected for modeling with MSEPM
182 if the site methylation values were correlated with age ($n = 276$), sex
183 ($n = 49$), CT-PC1 ($n = 120$), CT-PC2 ($n = 116$) or a combination of
184 factors ($n = 238$) by absolute pearson correlation coefficient. Where a
185 absolute pearson correlation coefficient greater than 0.7, 0.995, 0.92 and
186 0.64 for age, sex, CT-PC1 and CT-PC2 respectively. Sites with a sum of
187 absolute pearson coefficients across the four factors greater than 1.8 were
188 also included ($n = 238$) for a total of 778 methylation sites. Min-max,
189 (0-1), scalers were fit using the training input features. Validation and
190 testing sample features were transformed with the trained scalers. Age
191 was min-max scaled on a range from 0-100 years. MSEPM models were
192 trained with a learning rate of 0.01 with an iteration limit of 200.

193

194 **2.4 Blood MSEPM Model Cluster Transcription**

195 **Factor Overlap Analysis**

196 We evaluated the relationship between modeled sites, input factors and
197 regulatory transcription factors using overlap enrichment analysis. We
198 built a custom transcription factor reference set using ENCODE V4 trans-
199 scription factor chromatin immunoprecipitation[**ENCODE Project Consortium2012-oe**,
200 **Davis2018-ha**] (release 1.4.0 - 2.1.2) irreproducible discovery rate nar-
201 row bed peaks, which contains peaks with high rank consistency be-
202 tween replicates, that were not audited for non-compliance or errors.
203 GRCh38 region coordinates were lifted to GRCh37 coordinates using
204 liftOver[**Hinrichs2006-oq**]. The overlap reference contains 714 trans-
205 scription factor targets from 1621 accession IDs.

206 We then performed hierarchical clustering of the four factor MSEPM
207 model sites based on the similarity of their regression coefficients. In-
208 dividual methylation site coefficients were first normalized by the stan-
209 dard deviation of methylation values of the site among the training sam-
210 ples, $r_{i,k}/\sigma_i$. A distance matrix was then created by taking the Eu-
211 clidean distance between the normalized site model coefficients. Sites
212 were then clustered using Ward's method which seeks to minimize within
213 cluster variance by minimizing the increase in the error sum of squares
214 (ESS) through successive cluster fusions. Clusters label by tree cutting
215 at a height of 18. All clustering analysis was carried out using SciPy
216 v1.6.3[**Virtanen2020-mm**].

217 Transcription factor enrichment analysis was performed with LOLA[**Sheffield2016-wg**]
218 which assesses the genomic region set overlap between a set of query re-
219 gions and a set of reference regions, within a specified shared background
220 set, using Fisher's exact test. Overlap analysis was performed for sites
221 within a cluster against the ENCODE V4 reference region (1BP minimum
222 overlap) using all sites assayed with Infinium HumanMethylation450K
BeadChip as background.

223

224 **2.5 Clustering sites with age-associated increases**

225 **in methylation**

226 To better understand age associated methylation in whole blood, we ex-
227 amined each site within MSEPM four factor blood model cluster 7 indi-
228 vidually, as this cluster contains sites that have methylation that increases
229 with age but is not strongly affected by other factors. Using the EWAS
230 Data Hub (Xiong, et al. 2016), we validated our results by obtaining ad-
231 ditional methylation by age data in whole blood for each site in the cluster
232 (McCartney, et al. 2019). We created a matrix with every sample and its
233 associated methylation and age from cluster 7, then used age associated
234 methylation levels to create a clustered heatmap using the Matlab func-
235 tion Clustergram. We then clustered the tree into four groups which were
analyzed separately.

236 We also identified the genes that were proximal to each site using
237 Cistrome-GO (Li et al. 2019). We then examined the expression of the
238 genes across tissues in the Genotype-Tissue Expression (GTEx) database

239 database. We used the GTEx Multi Gene Query to find which tissues
240 those genes belonged to.

241 We utilized the Toolkit for Cistrome Data Browser [Zheng2019-vm,
242 Mei2017-yq] for the analysis of significant factors in each cluster. This
243 allowed us to input .bed files of each sub-cluster and generate a GIG-
244 GLE score for specific transcription factors, histone marks, and chromatin
245 regions to assess significance of these elements. A GIGGLE score tai-
246 lored ranking of loci based on overlap of genomic features provided by the
247 user[Layer2018-gr].

248 2.6 H3K4me3 enrichment analysis

249 Enrichment of analysis for H3K4me3 (Figure 7A) was carried out by
250 downloading rpm normalized bigwig files of H3K4me3 ChIP-seq data from
251 epigenomesportal[Bujold2016-vk] for CD38+ B Cells and CD56+ NTK
252 Cells (for both 0-5 years old and 60-65 years old individuals). Heatmaps
253 of H3K4me3 were generated using deepTools2[Ramirez2016-xl] using
254 the computeMatrix and plotHeatmap function to plot the bigwig signal
255 over genomic regions of cluster 7 as the BED input. The IGV genome
256 browser[Robinson2011-he] was used to generate an image of the KCTD1
257 and IRS2 promoter regions shown in Figure 7B using downloaded bigwig
258 tracks.

259 2.7 Analysis Environment

260 Analysis was carried out in a Jupyter[Basu'undated-vq] analysis envi-
261 ronment. Joblib[Varoquaux2009-al], SciPy[Virtanen2020-wt], Matplotlib[Hunter2007-nq],
262 Seaborn[Waskom2021-gj], Pandas[McKinney2012-ta] and TQDM[Da'Costa-Luis2019-lr]
263 packages were utilized during analysis.

264 3 Results

265 3.1 Simulated Methylation Associated Traits

We simulated individuals whose methylation is determined by four factors
and their associated epigenetic factors: a uniformly distributed factor
approximating age with a non-linear association with methylation status

$$q \sim \mathcal{U}(0, 100), s_{Age} = q^{0.5}, \text{Figure 1A-B}$$

a binary distributed trait resembling sex, linearly associated with methy-
lation status

$$q \sim B(1, .5), s_{Sex} = q, \text{Figure 1C-D}$$

a continuous normal (CN) phenotype a linear association with methylation
status

$$q \sim \mathcal{N}(1, 0.1), s_{CN} = q, \text{Figure 1E-F}$$

and a continuous exponentially (CE) distributed trait with a linear asso-
ciation with methylation status

$$q \sim \frac{1}{20} e^{-x/20}, s_{CE} = q, \text{Figure 1G-H}$$

266 We simulated 90 methylation sites (Figure 1I). We then evaluated
267 the MSEPM model as follows. We simulated 1000 samples with the four
268 epigenetic factors described above. We then simulated methylation values
269 using the simulated site rates. Simulated samples were then split for training
270 ($n = 500$) and testing ($n = 500$). MSEPM models were then fitted
271 using the values of the input factors, $p_{k,j}^0$. We generated 1000 simulated
272 datasets and fit MSEPM models using four combinations of input factors
273 (Age, Age-Sex, Age-Sex-CN, Age-Sex-CN-CE). Within each simulation,
274 epigenetic state predictions and methylation site predictions were made
275 for all testing samples. All models captured the nonlinear association
276 between simulated age and methylation (Supp. Figure 1). As the number of
277 factors in the model is increased the mean absolute error (MAE) between
278 the predicted epigenetic states and the simulated epigenetic factors de-
279 creases (Figure 2A). Importantly, to accurately assess simulated age it is
280 necessary to account for the influence of the other simulated factors (Sex,
281 CN, CE).

282 The MSEPM model generated using all four simulated factors can
283 capture the relative magnitude of the simulated site-specific rates (Figure
284 2C-F). However, the model has difficulty capturing the exact relationship
285 between the simulated factors (age, CN and CE) and the inferred factors
286 (Figure 2C, E-F). This is likely due to limitations of the model at cap-
287 turing nonlinear methylation association and a limited training range for
288 normally and exponentially distributed traits. Regardless, the four-factor
289 model can accurately predict the simulated methylation value (Figure 2
290 D) and site intercept (Supp. Figure 1A). We also assessed the model ro-
291 bustness to variation in the number of samples and sites used for model
292 training by randomly selecting a reduced subset of samples or sites for
293 model training. MSEPM models trained with age, sex, CN, and CE can
294 accurately assess all simulated phenotypes with few samples and sites
295 (Supp. Figure2 B-E).

296 3.2 Blood MSEPM Model

297 We next applied the MSEPM to real data. We utilized a large aggregated
298 dataset composed of Illumina 450k array data from 17 publicly available
299 datasets[Venthram2016-qj, Demetriou2013-wb, Polidoro2013-xv, Johansson2013-of,
300 Arloth2020-lo, Liu2013-dg, Soriano-Tarraga2016-uq, Chuang2017-nk,
301 Zannas2019-me, Kurushima2019-pe, Voisin2015-lh, Tan2014-sg,
302 Tserel2015-ro, Butcher2017-oz, Dabin2020-iw, Marabita2013-cj,
303 ValleG2019-xi] deposited in the Gene Expression Omnibus[Barrett2012-gu]
304 (GEO) generated from blood derived samples (whole blood, peripheral
305 blood lymphocytes, and peripheral blood mononuclear cells). The aggre-
306 gate data spanned a wide age range (0.0 - 99.0 years, Figure 3A), con-
307 tained more predicted females ($n = 3392$) than males ($n = 2295$, Figure
308 3B) and reasonable predicted cell type abundance estimates (Figure 3C).
309 The first principal component of a PCA mode trained cell type abun-
310 dance estimates (CT-PC1) is largely driven by the relative abundance of
311 granulocytes (Figure 3D), while the second PC (CT-PC2) captures rel-
312 ative differences in the abundance of differentiated lymphocytes (Figure
313 3D).

314 We trained MSEPM models using methylation sites ($n = 778$) that
315 were correlated with the observable input factors. MSEPM models were
316 fit using four combinations of input factors (Age, Age Sex, Age Sex CT-
317 PC1, and Age Sex CT-PC1 CT-PC2). The association between the fit epi-
318 genetic factor predictions against the input modeled factors was assessed
319 by fitting a trendline between epigenetic state predictions and scaled con-
320 tinuous input factors using the state prediction made for the MSEPM
321 model trained with all four input factors. Performance of the MSEPM
322 model was then evaluated using the testing samples ($n = 4,082$). The per-
323 formance of the MSEPM largely closely resembles the simulation results.
324 All four MSEPM models capture the nonlinear relationship between age
325 and methylation status (Supp. Figure 6). The epigenetic state prediction
326 associated with age improves as the underlying methylation data are more
327 fully explained through the addition of epigenetic factors (Supp. Figure
328 6). The MSEPM model fit with Age, Sex, CT-PC1 and CT-PC2 can
329 accurately model the associated epigenetic state for each factor (Figure
330 4 A-D) and accurately predicts the methylation levels at individual sites
331 ($R^2 = 0.935$, $MAE = 0.035$, Figure 4 E). The trained MSEPM produces a
332 collection of methylation site models that can help explain the association
333 between modeled factors and methylation status.

334 **3.3 Analysis of chromatin regulators of site clus-
335 ters**

336 We evaluated the relationship between sites that are influenced by age,
337 sex, CT-PC1 or CT-PC2 and potential regulatory factors by performing
338 overlap enrichment analysis of these sites with transcription factor chro-
339 matin immunoprecipitation peaks present in the ENCODE V4[Davis2018-ha,
340 ENCODE'Project'Consortium2012-oe] release. We first identified
341 sites with similar coefficients of epigenetic factors through hierarchical
342 clustering. The resulting tree was cut at a height of 18 to produce 10
343 distinct clusters with clear associations to the modeled factors (Figure
344 5A).

345 The site clusters largely conform to underlying biological expecta-
346 tions. Cluster one contains sites that are wholly associated with sex
347 status and localized to the X chromosome (Supp. Table 1) and is en-
348 riched for peaks of transcription factors associated with sex specific regu-
349 lation such as MAZ[Lopes-Ramos2020-ex]. Clusters nine and ten con-
350 tain sites whose methylation status is largely driven by CT-PC1, and are
351 enriched for transcription factors associated with granulocyte develop-
352 ment (CEPB, CEBPA, EP300, ETV6)[Theilgaard-Monch2022-zw,
353 Guerzoni2006-ii]. Similarly, clusters two, five and eight are associated
354 with CT-PC2 and are enriched for transcription factor peaks associated
355 with immune development (ZBED1, ETV6, FOSL2, FOS, TBX21). Clus-
356 ters four and six are associated with loss of methylation with age. Cluster
357 six is highly enriched for CTCF binding sites; CTCF is known to increase
358 at sites where methylation is lost during aging[Tharakan2020-pj]. Clus-
359 ter four is enriched for STAT3 whose activation during exercise is age
360 dependent[Mohamed2020-he, Trenerry2008-kj]. Cluster seven is as-

361 sociated with the accumulation of methylation with age and is enriched
362 for immunoprecipitation peaks for aging associated transcription factors
363 SMAD4 and RE1-Silencing Transcription Factor (REST). SMAD4 en-
364 codes a protein involved in the transforming growth factor beta (TGF- β)
365 signaling pathway. Age related dysregulation of TGF- β has been linked
366 to reduced skeletal muscle regeneration[**Carlson2009-uz**, **Paris2016-fo**]
367 and SMAD4 polymorphisms are associated with longevity[**Carrieri2004-by**].
368 REST is a transcriptional repressor of neuron specific genes in non-neuronal
369 cells[**Chong1995-dj**, **Coulson2005-pb**]. REST expression is upregu-
370 lated in aged prefrontal cortex tissue and the absence of REST expression
371 is associated with cognitive impairment[**Lu2014-dz**] and cellular senes-
372 cence in neurons[**Rocchi2021-od**].

373

374 **3.4 Analysis of sites with age-associated increases in methylation**

375 Because of our interest in the mechanisms that underlie ages associated
376 increase in methylation, we focused on cluster seven, as these sites have
377 methylation increases that depend primarily on age rather than sex and
378 cell types. Cluster 7 consisted of 93 CpG sites. To obtain an independent
379 measure of how these sites change with age, we obtained age associated
380 methylation

381 data from the EWAS Data Hub[**Xiong2020-fa**], with a focus on whole
382 blood methylation. The dataset consisted of about 1600 individuals with
383 ages ranging from 0 to 113 years old[**McCartney2019-qi**]. We clustered
384 the sites based on age associated methylation levels, meaning the rate of
385 methylation based on age for each marker. Each site was organized into
386 an ordered matrix with methylation levels at each age, then grouped into
387 four sub-clusters: A, B, C, and D. As seen in Figure 6A, Cluster A had the
388 highest average methylation across ages, and each consecutive cluster had
389 a decrease in average methylation. We next examined chromatin accessi-
390 bility, transcription factors, histone marks, and genes associated with each
391 cluster. As shown in Supp. Figure 7, genes proximal to Cluster 7 sites
392 were lowly expressed in blood compared to other tissues. We analyzed
393 chromatin accessibility, transcription factors, and histone marks associ-
394 ated with these four groups. We computed levels of H3K27ac, H3K27me3,
395 H3K4me3, and H3K9me3 across the four subclusters. As seen in Figure
396 6C, H3K4me3 increased from clusters A through D. Figure 6E shows that
397 H3K27ac increased from clusters A through C, but then decreased in D.
398 These results suggest that subcluster D is enriched for bivalent domains,
399 characterized by H3K4me3 and H3K27me3.

400 Based on these results we hypothesize that the mechanisms that underlie
401 the gain of methylation with age at these bivalent promoters is the
402 age-associated loss of H3K4me3. It is well established that the presence
403 of trimethylation on H3K4 inhibits de novo methylation, and this effect
404 explains the hypomethylation that is typical of promoters, including biva-
405 lent promoters. We therefore hypothesize that the gain of methylation at
406 these sites may be caused by an age associated loss of H3K4me3. In or-
407 der to demonstrate that H3K4me3 decreases with age for genomic regions

408 where DNA methylation increases, we used published H3K4me3 ChIP-seq
409 data from epigenomesportal[Bujold2016-vk]. We selected two different
410 blood cell types CD38+ B Cells and CD56+ NTK Cells and plotted the
411 H3K4me3 signal of young (0 to 5 years old) versus old individuals (60 to
412 65 years old) over genomic regions of cluster 7 (Figure 7A). Our analysis
413 shows that younger individuals have higher levels of H3K4me3 compared
414 to older ones (Figure 7A) as also shown for two selected genomic loci of
415 cluster 7 (the promoters of KCTD1 and IRS2 genes) where we can observe
416 a marked decrease in the levels of H3K4me3 as age increases (Figure 7B).
417 All together these data suggest that genomic regions whose DNA methyla-
418 tion is increased with age exhibit an age dependent loss of H3K4me3, thus
419 showing an inverse correlation between DNA methylation and H3K4me3
420 at these genomic loci.

421 4 Discussion

422 Epigenetic clocks are widely used tools to study human aging and health.
423 Despite their widespread use, the biological interpretability of the mod-
424 els is limited. A methylome is influenced by many different biological
425 processes occurring simultaneously over time that may differ among indi-
426 viduals. Epigenetic clocks, while producing accurate predictions of age,
427 attempt to capture this complexity through a single dependent variable.
428 Additionally, the penalized regression based methods used to fit most epi-
429 genetic clocks select sites that minimize, or regress out, the influence of
430 other factors and omit groups of sites that are correlated. To overcome
431 these limitations, here we propose a multidimensional extension of the
432 EPM model, the MSEPM.

433 In contrast to previous methods, the MSEPM aims to simultaneously
434 model the effect of multiple factors on the methylome. The simulation
435 and blood MSEPM models show that concurrently modeling age, cell
436 type composition and sex can minimize model residuals when compared
437 with the MSEPM model fit with age only. The residual of the age only
438 model is often interpreted as a measure of age acceleration. When multiple
439 methylome associated traits are modeled simultaneously this residual can
440 be explained directly by other factors and the association between the
441 methylome and a trait of interest can be inferred.

442 Additionally, the individual methylation site linear models fit as part of
443 the MSEPM optimization can provide information about the relationship
444 between modeled factors and site specific biology. To this end, we find
445 that the blood MSEPM model conforms to expected biology. Sites with a
446 strong sex association localize to the X chromosome and sites associated
447 with cell types are enriched for transcription factors associated with the
448 development of immune cells.

449 CpG sites that are primarily affected only by age in the blood MSEPM
450 model are of particular interest. As others have previously described,
451 sites that progressively lose methylation over time are strongly enriched
452 for CTCF[De'Lima'Camillo2022-lu, Han2020-zj]. As CTCF plays
453 a key role in long range chromatin interactions, this may suggest that
454 there are age-associated changes in three dimensional chromatin structure,

455 and that the structure may become more disordered with age. In fact,
456 alterations in CTCF binding and function with age have been implicated
457 in the pathogenesis of various age-related diseases, including cancer. For
458 example, changes in the chromatin structure and gene expression due
459 to altered CTCF binding can contribute to the genomic instability and
460 altered cell proliferation characteristic of cancerous cells (Hnisz et al.,
461 2016; Phillips et al., 2009).

462 We identified a cluster of sites that showed increasing methylation with
463 age and that were not significantly affected by other factors. We found
464 that these sites are enriched for the transcription factor REST. The RE1-
465 Silencing Transcription Factor (REST), also known as Neuron-Restrictive
466 Silencer Factor (NRSF), is a key regulatory protein involved in the de-
467 velopment and differentiation of neurons. It plays a crucial role in neuro-
468 genesis, neuronal differentiation, and in the maintenance of the neuronal
469 phenotype by regulating gene expression[**Schoenherr1995-fc**]. REST
470 achieves this by binding to the neuron-restrictive silencer element (NRSE)
471 or RE1 sites in the DNA, leading to the repression of gene transcription
472 in non-neuronal cells or in neuronal progenitor cells, ensuring that neu-
473 ronal genes are expressed only in neurons[**Chong1995-dj**, **Ooi2007-kk**,
474 **Bruce2004-je**]. The fact that this factor is enriched at the positively
475 age-associated sites suggests that these sites are likely expressed in neu-
476 ronal cells but not in blood. In fact this is what we find when we examine
477 the tissue specific expression of the genes proximal o these sites.

478 We also examined the histone modifications associated with the pos-
479 itively age-associated sites and found that they were enriched for H3K4me3
480 and H3K27me3. These sites are characteristic of bivalent promoters. Bi-
481 valent promoters play a crucial role in the regulation of gene expression
482 during development and differentiation. Characterized by the simultane-
483 ous presence of both activating (H3K4me3) and repressive (H3K27me3)
484 histone modifications, bivalent promoters mark genes that are poised for
485 transcription but are not actively transcribed. This dual modification
486 serves as a regulatory mechanism, ensuring that genes essential for dif-
487 ferentiation and development are ready to be activated at the appropri-
488 ate time. Bivalent domains are predominantly found in embryonic stem
489 cells and are crucial for maintaining the cells in a pluripotent state, al-
490 lowing for the rapid activation or repression of gene expression in re-
491 sponse to developmental cues. The significance of bivalent promoters
492 extends to their role in cell fate decisions, where they contribute to the
493 tight control of developmental pathways and the maintenance of stem cell
494 identity[**Bernstein2006-wt**, **Voigt2013-fe**]. Our results suggest that
495 the bivalent promoters we identified in blood are inactive (as seen by the
496 fact that the proximal genes are not expressed). However, the fact that
497 DNA methylation at these sites increases with age suggests that they
498 may be losing H3K4me3 with age. H3K4me3 is a critical regulator of
499 DNA methylation as it inhibits the binding of DNMT3 to histones, as the
500 DNMT3 ADD domain preferentially binds to the unmethylated H3K4
501 residue[**Ooi2007-dw**]. This explains why promoters, which are enriched
502 for H3K4me3, are generally hypomethylated. Our results suggests that
503 there must therefore be an age associated loss of H3K4me3 at these bi-
504 valent promoters. That is in fact what we saw when we examined these

505 marks in B cells and Nk cells of both young and old individuals. These
506 mechanisms further suggest that the age associated DNA methylation in-
507 creases may not have a functional consequence in blood and that their
508 proximal genes remain repressed throughout life.

509 In conclusion, we introduced a multi-dimensional extension of the Epi-
510 genetic Pacemaker, the MSEPM. The MSEPM is capable of accurately
511 modeling multiple methylation associated factors simultaneously. This
512 paradigm can elucidate the site specific regulation underpinning methy-
513 lome dynamics. It allows us to characterize the mechanisms underlying
514 age associated increases in methylation sites, suggesting that these were
515 caused by the loss of H3K4Me3 at bivalent promoters of genes that are
516 silenced in blood. The MSEPM is available under the MIT license at
517 <https://github.com/NuttyLogic/MultistateEpigeneticPacemaker>.

518 4.1 Supplementary Information

519 All analysis code, data processing code, and supplementary material asso-
520 ciated with this manuscript can be found at <https://github.com/NuttyLogic/MSEPMManuscript>.
521 The methylation simulation utility can be found at <https://github.com/NuttyLogic/MethSim>.
522 The data supporting these findings are openly available at GEO un-
523 der the series GSE87640, GSE87648, GSE51057, GSE51032, GSE87571,
524 GSE125105, GSE42861, GSE69138, GSE111629, GSE128235, GSE121633,
525 GSE73103, GSE61496, GSE59065, GSE97362, GSE156994, GSE128064
526 and GSE43976.

527 5 Acknowledgments

528 This work has benefited from the equipment and framework of the COMP-
529 HUB and COMP-R Initiatives, funded by the ‘Departments of Excellence’
530 program of the Italian Ministry for University and Research (MIUR, 2018-
531 2022 and MUR, 2023-2027).

532 6 Ethical Statement/Conflict of Interest

533 We have no conflicts of interest to disclose.

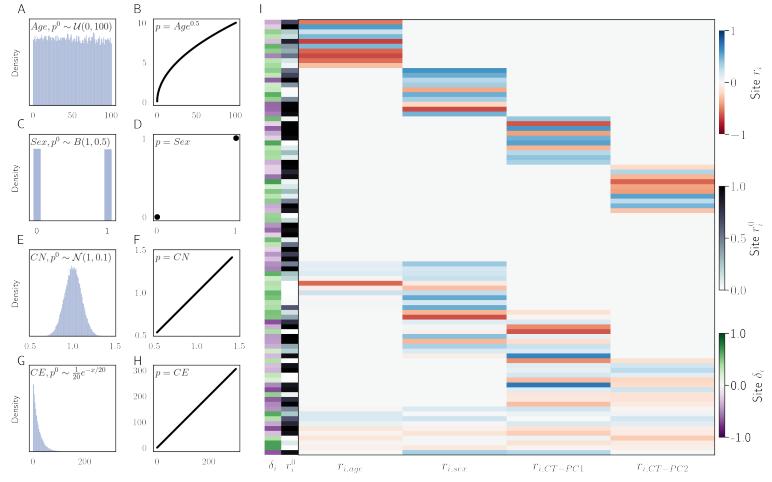


Figure 1: Simulated factors and the association with simulated methylation values. (A) Age with a non-linear association with methylation (B). Sex (C) with a binary association with methylation (D). Normal factor (E) with a linear relationship with methylation (F). Continuous exponential trait (G) with a linear relationship with methylation. (I) Simulated methylation sites. Each simulation site has a starting methylation value r_i^0 , rate of change associated with each simulated factor $r_{i,factor}$ and range of variation δ_i .

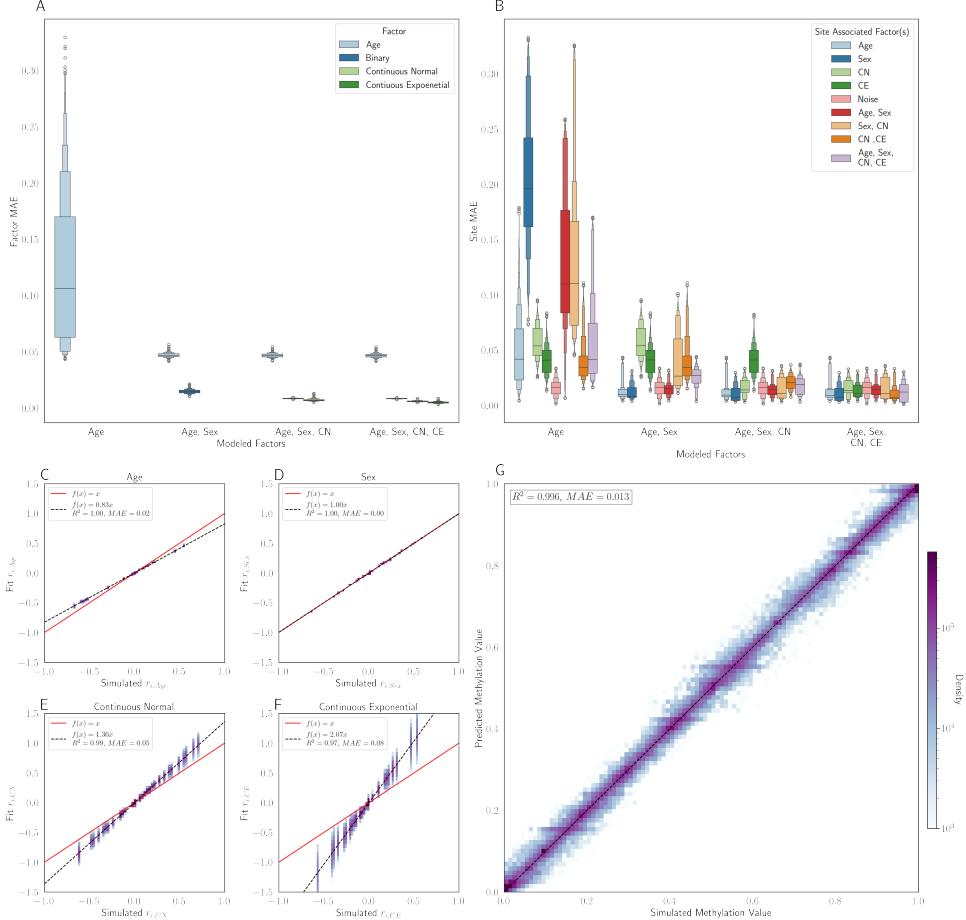


Figure 2: (A) The MAE of the factor predictions on the testing set as multiple factors are modeled simultaneously and (B) predicted methylation MAE as factors are included in the MSEPM model where the centerline is the 50th quantile and the box with greatest width contains 50% of the underlying data with each smaller box containing 50% of the remaining data with 6 levels of box width. (C) Model coefficients for Age, Sex, Continuous Normal and Continuous Exponential factors for models trained ($n = 500$) with all four simulated factors. (D). Simulated and predicted methylation values for all simulated testing sites across all training fold

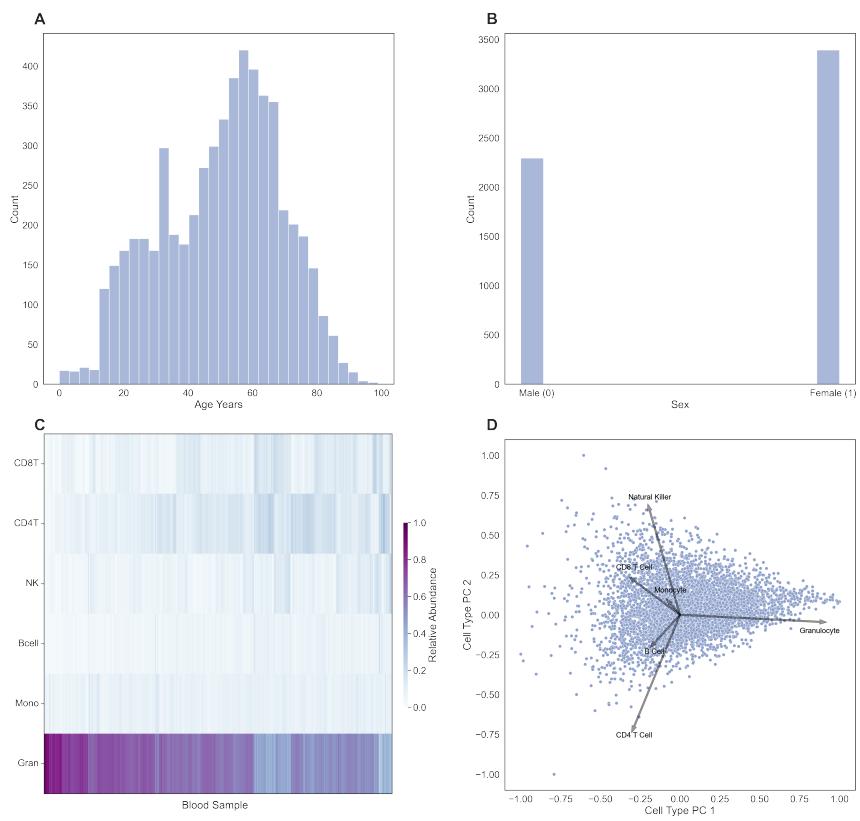


Figure 3: Distribution of age (A) and (B) sex in aggregate blood dataset. (C) Calculated cell type composition and (D) loading plot of principal components of cell type composition in aggregate blood data set.

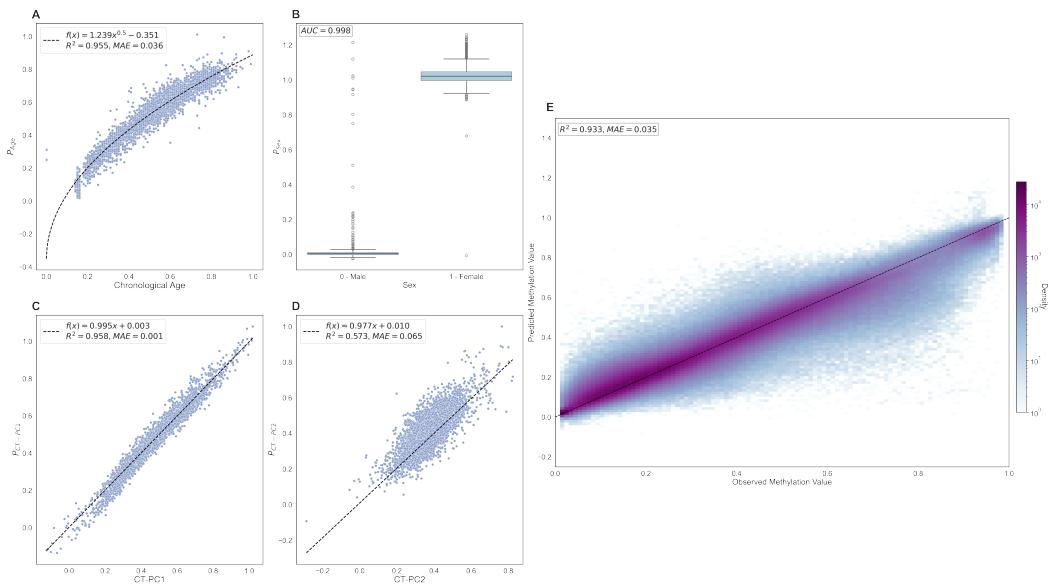


Figure 4: MSEPM model trained with age, sex, CT-PC1 and CT-PC2 predictions within testing set for epigenetic factors (A) age, (B) sex, (C) CT-PC1 and (D) CT-PC2. (E) Observed and predicted methylation values for training set has high concordance ($R^2 = 0.933$, $MAE = 0.035$)

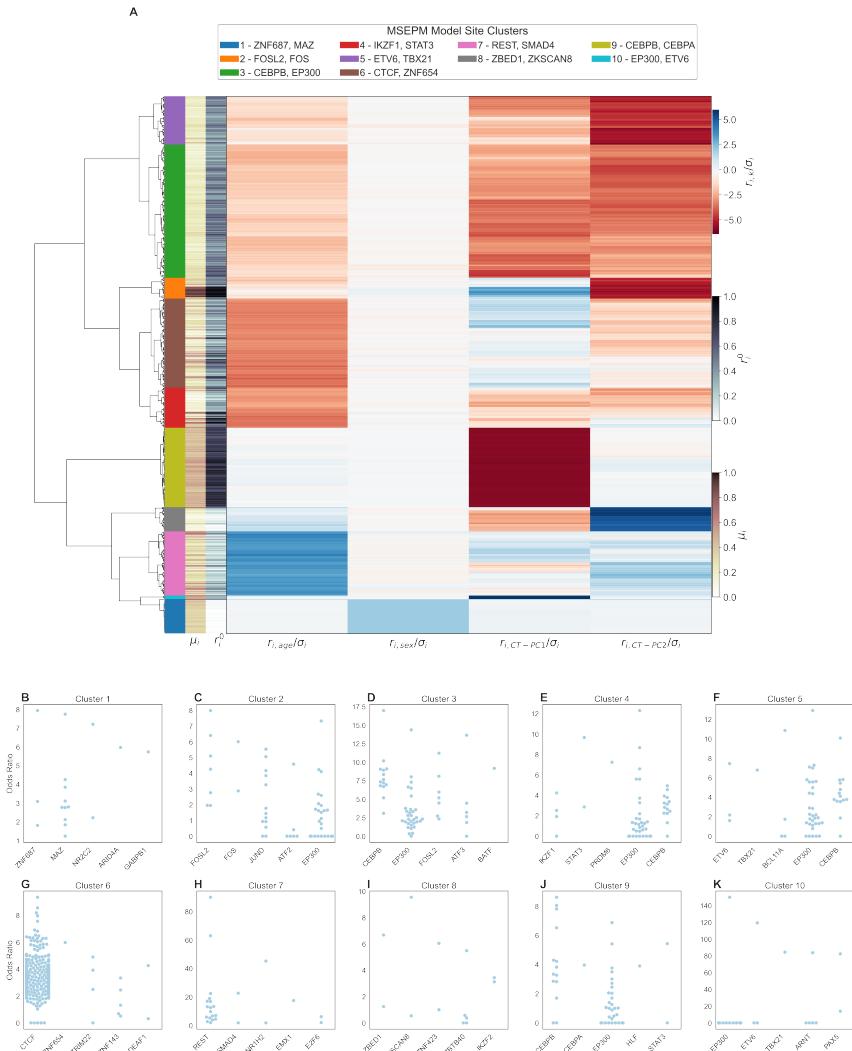


Figure5: (A) Site clustering by standardized model coefficients. Sites clusters show distinct relationships with modeled traits. (B-K) Top five enriched transcription factors for clusters 1 - 10.

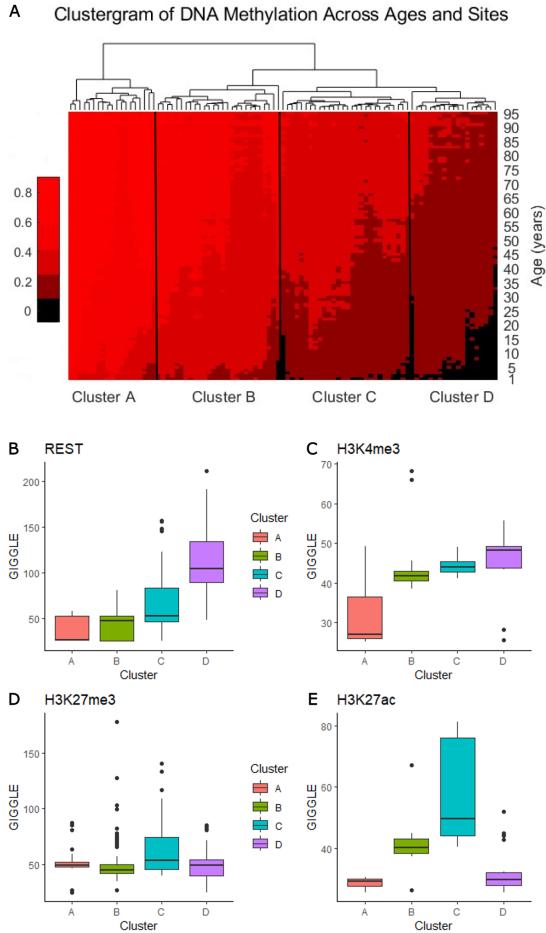


Figure6: (A) Heatmap of H3K4me3 ChIP-seq enrichment for two different blood cell types ($CD38^+$ B Cells and $CD56^+$ NTK Cells) in two cohorts of individual within 0 to 5 years old and 60 to 65 years old. The average level within 2kb up and downstream for centered genomic regions of cluster 7 is represented above the heatmap. (B) Genome browser view of H3K4me3 levels in each cohort at the promoter regions of *KCTD1* and *IRS2* genes.

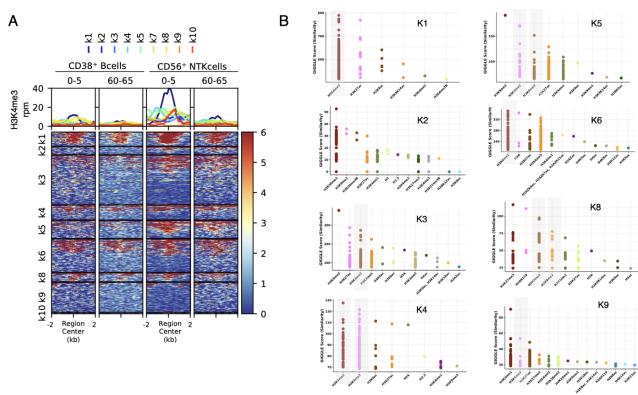


Figure7: (A) Heatmap of H3K4me3 ChIP-seq enrichment for two different blood cell types (CD38⁺ B Cells and CD56⁺ NTK Cells) in two cohorts of individual within 0 to 5 years old and 60 to 65 years old. The average level within 2kb up and downstream for centered genomic regions of cluster 7 is represented above the heatmap. (B) Genome browser view of H3K4me3 levels in each cohort at the promoter regions of *KCTD1* and *IRS2* genes.