

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

Colin Farrell<sup>1,4</sup>, Keshiv Tandon<sup>1</sup>, Roberto Ferrari<sup>2</sup>, Kalsuda Lapborisuth<sup>1</sup>, Rahil Modi<sup>1</sup>, Sagi Snir<sup>3</sup>, and Matteo Pellegrini<sup>1,4</sup>

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

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

<sup>3</sup>Dept. of Evolutionary Biology, University of Haifa, Israel;

<sup>4</sup>Corresponding Authors; colinpfarrell@gmail.com, matteop@mcdb.ucla.edu

## <sup>1</sup> 1 Abstract

<sup>2</sup> Epigenetic clocks, DNA methylation based predictive models of chronological age, are often  
<sup>3</sup> utilized to study aging associated biology. Despite their widespread use, these methods do  
<sup>4</sup> not account for other factors that also contribute to the variability of DNA methylation data.  
<sup>5</sup> For example, many CpG sites show strong sex-specific or cell type specific patterns that likely  
<sup>6</sup> impact the predictions of epigenetic age. To overcome these limitations, we developed a mul-  
<sup>7</sup> tidimensional extension of the Epigenetic Pacemaker, the Multi-State Epigenetic Pacemaker  
<sup>8</sup> (MSEPM). We show that the MSEPM is capable of accurately modeling multiple methyla-  
<sup>9</sup> tion associated factors simultaneously, while also providing site specific models that describe  
<sup>10</sup> the per site relationship between methylation and these factors. We utilized the MSEPM  
<sup>11</sup> with a large aggregate cohort of blood methylation data to construct models of the effects of  
<sup>12</sup> age, sex and cell type heterogeneity on DNA methylation. We found that these models cap-  
<sup>13</sup> ture a large fraction of the variability at thousands of DNA methylation sites. Moreover, this  
<sup>14</sup> approach allows us to identify sites that are primarily affected by aging and no other factors.  
<sup>15</sup> An analysis of these sites reveals that those that lose methylation over time are enriched  
<sup>16</sup> for CTCF transcription factor chip peaks, while those that gain methylation over time are  
<sup>17</sup> associated with bivalent promoters of genes that are not expressed in blood. These obser-  
<sup>18</sup> vations suggest mechanisms that underlie age associated methylation changes and suggest  
<sup>19</sup> that age associated increases in methylation may not have strong functional consequences on  
<sup>20</sup> cell states. In conclusion, the MSEPM is capable of accurately modeling multiple methyla-  
<sup>21</sup> tion associated factors and the models produced can illuminate site specific combinations of  
<sup>22</sup> factors that affect methylation dynamics.

## <sup>23</sup> 2 Introduction

<sup>24</sup> DNA methylation, the addition of a methyl group to the fifth carbon of the cytosine pyrim-  
<sup>25</sup> idine ring, is associated with the topological organization of the cellular genome, gene  
<sup>26</sup> expression and the state of a cell. Within a population of cells the methylation pattern at  
<sup>27</sup> certain sites can change predictably with the age of the individual from which the cells are  
<sup>28</sup> drawn. This predictable nature of DNA methylation has led to the development of accurate  
<sup>29</sup> DNA methylation based predictive models for age and health, termed epigenetic clocks. The  
<sup>30</sup> difference between the predicted and the expected epigenetic age given an individual's chrono-  
<sup>31</sup> logical age has been interpreted as a measure of age acceleration[1], and has been associated  
<sup>32</sup> with mortality[2, 3] and other adverse health outcomes[4–8].

33        However, epigenetic clocks suffer from several limitations that limit the interpretability of  
34        their predictions and the underlying mechanisms. Epigenetic clocks are generally trained by  
35        using penalized regression based methods that attempt to minimize the difference between the  
36        predicted and observed value of age. As a result, as the error between predicted and observed  
37        age is decreased, the associations between age acceleration and mortality disappears[9]. Sec-  
38        ond generation epigenetic clocks attempt to resolve this issue by fitting a measure of human  
39        health, rather than age, and as a result these clocks are generally more sensitive to individ-  
40        ual health status[10–12]. However, while the response variable is modified in these clocks the  
41        method used to fit the clock is largely the same. Epigenetic clocks are generally trained using  
42        regularized regression models, where the likelihood is maximized by minimizing the difference  
43        between the observed and predicted response variable subject to the elastic net penalty,  $\lambda_1$   
44        and  $\lambda_2$ . Methylation sites that increase model error and are influenced by other relevant fac-  
45        tors such as smoking or obesity, may be discarded during model fitting, thus limiting the  
46        ability of this approach to account for the effects of these extraneous factors on epigenetic  
47        aging.

48        As an alternative to penalized regression based methods we previously developed an  
49        evolutionary based model for epigenetic dynamics, the Epigenetic Pacemaker (EPM)[13, 14].  
50        The EPM attempts to minimize the difference between observed and predicted methylation  
51        values amongst a collection of sites through the implementation of a conditional expectation  
52        maximization algorithm[15]. Under the EPM the observed methylation status of a collection  
53        of sites is modeled linearly with respect to an input factor of interest, such as age. A hidden  
54        epigenetic state, that is related to the initial factor, but not necessarily linearly, is learned  
55        through the course of model fitting. The EPM can capture the non-linear relationship between  
56        methylation and age[16] and outputs an interpretable model for each site. However, both the  
57        EPM and regression based methods suffer from the same limitation, which is that they are  
58        limited to a single trait predicted by, or used to model, observed methylation patterns. In  
59        reality, the observed methylation landscape is likely impacted by a variety of factors that act  
60        simultaneously to produce the observed methylome of an individual.

61        To overcome this limitation, we have developed a multidimensional extension of the  
62        EPM, the Multi-State Epigenetic Pacemaker (MSEPM). We show that the (MSEPM) can  
63        accurately model site specific methylation variation driven by several factors, and given a  
64        trained model, accurately predict the values of the factors associated with an individual's  
65        observed methylation profile in both simulated methylation datasets and a large aggregate  
66        blood tissue methylation dataset. Importantly, as factors that explain the observed methy-  
67        lation profile of an individual are added to the model the ability to model the factors and  
68        methylation values improves. Additionally, we show that sites with similar associations to  
69        modeled factors cluster together and are enriched for specific transcription factors. There-  
70        fore, unlike traditional epigenetic clocks, the MSEPM allows us to study mechanisms that  
71        may underlie age associated methylation changes. In our large dataset of blood samples, we  
72        find that sites that increase methylation with age are enriched for bivalent promoters, and  
73        are proximal to genes that are lowly expressed in blood. These results suggest that pos-  
74       itively age associated sites may not have a significant functional impact on aging traits. The  
75        MSEPM is available as a Python package with scikit-learn style syntax under a MIT license  
76        at <https://github.com/NuttyLogic/MultistateEpigeneticPacemaker>.

77 **3 Methods**

78 **3.1 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  $p_{j,k}$ .

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

79 Where  $k$  epigenetic factors are weighted by  $k$  site specific epigenetic rates of change,  $r_{i,k}$ , and  
80 offset by a sites specific intercept term,  $r_i^0$ . Site parameters,  $r_{i,k}$  and  $r_i^0$ , are characteristic of  
81 the site and shared amongst all individuals while epigenetic factors,  $p_{j,k}$ , are characteristic of  
82 an individual and are the same across all sites for that individual. In practice, the observed  
83 methylation value is also dependent on a normally distributed error term  $\epsilon_{i,j}$ .

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

84 Under this model epigenetic factors are related to observable individual factors  $p_{k,j}^0$ , such  
85 as chronological age, sex and cell types, but may be transformed relative to observable factors.  
86 The epigenetic age factor, for example, often has a non-linear relationship with the observed  
87 age[16]. The MSEPM learns the appropriate transformation during model fitting to describe  
88 the observed methylation status linearly in terms of the epigenetic age factor, but not linearly  
89 with age.

90 Given an input matrix  $\hat{M} = [m_{i,j}]$  of methylation values for  $i$  sites and matched observable  
91 epigenetic factors  $\hat{P}^0 = [p_{j,k}^0]$  for  $j$  individuals the objective of the MSEPM is to find the  
92 optimal values of  $r_{i,k}$  and  $p_{j,k}$  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,

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

93 where  $\lambda$  is a specified learning rate. The optimization is accomplished by alternating between  
94 optimizing  $r_{i,K}$  and  $p_{j,k}$  until the reduction in sum of the site RSS is below a specified  
95 threshold or a set number of iterations is reached. Importantly, while the ML values of  $p_{j,k}$   
96 are by definition linear with the methylation status at any site, the original input factors for  
97  $p_{j,k}^0$  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

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

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

<sub>99</sub> **3.2 MSEPM Simulation Framework**

<sub>100</sub> We implemented a simulation framework using the MSEPM formulation to evaluate the  
<sub>101</sub> performance of the MSEPM model under various conditions. To simulate the association  
<sub>102</sub> between methylation status and an observable factor we modeled the epigenetic factor  $p_{k,j}$   
<sub>103</sub> as function of time, or age, and magnitude,  $p_{k,j}^0$  with a non-linear transformation  $\gamma_k$ , where  
<sub>104</sub>  $p_{k,j} = Age_j p_{k,j}^{0,\gamma_k}$ . In practice the value of the  $p_{k,j}$  is often unknown and the association  
<sub>105</sub> between methylation status and  $p_{k,j}$  is inferred through the observable factor  $p_{k,j}^0$ .

<sub>106</sub> Methylation sites were simulated by first randomly setting the range of the methylation  
<sub>107</sub> 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  
<sub>108</sub> range of the methylation site is described by the initial methylation value,  $m_0 \beta(.2, .2)$ , and  
<sub>109</sub> the target methylation value,  $m_t$ , where the range is  $\delta_i = m_t - m_0$ .  $\beta(.2, .2)$  is the beta  
<sub>110</sub> distribution with its parameters for randomly setting the sample. The value of  $m_t$  is set  
<sub>111</sub> conditionally to ensure site variability is always larger than some specified threshold,  $\theta$ , where  
<sub>112</sub>  $\theta \leq |\delta| \geq .r_i^0 \beta(.2, .2)$ .

<sub>113</sub> Simulated methylation sites are then randomly associated with a combination of zero,  
<sub>114</sub> one, or multiple epigenetic factors. Rates for sites associated with multiple factors were set  
<sub>115</sub> by sampling from a uniform distribution. The weighted factor rates are normalized so the  
<sub>116</sub> input combination of traits describes the range of the simulated site,  $\delta$ . If a site is associated  
<sub>117</sub> with no factors the observed methylation status of a site is described by a random normal  
<sub>118</sub> with a characteristic offset,  $\hat{m}_i = r_{0,i} + N(\mu, \sigma)$ .

<sub>119</sub> **3.3 Blood MSEPM Model Training**

<sub>120</sub> MSEPM models were trained using a large aggregate dataset of blood derived methylation  
<sub>121</sub> data from 17 publicly available datasets[7, 17–32]. Illumina methylation 450K Beadchip  
<sub>122</sub> methylation array IDAT files were processed using minfi[33] (v1.34.0). Sample IDAT files  
<sub>123</sub> were processed in batches according to GEO series and Beadchip identification. Methylation  
<sub>124</sub> values within each batch were normal-exponential normalized using out-of-band probes[34].  
<sub>125</sub> Blood cell types counts were estimated using a regression calibration approach[35] and sex  
<sub>126</sub> predictions were made using the median intensity measurements of the X and Y chromosomes  
<sub>127</sub> as implemented in minfi[33]. Samples were filtered for quality control using the the  
<sub>128</sub> relative intensity of the methylated and unmethylated probes. Samples were used for down-  
<sub>129</sub> stream analysis if the sample median methylation probe intensity was greater than 10.5  
<sub>130</sub> and the difference between the observed and expected median unmethylation probe inten-  
<sub>131</sub> sity is less than 0.4, where the expected median unmethylated intensity is described by  
<sub>132</sub>  $E[intensity_{unmethylated}] = 0.66intensity_{methylated} + 3.718$ . This resulted in a total of 5687  
<sub>133</sub> samples.

<sub>134</sub> We trained MSEPM models using data assembled from four GEO series[20, 22, 29, 36]  
<sub>135</sub> ( $n = 1605$ ). The samples were randomly split into training ( $n = 1203$ ) and validation ( $n =$   
<sub>136</sub> 402) sets stratified by age. Methylation values for all samples were quantile normalized by  
<sub>137</sub> probe type[37] using the median site methylation values across all training samples for each  
<sub>138</sub> methylation site. Training set blood cell type abundance estimates were used to train a  
<sub>139</sub> principal component analysis (PCA) model which was then used to calculate cell type PCA  
<sub>140</sub> estimates for the validation and testing sets. Methylation sites were selected for modeling  
<sub>141</sub> with MSEPM if the site methylation values were correlated with age ( $n = 276$ ), sex ( $n = 49$ ),  
<sub>142</sub> CT-PC1 ( $n = 120$ ), CT-PC2 ( $n = 116$ ) or a combination of factors ( $n = 238$ ) by absolute  
<sub>143</sub> pearson correlation coefficient. Where a absolute pearson correlation coefficient greater than  
<sub>144</sub> 0 .7, 0.995, 0.92 and 0.64 for age, sex, CT-PC1 and CT-PC2 respectively. Sites with a sum  
<sub>145</sub> of absolute pearson coefficients across the four factors greater than 1.8 were also included

146 ( $n = 238$ ) for a total of 778 methylation sites. Min-max, (0-1), scalers were fit using the  
147 training input features. Validation and testing sample features were transformed with the  
148 trained scalers. Age was min-max scaled on a range from 0-100 years. MSEPM models were  
149 trained with a learning rate of 0.01 with an iteration limit of 200.

150 **3.4 Blood MSEPM Model Cluster Transcription Factor  
151 Overlap Analysis**

152 We evaluated the relationship between modeled sites, input factors and regulatory tran-  
153 scription factors using overlap enrichment analysis. We built a custom transcription factor  
154 reference set using ENCODE V4 transcription factor chromatin immunoprecipitation[38, 39]  
155 (release 1.4.0 - 2.1.2) irreproducible discovery rate narrow bed peaks, which contains peaks  
156 with high rank consistency between replicates, that were not audited for non-compliance or  
157 errors. GRCh38 region coordinates were lifted to GRCh37 coordinates using liftOver[40]. The  
158 overlap reference contains 714 transcription factor targets from 1621 accession IDs.

159 We then performed hierarchical clustering of the four factor MSEPM model sites based  
160 on the similarity of their regression coefficients. Individual methylation site coefficients were  
161 first normalized by the standard deviation of methylation values of the site among the train-  
162 ing samples,  $r_{i,k}/\sigma_i$ . A distance matrix was then created by taking the Euclidean distance  
163 between the normalized site model coefficients. Sites were then clustered using Ward's method  
164 which seeks to minimize within cluster variance by minimizing the increase in the error sum  
165 of squares (ESS) through successive cluster fusions. Clusters label by tree cutting at a height  
166 of 18. All clustering analysis was carried out using SciPy v1.6.3[41].

167 Transcription factor enrichment analysis was performed with LOLA[42] which assesses the  
168 genomic region set overlap between a set of query regions and a set of reference regions, within  
169 a specified shared background set, using Fisher's exact test. Overlap analysis was performed  
170 for sites within a cluster against the ENCODE V4 reference region (1BP minimum overlap)  
171 using all sites assayed with Infinium HumanMethylation450K BeadChip as background.

172 **3.5 Clustering sites with age-associated increases in  
173 methylation**

174 To better understand age associated methylation in whole blood, we examined each site  
175 within MSEPM four factor blood model cluster 7 individually, as this cluster contains sites  
176 that have methylation that increases with age but is not strongly affected by other factors.  
177 Using the EWAS Data Hub (Xiong, et al. 2016), we validated our results by obtaining addi-  
178 tional methylation by age data in whole blood for each site in the cluster (McCartney, et al.  
179 2019). We created a matrix with every sample and its associated methylation and age from  
180 cluster 7, then used age associated methylation levels to create a clustered heatmap using  
181 the Matlab function Clustergram. We then clustered the tree into four groups which were  
182 analyzed separately.

183 We also identified the genes that were proximal to each site using Cistrome-GO (Li et al.  
184 2019). We then examined the expression of the genes across tissues in the Genotype-Tissue  
185 Expression (GTEx) database database. We used the GTEx Multi Gene Query to find which  
186 tissues those genes belonged to.

187 We utilized the Toolkit for Cistrome Data Browser [43, 44] for the analysis of significant  
188 factors in each cluster. This allowed us to input .bed files of each sub-cluster and generate  
189 a GIGGLE score for specific transcription factors, histone marks, and chromatin regions  
190 to assess significance of these elements. A GIGGLE score tailored ranking of loci based on  
191 overlap of genomic features provided by the user[45].

192 **3.6 H3K4me3 enrichment analysis**

193 Enrichment of analysis for H3K4me3 (Figure 7A) was carried out by downloading rpm nor-  
194 malized bigwig files of H3K4me3 ChIP-seq data from epigenomesportal[46] for CD38+ B Cells  
195 and CD56+ NTK Cells (for both 0-5 years old and 60-65 years old individuals). Heatmaps of  
196 H3K4me3 were generated using deepTools2[47] using the computeMatrix and plotHeatmap  
197 function to plot the bigwig signal over genomic regions of cluster 7 as the BED input. The  
198 IGV genome browser[48] was used to generate an image of the KCTD1 and IRS2 promoter  
199 regions shown in Figure 7B using downloaded bigwig tracks.

200 **3.7 Analysis Environment**

201 Analysis was carried out in a Jupyter[49] analysis environment. Joblib[50], SciPy[51],  
202 Matplotlib[52], Seaborn[53], Pandas[54] and TQDM[55] packages were utilized during analy-  
203 sis.

204 **4 Results**

205 **4.1 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 methylation 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 association with methylation status

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

206 We simulated 90 methylation sites (Figure 1I). We then evaluated the MSEPM model  
207 as follows. We simulated 1000 samples with the four epigenetic factors described above.  
208 We then simulated methylation values using the simulated site rates. Simulated samples  
209 were then split for training ( $n = 500$ ) and testing ( $n = 500$ ). MSEPM models were then  
210 fitted using the values of the input factors,  $p_{k,j}^0$ . We generated 1000 simulated datasets and  
211 fit MSEPM models using four combinations of input factors (Age, Age-Sex, Age-Sex-CN,  
212 Age-Sex-CN-CE). Within each simulation, epigenetic state predictions and methylation site  
213 predictions were made for all testing samples. All models captured the nonlinear association  
214 between simulated age and methylation (Supp. Figure 1). As the number of factors in the  
215 model is increased, the mean absolute error (MAE) between the predicted epigenetic states  
216 and the simulated epigenetic factors decreases (Figure 2A). Importantly, to accurately assess  
217 simulated age it is necessary to account for the influence of the other simulated factors (Sex,  
218 CN, CE).

219 The MSEPM model generated using all four simulated factors can capture the relative  
220 magnitude of the simulated site-specific rates (Figure 2C-F). However, the model has difficulty  
221 capturing the exact relationship between the simulated factors (age, CN and CE) and the

222 inferred factors (Figure 2C, E-F). This is likely due to limitations of the model at capturing  
223 nonlinear methylation association and a limited training range for normally and exponentially  
224 distributed traits. Regardless, the four-factor model can accurately predict the simulated  
225 methylation value (Figure 2 D) and site intercept (Supp. Figure 1A). We also assessed the  
226 model robustness to variation in the number of samples and sites used for model training by  
227 randomly selecting a reduced subset of samples or sites for model training. MSEPM models  
228 trained with age, sex, CN, and CE can accurately assess all simulated phenotypes with few  
229 samples and sites (Supp. Figure2 B-E).

## 230 4.2 Blood MSEPM Model

231 We next applied the MSEPM to real data. We utilized a large aggregated dataset composed  
232 of Illumina 450k array data from 17 publicly available datasets[7, 17–32] deposited in the  
233 Gene Expression Omnibus[56] (GEO) generated from blood derived samples (whole blood,  
234 peripheral blood lymphocytes, and peripheral blood mononuclear cells). The aggregate data  
235 spanned a wide age range (0.0 - 99.0 years, Figure 3A), contained more predicted females  
236 ( $n = 3392$ ) than males ( $n = 2295$ , Figure 3B) and reasonable predicted cell type abundance  
237 estimates (Figure 3C). The first principal component of a PCA model trained cell type  
238 abundance estimates (CT-PC1) is largely driven by the relative abundance of granulocytes  
239 (Figure 3D), while the second PC (CT-PC2) captures relative differences in the abundance  
240 of differentiated lymphocytes (Figure 3D).

241 We trained MSEPM models using methylation sites ( $n = 778$ ) that were correlated with  
242 the observable input factors. MSEPM models were fit using four combinations of input factors  
243 (Age, Age Sex, Age Sex CT-PC1, and Age Sex CT-PC1 CT-PC2). The association between  
244 the fit epigenetic factor predictions against the input modeled factors was assessed by fitting  
245 a trendline between epigenetic state predictions and scaled continuous input factors using the  
246 state prediction made for the MSEPM model trained with all four input factors. Performance  
247 of the MSEPM model was then evaluated using the testing samples ( $n = 4,082$ ). The per-  
248 formance of the MSEPM largely closely resembles the simulation results. All four MSEPM  
249 models capture the nonlinear relationship between age and methylation status (Supp. Figure  
250 6). The epigenetic state prediction associated with age improves as the underlying methyla-  
251 tion data are more fully explained through the addition of epigenetic factors (Supp. Figure 6).  
252 The MSEPM model fit with Age, Sex, CT-PC1 and CT-PC2 can accurately model the asso-  
253 ciated epigenetic state for each factor (Figure 4 A-D) and accurately predicts the methylation  
254 levels at individual sites ( $R^2 = 0.935$ ,  $MAE = 0.035$ , Figure 4 E). The trained MSEPM pro-  
255 duces a collection of methylation site models that can help explain the association between  
256 modeled factors and methylation status.

## 257 4.3 Analysis of chromatin regulators of site clusters

258 We evaluated the relationship between sites that are influenced by age, sex, CT-PC1 or CT-  
259 PC2 and potential regulatory factors by performing overlap enrichment analysis of these  
260 sites with transcription factor chromatin immunoprecipitation peaks present in the ENCODE  
261 V4[38, 39] release. We first identified sites with similar coefficients of epigenetic factors  
262 through hierarchical clustering. The resulting tree was cut at a height of 18 to produce 10  
263 distinct clusters with clear associations to the modeled factors (Figure 5A).

264 The site clusters largely conform to underlying biological expectations. Cluster one con-  
265 tains sites that are wholly associated with sex status and localized to the X chromosome  
266 (Supp. Table 1) and is enriched for peaks of transcription factors associated with sex spe-  
267 cific regulation such as MAZ[57]. Clusters nine and ten contain sites whose methylation

status is largely driven by CT-PC1, and are enriched for transcription factors associated with granulocyte development (CEBPB, CEBPA, EP300, ETV6)[58, 59]. Similarly, clusters two, five and eight are associated with CT-PC2 and are enriched for transcription factor peaks associated with immune development (ZBED1, ETV6, FOSL2, FOS, TBX21). Clusters four and six are associated with loss of methylation with age. Cluster six is highly enriched for CTCF binding sites; CTCF is known to increase at sites where methylation is lost during aging[60]. Cluster four is enriched for STAT3 whose activation during exercise is age dependent[61, 62]. Cluster seven is associated with the accumulation of methylation with age and is enriched for immunoprecipitation peaks for aging associated transcription factors SMAD4 and RE1-Silencing Transcription Factor (REST). SMAD4 encodes a protein involved in the transforming growth factor beta (TGF- $\beta$ ) signaling pathway. Age related dysregulation of TGF- $\beta$  has been linked to reduced skeletal muscle regeneration[63, 64] and SMAD4 polymorphisms are associated with longevity[65]. REST is a transcriptional repressor of neuron specific genes in non-neuronal cells[66, 67]. REST expression is upregulated in aged prefrontal cortex tissue and the absence of REST expression is associated with cognitive impairment[68] and cellular senescence in neurons[69].

#### 4.4 Analysis of sites with age-associated increases in methylation

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

data from the EWAS Data Hub[70], with a focus on whole blood methylation. The dataset consisted of about 1600 individuals with ages ranging from 0 to 113 years old[71]. We clustered the sites based on age associated methylation levels, meaning the rate of methylation based on age for each marker. Each site was organized into an ordered matrix with methylation levels at each age, then grouped into four sub-clusters: A, B, C, and D. As seen in Figure 6A, Cluster A had the highest average methylation across ages, and each consecutive cluster had a decrease in average methylation. We next examined chromatin accessibility, transcription factors, histone marks, and genes associated with each cluster. As shown in Supp. Figure 7, genes proximal to Cluster 7 sites were lowly expressed in blood compared to other tissues. We analyzed chromatin accessibility, transcription factors, and histone marks associated with these four groups. We computed levels of H3K27ac, H3K27me3, H3K4me3, and H3K9me3 across the four subclusters. As seen in Figure 6C, H3K4me3 increased from clusters A through D. Figure 6E shows that H3K27ac increased from clusters A through C, but then decreased in D. These results suggest that subcluster D is enriched for bivalent domains, characterized by H3K4me3 and H3K27me3.

Based on these results we hypothesize that the mechanisms that underlie the gain of methylation with age at these bivalent promoters is the age-associated loss of H3K4me3. It is well established that the presence of trimethylation on H3K4 inhibits de novo methylation, and this effect explains the hypomethylation that is typical of promoters, including bivalent promoters. We therefore hypothesize that the gain of methylation at these sites may be caused by an age associated loss of H3K4me3. In order to demonstrate that H3K4me3 decreases with age for genomic regions where DNA methylation increases, we used published H3K4me3 ChIP-seq data from epigenomesportal[46]. We selected two different blood cell types CD38+ B Cells and CD56+ NTK Cells and plotted the H3K4me3 signal of young (0 to 5 years old) versus old individuals (60 to 65 years old) over genomic regions of cluster 7 (Figure 7A). Our

316 analysis shows that younger individuals have higher levels of H3K4me3 compared to older  
317 ones (Figure 7A) as also shown for two selected genomic loci of cluster 7 (the promoters of  
318 KCTD1 and IRS2 genes) where we can observe a marked decrease in the levels of H3K4me3 as  
319 age increases (Figure 7B). All together these data suggest that genomic regions whose DNA  
320 methylation is increased with age exhibit an age dependent loss of H3K4me3, thus showing  
321 an inverse correlation between DNA methylation and H3K4me3 at these genomic loci.

## 322 5 Discussion

323 Epigenetic clocks are widely used tools to study human aging and health. Despite their  
324 widespread use, the biological interpretability of the models is limited. A methylome is influ-  
325 enced by many different biological processes occurring simultaneously over time that may  
326 differ among individuals. Epigenetic clocks, while producing accurate predictions of age,  
327 attempt to capture this complexity through a single dependent variable. Additionally, the  
328 penalized regression based methods used to fit most epigenetic clocks select sites that mini-  
329 mize, or regress out, the influence of other factors and omit groups of sites that are correlated.  
330 To overcome these limitations, here we propose a multidimensional extension of the EPM  
331 model, the MSEPM.

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

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

345 CpG sites that are primarily affected only by age in the blood MSEPM model are of  
346 particular interest. As others have previously described, sites that progressively lose methy-  
347 lation over time are strongly enriched for CTCF[[72](#), [73](#)]. As CTCF plays a key role in long  
348 range chromatin interactions, this may suggest that there are age-associated changes in three  
349 dimensional chromatin structure, and that the structure may become more disordered with  
350 age. In fact, alterations in CTCF binding and function with age have been implicated in the  
351 pathogenesis of various age-related diseases, including cancer. For example, changes in the  
352 chromatin structure and gene expression due to altered CTCF binding can contribute to the  
353 genomic instability and altered cell proliferation characteristic of cancerous cells (Hnisz et  
354 al., 2016; Phillips et al., 2009).

355 We identified a cluster of sites that showed increasing methylation with age and that  
356 were not significantly affected by other factors. We found that these sites are enriched  
357 for the transcription factor REST. The RE1-Silencing Transcription Factor (REST), also  
358 known as Neuron-Restrictive Silencer Factor (NRSF), is a key regulatory protein involved  
359 in the development and differentiation of neurons. It plays a crucial role in neurogenesis,  
360 neuronal differentiation, and in the maintenance of the neuronal phenotype by regulating  
361 gene expression[[74](#)]. REST achieves this by binding to the neuron-restrictive silencer element

362 (NRSE) or RE1 sites in the DNA, leading to the repression of gene transcription in non-  
363 neuronal cells or in neuronal progenitor cells, ensuring that neuronal genes are expressed only  
364 in neurons[66, 75, 76]. The fact that this factor is enriched at the positively age-associated  
365 sites suggests that these sites are likely expressed in neuronal cells but not in blood. In fact  
366 this is what we find when we examine the tissue specific expression of the genes proximal o  
367 these sites.

368 We also examined the histone modifications associated with the positively age-associated  
369 sites and found that they were enriched for H3K4me3 and H3K27me3. These sites are char-  
370 acteristic of bivalent promoters. Bivalent promoters play a crucial role in the regulation of  
371 gene expression during development and differentiation. Characterized by the simultaneous  
372 presence of both activating (H3K4me3) and repressive (H3K27me3) histone modifications,  
373 bivalent promoters mark genes that are poised for transcription but are not actively tran-  
374 scribed. This dual modification serves as a regulatory mechanism, ensuring that genes  
375 essential for differentiation and development are ready to be activated at the appropriate  
376 time. Bivalent domains are predominantly found in embryonic stem cells and are crucial for  
377 maintaining the cells in a pluripotent state, allowing for the rapid activation or repression  
378 of gene expression in response to developmental cues. The significance of bivalent promot-  
379 ers extends to their role in cell fate decisions, where they contribute to the tight control of  
380 developmental pathways and the maintenance of stem cell identity[77, 78]. Our results sug-  
381 gest that the bivalent promoters we identified in blood are inactive (as seen by the fact that  
382 the proximal genes are not expressed). However, the fact that DNA methylation at these  
383 sites increases with age suggests that they may be losing H3K4me3 with age. H3K4me3 is a  
384 critical regulator of DNA methylation as it inhibits the binding of DNMT3 to histones, as  
385 the DNMT3 ADD domain preferentially binds to the unmethylated H3K4 residue[79]. This  
386 explains why promoters, which are enriched for H3K4me3, are generally hypomethylated.  
387 Our results suggests that there must therefore be an age associated loss of H3K4me3 at these  
388 bivalent promoters. That is in fact what we saw when we examined these marks in B cells  
389 and Nk cells of both young and old individuals. These mechanisms further suggest that the  
390 age associated DNA methylation increases may not have a functional consequence in blood  
391 and that their proximal genes remain repressed throughout life.

392 In conclusion, we introduced a multi-dimensional extension of the Epigenetic Pacemaker,  
393 the MSEPM. The MSEPM is capable of accurately modeling multiple methylation associated  
394 factors simultaneously. This paradigm can elucidate the site specific regulation underpinning  
395 methylome dynamics. It allows us to characterize the mechanisms underlying age associated  
396 increases in methylation sites, suggesting that these were caused by the loss of H3K4Me3 at  
397 bivalent promoters of genes that are silenced in blood. The MSEPM is available under the  
398 MIT license at <https://github.com/NuttyLogic/MultistateEpigeneticPacemaker>.

## 399 5.1 Supplementary Information

400 All analysis code, data processing code, and supplementary material associated with  
401 this manuscript can be found at <https://github.com/NuttyLogic/MSEPMManuscript>. The  
402 methylation simulation utility can be found at <https://github.com/NuttyLogic/MethSim>.  
403 The data supporting these findings are openly available at GEO under the  
404 series GSE87640, GSE87648, GSE51057, GSE51032, GSE87571, GSE125105, GSE42861,  
405 GSE69138, GSE111629, GSE128235, GSE121633, GSE73103, GSE61496, GSE59065,  
406 GSE97362, GSE156994, GSE128064 and GSE43976.

## 407 6 Acknowledgments

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

## 411 7 Ethical Statement/Conflict of Interest

412 We have no conflicts of interest to disclose.

## References

- [1] Horvath, S. & Raj, K. DNA methylation-based biomarkers and the epigenetic clock theory of ageing. *Nat. Rev. Genet.* **19**, 371–384 (2018).
- [2] Perna, L. *et al.* Epigenetic age acceleration predicts cancer, cardiovascular, and all-cause mortality in a german case cohort (2016).
- [3] Marioni, R. E. *et al.* DNA methylation age of blood predicts all-cause mortality in later life. *Genome Biol.* **16**, 25 (2015).
- [4] Dugué, P.-A. *et al.* DNA methylation-based biological aging and cancer risk and survival: Pooled analysis of seven prospective studies. *Int. J. Cancer* **142**, 1611–1619 (2018).
- [5] Huang, R.-C. *et al.* Epigenetic age acceleration in adolescence associates with BMI, inflammation, and risk score for middle age cardiovascular disease. *J. Clin. Endocrinol. Metab.* **104**, 3012–3024 (2019).
- [6] Armstrong, N. J. *et al.* Aging, exceptional longevity and comparisons of the hannum and horvath epigenetic clocks. *Epigenomics* **9**, 689–700 (2017).
- [7] Chuang, Y.-H. *et al.* Parkinson’s disease is associated with DNA methylation levels in human blood and saliva. *Genome Med.* **9**, 76 (2017).
- [8] Horvath, S. & Ritz, B. R. Increased epigenetic age and granulocyte counts in the blood of parkinson’s disease patients. *Aging* **7**, 1130–1142 (2015).
- [9] Zhang, Q. *et al.* Improved precision of epigenetic clock estimates across tissues and its implication for biological ageing (2019).
- [10] Lu, A. T. *et al.* DNA methylation GrimAge strongly predicts lifespan and healthspan. *Aging* **11**, 303–327 (2019).
- [11] Levine, M. E. *et al.* An epigenetic biomarker of aging for lifespan and healthspan. *Aging* **10**, 573–591 (2018).

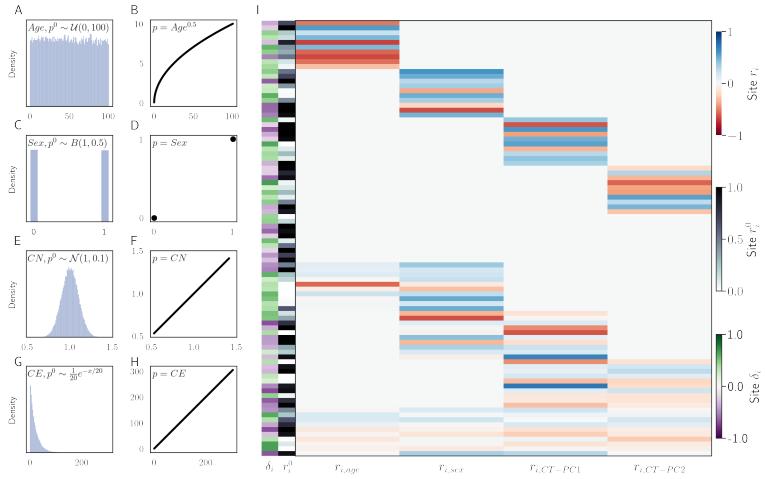
- [12] Belsky, D. W. *et al.* Quantification of the pace of biological aging in humans through a blood test, the DunedinPoAm DNA methylation algorithm. *Elife* **9** (2020).
- [13] Farrell, C., Snir, S. & Pellegrini, M. The epigenetic pacemaker: modeling epigenetic states under an evolutionary framework. *Bioinformatics* **36**, 4662–4663 (2020).
- [14] Snir, S., vonHoldt, B. M. & Pellegrini, M. A statistical framework to identify deviation from time linearity in epigenetic aging. *PLoS Comput. Biol.* **12**, e1005183 (2016).
- [15] Snir, S. Epigenetic pacemaker: closed form algebraic solutions. *BMC Genomics* **21**, 257 (2020).
- [16] Snir, S., Farrell, C. & Pellegrini, M. Human epigenetic ageing is logarithmic with time across the entire lifespan. *Epigenetics* **14**, 912–926 (2019).
- [17] Venthram, N. T. *et al.* Integrative epigenome-wide analysis demonstrates that DNA methylation may mediate genetic risk in inflammatory bowel disease. *Nat. Commun.* **7**, 13507 (2016).
- [18] Demetriou, C. A. *et al.* Methylome analysis and epigenetic changes associated with menarcheal age. *PLoS One* **8**, e79391 (2013).
- [19] Polidoro, S. *et al.* EPIC-Italy at HuGeF. GSE51032. *Gene Expression Omnibus* (2013).
- [20] Johansson, A., Enroth, S. & Gyllensten, U. Continuous aging of the human DNA methylome throughout the human lifespan. *PLoS One* **8**, e67378 (2013).
- [21] Arloth, J. *et al.* DeepWAS: Multivariate genotype-phenotype associations by directly integrating regulatory information using deep learning. *PLoS Comput. Biol.* **16**, e1007616 (2020).
- [22] Liu, Y. *et al.* Epigenome-wide association data implicate DNA methylation as an intermediary of genetic risk in rheumatoid arthritis (2013).
- [23] Soriano-Tárraga, C. *et al.* Epigenome-wide association study identifies TXNIP gene associated with type 2 diabetes mellitus and sustained hyperglycemia. *Hum. Mol. Genet.* **25**, 609–619 (2016).
- [24] Zannas, A. S. *et al.* Epigenetic upregulation of FKBP5 by aging and stress contributes to NF- $\kappa$ B–driven inflammation and cardiovascular risk (2019).
- [25] Kurushima, Y. *et al.* Epigenetic findings in periodontitis in UK twins: a cross-sectional study. *Clin. Epigenetics* **11**, 27 (2019).

- [26] Voisin, S. *et al.* Many obesity-associated SNPs strongly associate with DNA methylation changes at proximal promoters and enhancers. *Genome Med.* **7**, 103 (2015).
- [27] Tan, Q. *et al.* Epigenetic signature of birth weight discordance in adult twins. *BMC Genomics* **15**, 1062 (2014).
- [28] Tserel, L. *et al.* Age-related profiling of DNA methylation in CD8+ T cells reveals changes in immune response and transcriptional regulator genes. *Sci. Rep.* **5**, 13107 (2015).
- [29] Butcher, D. T. *et al.* CHARGE and kabuki syndromes: Gene-Specific DNA methylation signatures identify epigenetic mechanisms linking these clinically overlapping conditions. *Am. J. Hum. Genet.* **100**, 773–788 (2017).
- [30] Dabin, L. C. *et al.* Altered DNA methylation profiles in blood from patients with sporadic Creutzfeldt-Jakob disease. *Acta Neuropathol.* **140**, 863–879 (2020).
- [31] Marabita, F. *et al.* An evaluation of analysis pipelines for DNA methylation profiling using the illumina HumanMethylation450 BeadChip platform. *Epigenetics* **8**, 333–346 (2013).
- [32] del Valle J *et al.* GSE128064 (2019). Title of the publication associated with this dataset: Comprehensive constitutional genetic and epigenetic characterization of Lynch-like individuals.
- [33] Aryee, M. J. *et al.* Minfi: a flexible and comprehensive bioconductor package for the analysis of infinium DNA methylation microarrays. *Bioinformatics* **30**, 1363–1369 (2014).
- [34] Triche, T. J., Jr, Weisenberger, D. J., Van Den Berg, D., Laird, P. W. & Siegmund, K. D. Low-level processing of illumina infinium DNA methylation BeadArrays. *Nucleic Acids Res.* **41**, e90 (2013).
- [35] Houseman, E. A. *et al.* DNA methylation arrays as surrogate measures of cell mixture distribution. *BMC Bioinformatics* **13**, 86 (2012).
- [36] Dámaso, E. *et al.* Comprehensive constitutional genetic and epigenetic characterization of Lynch-Like individuals. *Cancers* **12** (2020).
- [37] Horvath, S. DNA methylation age of human tissues and cell types (2013).
- [38] ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57–74 (2012).
- [39] Davis, C. A. *et al.* The encyclopedia of DNA elements (ENCODE): data portal update. *Nucleic Acids Res.* **46**, D794–D801 (2018).

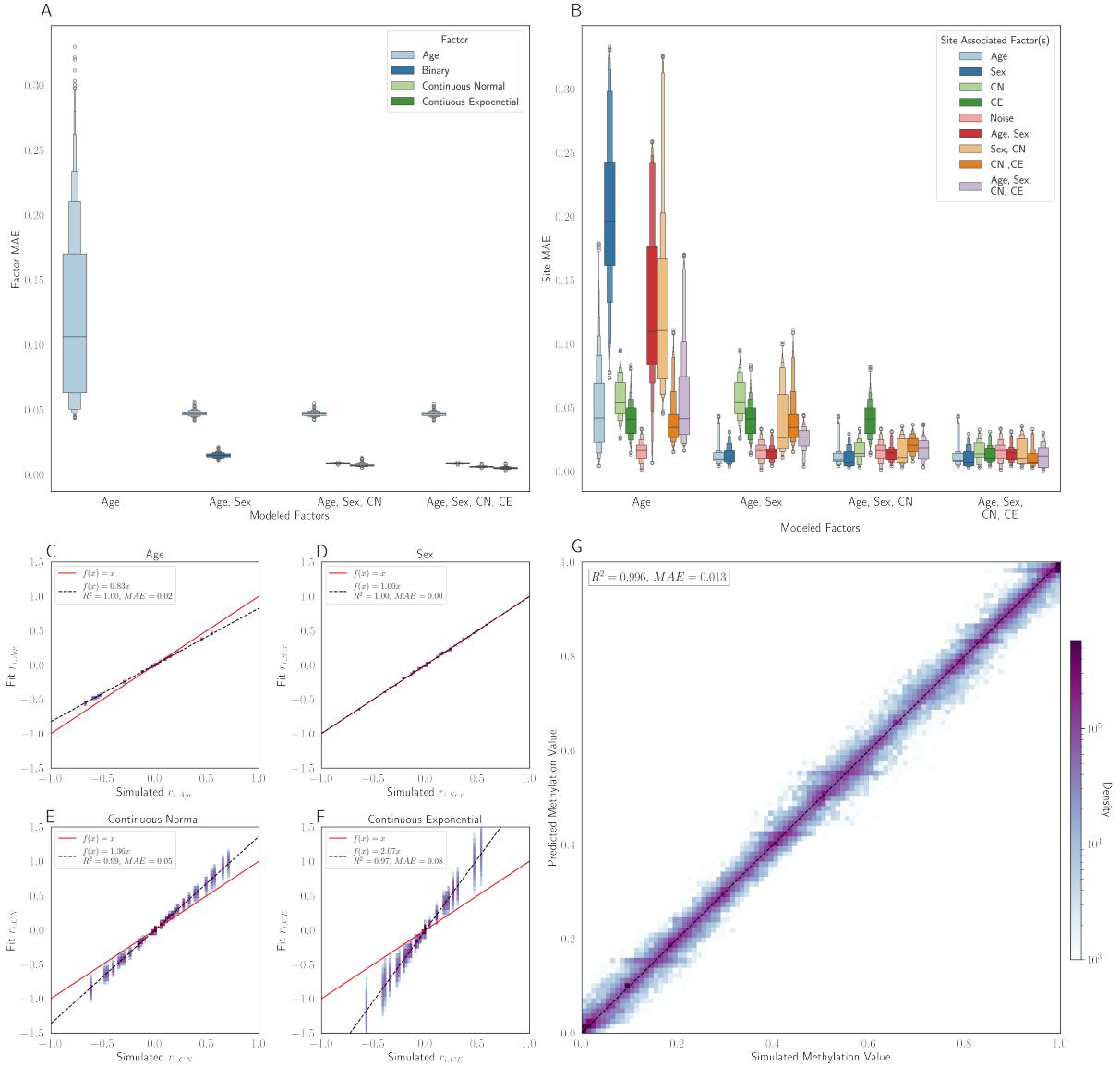
- [40] Hinrichs, A. S. *et al.* The UCSC genome browser database: update 2006. *Nucleic Acids Res.* **34**, D590–8 (2006).
- [41] Virtanen, P. *et al.* SciPy 1.0: fundamental algorithms for scientific computing in python. *Nat. Methods* (2020).
- [42] Sheffield, N. C. & Bock, C. LOLA: enrichment analysis for genomic region sets and regulatory elements in R and bioconductor. *Bioinformatics* **32**, 587–589 (2016).
- [43] Zheng, R. *et al.* Cistrome data browser: expanded datasets and new tools for gene regulatory analysis. *Nucleic Acids Res.* **47**, D729–D735 (2019).
- [44] Mei, S. *et al.* Cistrome data browser: a data portal for ChIP-Seq and chromatin accessibility data in human and mouse. *Nucleic Acids Res.* **45**, D658–D662 (2017).
- [45] Layer, R. M. *et al.* GIGGLE: a search engine for large-scale integrated genome analysis. *Nat. Methods* **15**, 123–126 (2018).
- [46] Bujold, D. *et al.* The international human epigenome consortium data portal. *Cell Syst* **3**, 496–499.e2 (2016).
- [47] Ramírez, F. *et al.* deeptools2: a next generation web server for deep-sequencing data analysis. *Nucleic Acids Res.* **44**, W160–5 (2016).
- [48] Robinson, J. T. *et al.* Integrative genomics viewer. *Nat. Biotechnol.* **29**, 24–26 (2011).
- [49] Basu, A. Reproducible research with jupyter notebooks.
- [50] Varoquaux, G. & Grisel, O. Joblib: running python function as pipeline jobs. *packages. python. org/joblib* (2009).
- [51] Virtanen, P. *et al.* SciPy 1.0: fundamental algorithms for scientific computing in python. *Nat. Methods* (2020).
- [52] Hunter, J. D. Matplotlib: A 2D graphics environment (2007).
- [53] Waskom, M. seaborn: statistical data visualization. *J. Open Source Softw.* **6**, 3021 (2021).
- [54] McKinney, W. *Python for Data Analysis: Data Wrangling with Pandas, NumPy, and IPython* (“O’Reilly Media, Inc.”, 2012).
- [55] da Costa-Luis, C. O. tqdm: A fast, extensible progress meter for python and CLI. *JOSS* **4**, 1277 (2019).

- [56] Barrett, T. *et al.* NCBI GEO: archive for functional genomics data sets—update. *Nucleic Acids Res.* **41**, D991–D995 (2012).
- [57] Lopes-Ramos, C. M. *et al.* Sex differences in gene expression and regulatory networks across 29 human tissues. *Cell Rep.* **31**, 107795 (2020).
- [58] Theilgaard-Mönch, K. *et al.* Transcription factor-driven coordination of cell cycle exit and lineage-specification in vivo during granulocytic differentiation : In memoriam professor niels borregaard. *Nat. Commun.* **13**, 3595 (2022).
- [59] Guerzoni, C. *et al.* Inducible activation of CEBPB, a gene negatively regulated by BCR/ABL, inhibits proliferation and promotes differentiation of BCR/ABL-expressing cells. *Blood* **107**, 4080–4089 (2006).
- [60] Tharakan, R. *et al.* Blood DNA methylation and aging: A Cross-Sectional analysis and longitudinal validation in the InCHIANTI study. *J. Gerontol. A Biol. Sci. Med. Sci.* **75**, 2051–2055 (2020).
- [61] Mohamed, E. A. & Sayed, W. M. Implication of JAK1/STAT3/SOCS3 pathway in aging of cerebellum of male rat: Histological and molecular study. *Sci. Rep.* **10**, 8840 (2020).
- [62] Trenerry, M. K., Carey, K. A., Ward, A. C., Farnfield, M. M. & Cameron-Smith, D. Exercise-induced activation of STAT3 signaling is increased with age. *Rejuvenation Res.* **11**, 717–724 (2008).
- [63] Carlson, M. E. *et al.* Relative roles of TGF- $\beta$ 1 and wnt in the systemic regulation and aging of satellite cell responses. *Aging Cell* **8**, 676–689 (2009).
- [64] Paris, N. D., Soroka, A., Klose, A., Liu, W. & Chakkalakal, J. V. Smad4 restricts differentiation to promote expansion of satellite cell derived progenitors during skeletal muscle regeneration. *Elife* **5** (2016).
- [65] Carrieri, G. *et al.* The G/C915 polymorphism of transforming growth factor beta1 is associated with human longevity: a study in italian centenarians. *Aging Cell* **3**, 443–448 (2004).
- [66] Chong, J. A. *et al.* REST: a mammalian silencer protein that restricts sodium channel gene expression to neurons. *Cell* **80**, 949–957 (1995).
- [67] Coulson, J. M. Transcriptional regulation: cancer, neurons and the REST. *Curr. Biol.* **15**, R665–8 (2005).
- [68] Lu, T. *et al.* REST and stress resistance in ageing and alzheimer’s disease (2014).
- [69] Rocchi, A. *et al.* REST/NRSF deficiency impairs autophagy and leads to cellular senescence in neurons. *Aging Cell* **20**, e13471 (2021).

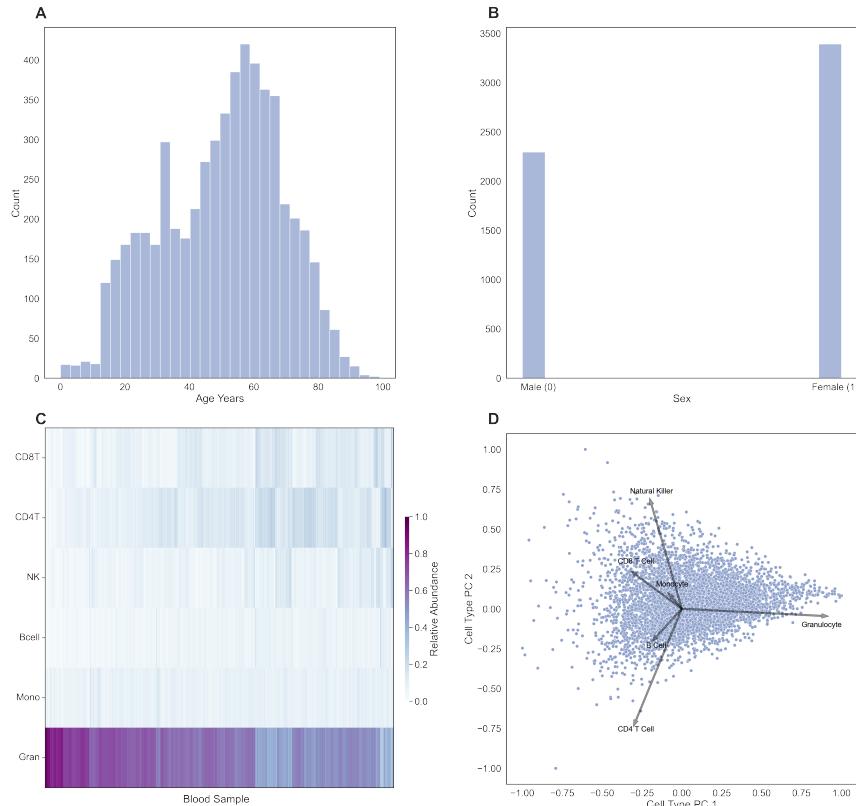
- [70] Xiong, Z. *et al.* EWAS data hub: a resource of DNA methylation array data and metadata. *Nucleic Acids Res.* **48**, D890–D895 (2020).
- [71] McCartney, D. L. *et al.* An epigenome-wide association study of sex-specific chronological ageing. *Genome Med.* **12**, 1 (2019).
- [72] de Lima Camillo, L. P., Lapierre, L. R. & Singh, R. A pan-tissue DNA-methylation epigenetic clock based on deep learning. *npj Aging* **8**, 1–15 (2022).
- [73] Han, Y. *et al.* New targeted approaches for epigenetic age predictions. *BMC Biol.* **18**, 71 (2020).
- [74] Schoenherr, C. J. & Anderson, D. J. The neuron-restrictive silencer factor (NRSF): a coordinate repressor of multiple neuron-specific genes. *Science* **267**, 1360–1363 (1995).
- [75] Ooi, L. & Wood, I. C. Chromatin crosstalk in development and disease: lessons from REST. *Nat. Rev. Genet.* **8**, 544–554 (2007).
- [76] Bruce, A. W. *et al.* Genome-wide analysis of repressor element 1 silencing transcription factor/neuron-restrictive silencing factor (REST/NRSF) target genes. *Proc. Natl. Acad. Sci. U. S. A.* **101**, 10458–10463 (2004).
- [77] Bernstein, B. E. *et al.* A bivalent chromatin structure marks key developmental genes in embryonic stem cells. *Cell* **125**, 315–326 (2006).
- [78] Voigt, P., Tee, W.-W. & Reinberg, D. A double take on bivalent promoters. *Genes Dev.* **27**, 1318–1338 (2013).
- [79] Ooi, S. K. T. *et al.* DNMT3L connects unmethylated lysine 4 of histone H3 to de novo methylation of DNA. *Nature* **448**, 714–717 (2007).



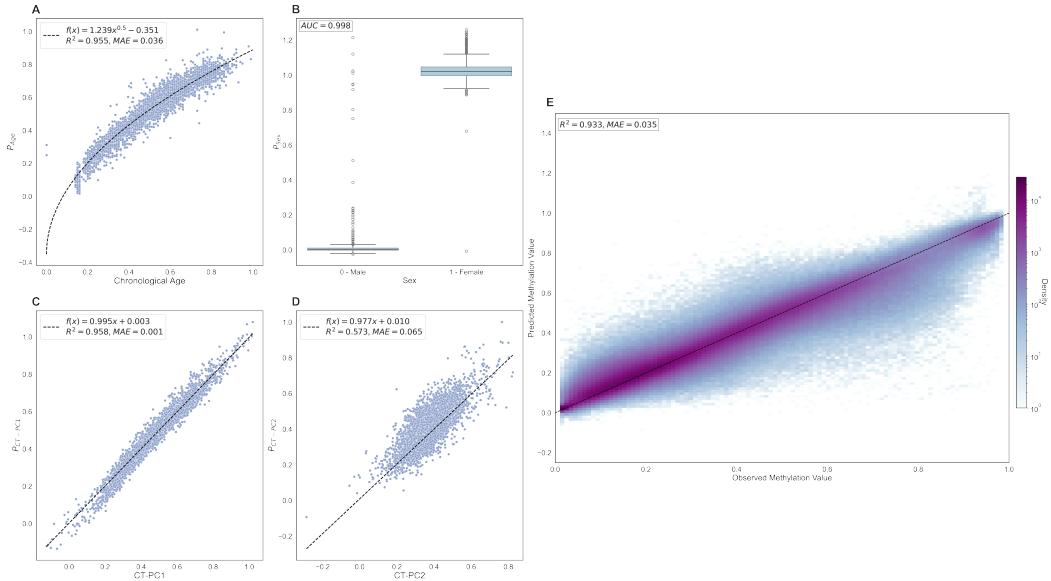
**Figure 1:** Simulated factors and the association with simulated methylation values. (A) Age with a non-linear association with methylation (B). Sex (C) with a binary association with methylation (D). Normal factor (E) with a linear relationship with methylation (F). Continuous exponential trait (G) 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  $\delta_i$ .



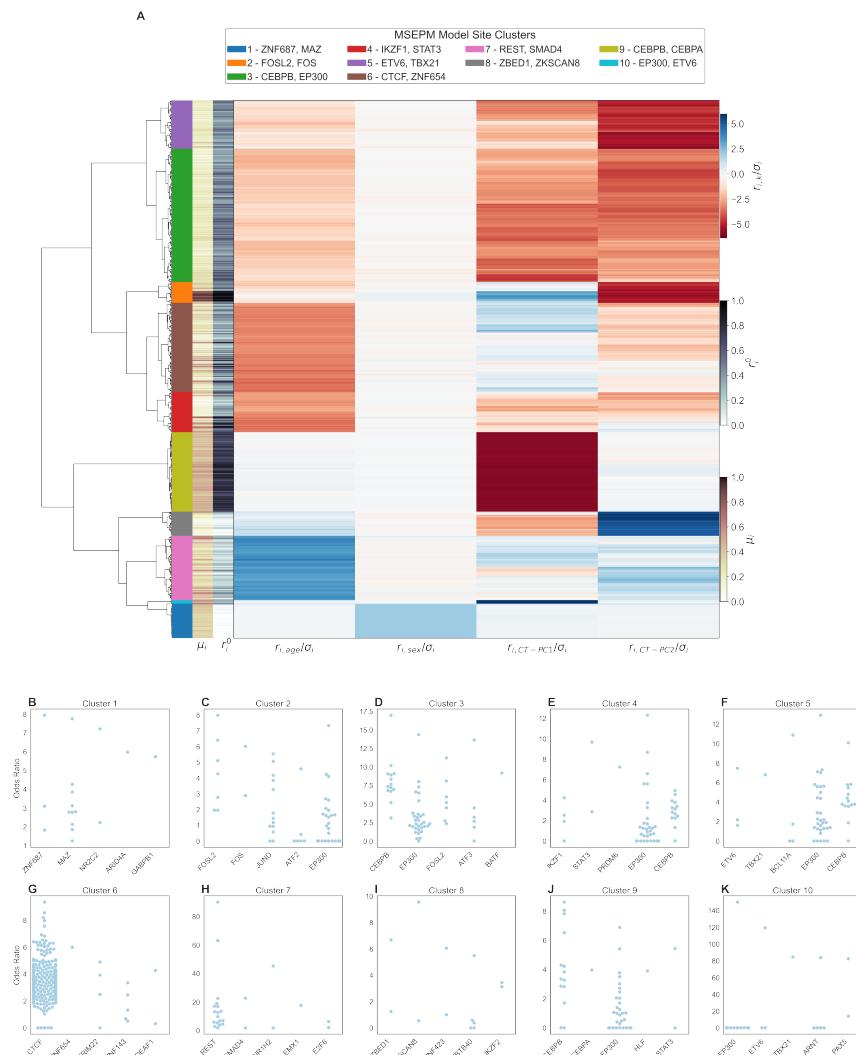
**Figure 2:** (A) The MAE of the factor predictions on the testing set as multiple factors are modeled simultaneously and (B) predicted methylation MAE as factors are included in the MSEPM model where the centerline is the 50th quantile and the box with greatest width contains 50% of the underlying data 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.

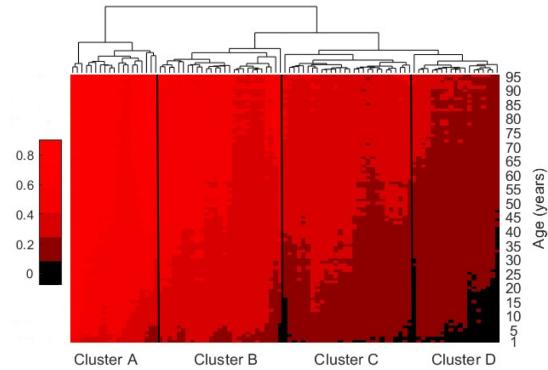


**Figure 4:** MSEPM model trained with age, sex, CT-PC1 and CT-PC2 predictions within testing set for epigenetic factors (A) age, (B) sex, (C) CT-PC1 and (D) CT-PC2. (E) Observed and predicted methylation values for training set has high concordance ( $R^2 = 0.933$ ,  $MAE = 0.035$ )

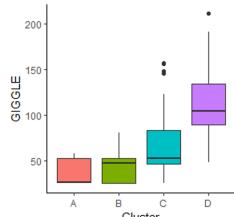


**Figure5:** (A) Site clustering by standardized model coefficients. Sites clusters show distinct relationships with modeled traits. (B-K) Top five enriched transcription factors for clusters 1 - 10.

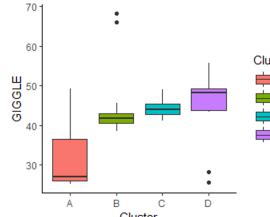
A Clustergram of DNA Methylation Across Ages and Sites



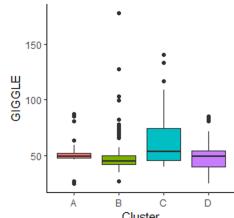
B REST



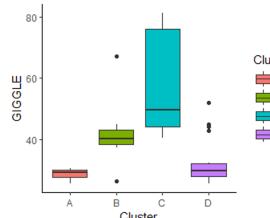
C H3K4me3



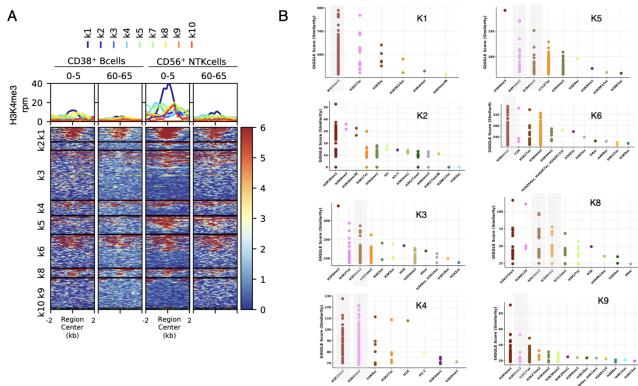
D H3K27me3



E H3K27ac



**Figure6:** (A) Heatmap of H3K4me3 ChIP-seq enrichment for two different blood cell types (CD38<sup>+</sup> B Cells and CD56<sup>+</sup> 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<sup>+</sup> B Cells and CD56<sup>+</sup> 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.