

1 **DNA methylation biomarker in blood predicts frailty**
2 **in an HIV-positive veteran population**

3
4 Chang Shu^{1,2}, Amy C. Justice^{1,2}, Xinyu Zhang^{1,2}, Vincent C. Marconi³, Dana B.
5 Hancock⁴, Eric O. Johnson^{4,5}, Ke Xu^{1,2}

6 1. Department of Psychiatry, Yale School of Medicine, New Haven, CT, 06516

7 2. Connecticut Veteran Healthcare System, West Haven, CT, 06515

8 3. Division of Infectious Disease, Emory University School of Medicine, Atlanta GA

9 4. Center for Omics Discovery and Epidemiology, Behavioral Health Research Division,

10 RTI International, Research Triangle Park, NC, USA

11 5. Fellow Program, RTI International, Research Triangle Park, NC, USA.

12
13 All corresponding address to

14 Ke Xu, M.D., Ph.D.

15 Associate Professor of Psychiatry

16 Yale School of Medicine

17 Email: ke.xu@yale.edu

18 Tel: 203-932-5711x7430

19 **Abstract**

20 **Background:** With the improved life expectancy in people living with HIV, predicting
21 individual frailty is important for clinical care. DNA methylation (DNAm) has emerged as
22 a robust biomarker in precision medicine and has been previously linked to aging and
23 mortality in non-HIV populations. Here, we aim to establish a panel of DNAm features
24 selected from blood methylome to predict frailty in HIV-positive men from a veteran
25 population, Veterans Aging Cohort Study (VACS).

26 **Methods:** We used a well-established score, VACS Index, as a measure of frailty.
27 Samples ($n_{\text{total}}=1,150$) were divided into a training set ($n=612$) and a validation set
28 ($n=538$). We first selected a panel of frailty-associated CpGs by conducting an
29 epigenome-wide association analysis on the VACS index in the training set. We then
30 applied four machine learning methods to build models to predict high and low frailty
31 individuals in the training set. The prediction models were tested in the validation set. A
32 methylation score constructed from the selected CpGs was tested an associated with
33 mortality by performing survival analysis. To assess the biological relevance of the
34 selected CpG sites, we performed a gene ontology enrichment analysis.

35 **Results:** A panel of 119 CpGs were identified from the training set (False Discovery
36 Rate <0.2) that showed excellent performance on predicting high frailty individuals with
37 Area Under Curve (AUC) of 0.835 (95%CI: 0.792, 0.879) and balanced accuracy of
38 0.693. The same panel showed good performance on predicting low frailty individuals
39 with AUC of 0.735 (95%CI: 0.688-0.782) and a balanced accuracy of 0.528. The
40 methylation score from the 119 CpGs was significantly associated with 5-year and 10-
41 year mortality with hazard ratio of 1.40 (95% CI:1.033, 1.897 $p=0.03$) and 1.48 (95%CI:

42 1.10, 2.02; $p=0.01$) respectively. These 119 CpGs were located within or near 73 genes
43 that were significantly enriched in 9 biological pathways relevant to immune and
44 inflammation response.

45 **Conclusions:** We identified a panel of predictive DNAm features associated with frailty
46 and mortality among HIV-positive population. These DNAm features may serve as
47 biomarkers to discriminate high and low frailty people who live with HIV.

48

49 **Keywords:** DNA methylation, HIV-positive, frailty, mortality, DNA methylation score,
50 machine learning prediction

51

52

53 **Background**

54 Combination antiretroviral therapy has significantly improved life expectancy
55 among HIV-positive individuals. The prolonged life expectancy increases frailty risk in
56 the HIV population. The prevalence of frailty in people living with HIV is significantly
57 higher and onset of frailty occurs at an earlier age compared to the general population
58 (1). High frailty is characterized by marked vulnerability and is associated with increased
59 mortality. Thus, prediction of frailty is important to identify vulnerable group and deliver
60 clinical care to highly vulnerable HIV-positive patients. At present, frailty is usually
61 defined by clinical symptoms (e.g., Fried's frailty phenotype, Rockwood and Mitnitski's
62 Frailty Index) or a combination of lab tests that indicate organ damage such as the
63 Veterans Aging Cohort Study index (VACS index) (2). No biomarkers available to
64 capture the early stage of the frailty to predict individual vulnerability in HIV-positive
65 patients.

66 A large body of evidence has demonstrated that epigenetic modification adapts
67 internal and external environmental changes and that is associated with early stage of
68 pathophysiological processes (3-6). DNAm, one type of the most widely studied
69 epigenetic marks, is strongly related with aging (7-9), substance use (e.g. cigarette
70 smoking and alcohol consumption) (10-16), and a variety of diseases (3-6, 17, 18).
71 Because DNAm is relatively stable and easy to detect from body fluids biospecimen
72 through non-invasive procedure, DNAm marks have emerged as robust biomarkers for
73 cancer diagnosis (19), disease subtype classification (20, 21) and treatment response
74 monitoring (22, 23).

75 DNAm may play an important role in frailty among HIV-positive individuals. Frailty
76 is associated with elevated inflammation markers such as IL-6 and hsCRP in HIV-
77 positive individuals (24). Genes involved in immune and inflammation processes are
78 also subject to epigenetic modification in immune cells. We previously reported
79 association of two CpG sites in the promoter region of *NLRC5* with HIV infection (25),
80 and *NLRC5* gene is a major transcriptional activator of MHC class I gene. DNAm was
81 also linked to HIV comorbid medical diseases such as HIV-positive diabetes and kidney
82 function (26, 27). Furthermore, aging, a critical contributor to frailty, is significantly
83 associated with thousands of CpGs in the epigenome, and the epigenetic clock and
84 DNAm age are becoming widely recognized (7-9). DNAm age is significantly correlated
85 with physical frailty in HIV-negative individuals (28) and HIV-positive individuals (29). In
86 HIV-positive individuals, the average DNAm age is accelerated 5 to 10 years faster than
87 HIV-negative individuals (30-33) and 10 years faster in the HIV-positive frailty
88 individuals in comparing to HIV-positive non-frailty individuals (29). These age-related
89 DNAm signatures are predictive of mortality (34-38). Additional CpGs in blood have
90 been identified to predict all-cause mortality (34, 39). Therefore, we hypothesize that
91 DNAm is associated with frailty and that DNAm signatures in blood can serve as
92 biomarkers to predict frailty and mortality among HIV-positive individuals.

93 In this study, our goal is to identify DNAm marks that can serve as biomarkers
94 on frailty and to link frailty-associated DNAm to mortality among HIV-positive
95 individuals. We also evaluated biological relevance of the identified CpG sites. We used
96 a well-established frailty score, VACS index, as a measure of frailty in the HIV-positive
97 population. The VACS index is a composite score constructed by a sum of pre-assigned

98 points on age, CD4 count, HIV-1 RNA, hemoglobin, platelets, aspartate and alanine
99 transaminase (AST and ALT), creatine, estimated glomerular filtration rate (eGFR), and
100 viral hepatitis C infection (40). Findings from this study provide a set of DNAm
101 biomarkers to predict frailty for future clinic use and new insights into the epigenetic
102 mechanism of frailty and mortality in HIV pathology.

103

104 **Methods**

105 As an overview, our analytical procedure is shown in the flowchart in **Figure 1**.
106 Our sample was divided into the training set (n=612) and the validation set (n=538)
107 which DNA methylation were processed at different time and using two different
108 platforms. We first selected a panel of CpG sites relevant to frailty based on EWAS
109 results in the training set. Then, we applied four commonly used machine learning
110 classification methods to develop prediction models: random forest, Extreme Gradient
111 Boosting Tree (XGBoost), Lasso and Elastic-Net Regularized Generalized Linear
112 Models (GLMNET), and Support Vector Machines (SVM). The selected CpG sites were
113 used as predictors to differentiate high and low frailty in the training set and each model
114 was evaluated in the validation set. Additionally, the DNAm score was constructed
115 based on the selected CpG sites and we conducted survival analysis to assess whether
116 the DNAm score was associated with mortality in the HIV-positive samples. Lastly, we
117 conducted a gene ontology enrichment analysis to reveal the underlying biological
118 pathways of the selected CpG sites.

119 ***Study population***

120 All samples in the training and validating sets were from the VACS cohort. The
121 VACS is a prospective cohort study of veterans focusing on clinical outcomes in HIV
122 infection (41). DNA samples were extracted from peripheral blood of 1,150 HIV-positive
123 men from the VACS. Demographic and clinical information on age, race, smoking
124 status, CD4 count, viral load, HIV medication adherence, VACS index, and mortality by
125 training set and validation set are presented in **Table 1**. The training set included slightly
126 older individuals, more African Americans, and greater VACS index than the validation
127 set. There was no significant difference in HIV medication adherence, CD4 count, or
128 \log_{10} HIV-1 viral load between the training set and the validation set.

Table 1: Study sample characteristics

	Training set (N=612)	Validation set (N=538)	p value*
Age (year, mean +/-sd)	49.4 (7.6)	48.1 (7.8)	0.005
Male (%)	612 (100)	538 (100)	
Race (%)			
Caucasian	58 (9.5)	48 (8.9)	0.001
African Americans	526 (85.9)	435 (80.9)	
Other	28 (4.6)	55 (10.2)	
Smokers (%)	360 (59.4)	309 (58.4)	0.78
HIV positive (%)	614 (100)	538 (100.0)	
HIV treatment adherence (%)	469 (78.4)	407 (76.6)	0.519
CD4 count	428.97 (286.33)	447.46 (279.27)	0.287
log 10 HIV-1 viral load	2.61 (1.19)	2.67 (1.23)	0.423
VACS index (mean +/-)	35.64 (21.99)	30.78 (20.23)	<0.001
Mortality (%)	162 (26.5)	129 (24.0)	0.367

*t test is used for continuous variables, chi-square test is used for categorical variables

129

130 ***Genome-wide DNAm profiling and quality control***

131 The DNA samples in the training set were profiled by Infinium Human Methylation 450K
132 BeadChip (HM450K) and the DNA samples in the validation set were profiled by the

133 Infinium Human Methylation EPIC BeadChip. DNAm for the training and validation sets
134 were evaluated using the same quality control (QC) protocol (42) in the R package minfi
135 (43). In detail, CpG sites on sex chromosomes and within 10 base pairs of a single
136 nucleotide polymorphism were removed. The detection p-value threshold was set at 10^{-12}
137 for both the training and validation sets. After QC, 408,583 shared CpG sites between
138 HM450K and EPIC arrays were used for analysis to ensure the same coverage
139 between two sets. Proportions of 6 cell types (CD4+ T cells, CD8+ T cells, Natural Killer
140 T cells, B cells, monocytes and granulocytes) were estimated using the established
141 method (Houseman et al., 2012) for all samples in the training and validation sets.

142 ***Selection of CpG sites from EWAS on frailty in the training set***

143 We performed an EWAS on VACS index in the training set using a two-step
144 linear model approach as previously described (42). The analytical model was adjusted
145 for systematic errors and clinical confounding factors. Here, log transformation was
146 applied to VACS index to ensure normality distribution assumption for the linear model.

147 First, the following linear model was used to extract principal components (PCs)
148 for potential confounding variables. Here, top 30 control PCs were derived from internal
149 control probes in *minfi*:

150 *Methylation β ~ age + race + smoking + self-reported HIV medication adherence*
151 *+ log viral load + WBC + CD4T + CD8T + Gran + NK + B cell + Mono + 30 Control PCs*

152 Second, the top 20 residual PCs were extracted in the model above to fully
153 adjust for unmeasured confounding. The final EWAS model on VACS index was:

154 $\log(\text{VACS}) \sim \text{Methylation } \beta + \text{age} + \text{race} + \text{smoking} + \text{self-reported HIV}$
155 $\text{medication adherence} + \log \text{viral load} + \text{WBC} + \text{CD4T} + \text{CD8T} + \text{Gran} + \text{NK} + \text{B cell} +$
156 $\text{Mono} + 30 \text{ Control PCs} + 20 \text{ Residual PCs}$

157 To ensure a sufficient number of CpGs as predictors for the prediction model
158 development, we selected CpG sites with false discovery rate (FDR) <0.2 in the EWAS
159 as predictors.

160 Additionally, to explore the underlying epigenetic mechanism of frailty, we also
161 detected differentially methylated regions (DMRs) using *bumpHunter* (44) in the training
162 set. DMRs with family-wise error rate (FWER)<0.2, which is equivalent to FDR<0.03,
163 were considered significant.

164 ***Developing machine learning prediction models for frailty***

165 Based on the distribution of the VACS index in the entire sample (**Figure S1**) and
166 clinical significance (45), high HIV frailty was defined as VACS index > 50, and low HIV
167 frailty was defined as VACS index < 16. Models were developed to predict high frailty
168 (VACS index >50 versus ≤ 50) and low frailty (VACS index <16 versus ≥ 16). CpGs
169 relevant to frailty were selected from EWAS as described above. The following steps
170 were taken to develop separate prediction models on high and low frailty:

171 1) Model development in the training set: Four machine learning models, random
172 forest, Extreme Gradient Boosting Tree (XGBoost), support vector machines(SVM),
173 Lasso and Elastic-Net Regularized Generalized Linear Models (GLMNET), were
174 separately applied to predict VACS index by using the R package *caret* (46). These four
175 models are commonly used in the supervised learning on classification (47-50). Each

176 model was developed separately and 10-fold cross-validation was used in model
177 training process to minimize data overfitting.

178 2) Model evaluation in the validation set: All four models built in the training set
179 were evaluated in the validation set. Area Under Curve (AUC) and balanced accuracy
180 were used to assess the prediction performance. We used balanced accuracy for model
181 evaluation because in the presence of imbalance samples in each class since there are
182 228 subjects with VACS index greater than or equal to 50 and 916 subjects with VACS
183 index less than 916. Balanced accuracy is defined as the average accuracy obtained on
184 each class as shown in the following formula (51). Balanced accuracy was used in this
185 study to avoid biased accuracy due to imbalanced samples (51).

186

187 *Balanced accuracy*

$$188 = \frac{1}{2} \left(\frac{\text{True positive}}{\text{True positive} + \text{False negative}} + \frac{\text{True negative}}{\text{True negative} + \text{False positive}} \right)$$

189

190 ***Association of a DNAm score for VACS index with mortality***

191 We constructed a DNAm score by the selected frailty-associated CpG sites. A
192 survival analysis was conducted to assess whether the DNAm score was associated
193 with 5-year and 10-year mortality respectively.

194 DNAm score was constructed based on a previously reported method (12). We
195 computed the DNAm score by normalizing the score with mean β value (μ_c) and
196 standard deviation, σ_c , from subjects with VACS index <50 (control group). For selected
197 n methylation sites, W_c is 1 for a hypermethylated methylation site c , and -1 for a

198 hypomethylated site c , β_{ci} is the β value for subject i and methylation site c . The
199 methylation score for subject i was constructed by the following formula:

$$200 \quad \text{Methylation Score}_i = \frac{1}{n} \sum_{c=1}^n W_c \frac{\beta_{ci} - \mu_c}{\sigma_c}$$

201 The mean methylation score was 0, and samples were divided into high
202 methylation score (>0) and low methylation score (≤ 0) groups.

203 We performed Kaplan-Meier survival curves during 10-year follow-up to visualize
204 the survival differences between high and low methylation score groups. Survival
205 analysis was conducted by cox proportional hazards model on 5-year and 10-year
206 mortality comparing high and low methylation groups. We used age as time scale t , and
207 our model was adjusted for race, smoking, self-reported HIV medication adherence, log
208 10 of HIV viral load and CD4 count.

$$209 \quad h(t) = h_0(t) \exp(\beta_1 \text{ methylation score} + \beta_2 \text{ race} + \beta_3 \text{ smoking} \\ 210 \quad + \beta_4 \text{ HIV medication adherence} + \beta_5 \log_{10}(\text{viral load}) + \beta_6 \text{ CD4 count})$$

213 ***Biological interpretation of the predictive panel of CpG sites on frailty***

214 We performed gene ontology (GO) enrichment analysis using Database for
215 Annotation, Visualization and Integrated Discovery (DAVID) (52). To avoid redundancy
216 in pathway names, we only used level 4 GO terms defined in DAVID in the enrichment
217 analysis. Genes with at least one frailty-associated CpG site was used for GO analysis.
218 We considered biological pathways with FDR <0.05 as statically significant pathways.

219

220 **Results**

221 ***Selection a panel of 119 frailty-associated CpG sites by EWAS in the training set***

222 In the training set, we conducted an EWAS on VACS index (**Figure 2**; $\lambda=1.038$) and
223 selected 119 CpGs with $FDR < 0.2$ to ensure sufficient number of predictors (**Table 2**).
224 Eighty out of 119 CpGs were negatively associated with VACS index, while 39 out of
225 119 CpGs were positively associated with VACS index. The majority of CpGs were
226 located within or near known genes, except for 19 CpGs that were intergenic. Only 40
227 CpGs were located in gene bodies, while 79 CpGs were located in promoter regions,
228 first exons, or 3' UTR. Thirty-two CpGs were located in CpG islands.
229 Notably, 22 out of 119 CpGs reached epigenome-wide significance ($FDR < 0.05$),
230 including 16 CpG sites negatively associated with the VACS index and 6 CpG sites with
231 positive associations. These 22 CpGs were located on 17 genes, including 13 of 22
232 CpGs located in promoter regions, 6 CpGs in gene body, and 3 CpGs in 3'UTR. Among
233 these genes, some harbored more than one significant CpG. For example, three CpGs
234 were annotated within or near *PSMB8* (cg01309328, $p=4.44 \times 10^{-9}$; cg08099136,
235 $p=4.63 \times 10^{-9}$; cg00533183, $p=3.18 \times 10^{-7}$), two CpGs were located near *PARP9*
236 (cg08122652, $p=2.18 \times 10^{-7}$; cg22930808, $p=6.78 \times 10^{-7}$), two CpG sites were located
237 near *IFITM1* (cg03038262, $p=5.83 \times 10^{-7}$; cg23570810, $p=1.58 \times 10^{-6}$), and the rest of the
238 significant CpG sites were annotated to 13 genes including *MX1*, *TAP1*, *ZNF32*,
239 *NLRC5*, *IFI44L*.

240 We found 9 significant DMRs in the training set (**Table S1**). Methylation β values
241 between high and low frailty ranged from -0.15 to 0.05 in these regions. As examples, 3

242 region plots showing DMRs in the *MX1*, *PARP9*, and *IFI44L* genes are shown in **Figure**

243 **3**.

244 ***Machine learning prediction on frailty by the selected 119 CpGs***

Table 3: Performance of machine learning model predicting frailty in an HIV-infected population

Method	High frailty (VACS index >50 vs. ≤50)				Low frailty (VACS index >16 vs. ≤16)			
	Training		Validation		Training		Validation	
	AUC	Balanced Accuracy	AUC	Balanced Accuracy	AUC	Balanced Accuracy	AUC	Balanced Accuracy
Random Forest	0.775	0.622	0.835 (0.792,0.879)	0.693	0.774	0.509	0.735 (0.688,0.782)	0.528
XGBoost	0.760	0.624	0.830 (0.786,0.874)	0.684	0.820	0.637	0.716 (0.669,0.762)	0.500
GLMNET	0.783	0.642	0.757 (0.701,0.813)	0.668	0.795	0.589	0.715 (0.667,0.763)	0.503
SVM	0.759	0.618	0.744 (0.685,0.802)	0.500	0.756	0.594	0.501 (0.442,0.560)	0.492

AUC: Area Under Curve

VACS: Veteran Aging Cohort Study

XGBoost: Extreme Gradient Boosting Tree

GLMNET: Lasso and Elastic-Net Regularized Generalized Linear Models

SVM: Support Vector Machines

245

246 In the training set, we developed our prediction models based on the selected
 247 119 CpGs in the training set using four commonly used machine learning classification
 248 models: XGBoost, random forest, SVM and GLMNET (47-50). Performances of four
 249 prediction models are presented in **Table 3**. We found that AUCs ranged from 0.756 to
 250 0.820 and balanced accuracies ranged from 0.509 to 0.642, suggesting that the 119
 251 CpGs had good to excellent predictive performance on frailty. The performances of four
 252 models were mostly comparable. For predicting high frailty, GLMNET outperformed the
 253 three other models (AUC=0.783, balanced accuracy=0.642) while SVM performed
 254 slightly worse (AUC=0.759, balanced accuracy = 0.618). For the prediction of low frailty,

255 XGBoost performed the best among the four models (AUC=0.820, balanced
256 accuracy=0.637) while SVM performed worst (AUC=0.756, balanced accuracy =0.594).

257 In the validation set, for the prediction of high frailty, random forest and XGBoost
258 showed the best performance. AUC of random forest model was 0.835 (95%CI: 0.792,
259 0.879) and XGBoost was 0.830 (95%CI: 0.786, 0.874) for predicting high frailty.
260 Balanced accuracy was 0.69 and 0.68, respectively, for random forest and XGBoost
261 (**Figure 4**). For the prediction of low frailty, XGBoost and random forest models also
262 outperformed GLMNET and SVM. AUC was 0.735 (95%CI: 0.688, 0.782) for the
263 random forest and 0.716 (95%CI: 0.669, 0.762) for the XGBoost. The balanced accuracy
264 was 0.528 and 0.500 for the random forest and XGBoost, respectively.

265 These results suggest that the selected panel of 119 CpG sites were able to
266 predict frailty, and the prediction models performed better at predicting high frailty than
267 predicting low frailty. Among four machine learning methods, SVM showed poor
268 performance of predicting frailty in our samples.

269

270 ***DNAm score by 119 CpGs was significant associated with mortality***

271 We constructed DNAm score based on the selected 119 CpGs, and we further
272 assessed whether VACS index-related DNAm score were associated with mortality. In
273 **Figure 5**, the Kaplan Meier curves during 10-year follow-up showed that the individuals
274 with high DNAm score were at higher risk of mortality than those with low DNAm
275 scores. After adjusting for confounding factors, our cox proportional hazards regression
276 model showed that the hazard ratio of 5-year and 10-year mortality comparing high and
277 low DNAm score groups were 1.40 (95%CI: 1.03-1.90, p=0.03) and 1.48 (95%CI: 1.10-

278 2.02, $p=0.01$) respectively. In conclusion, we found that the DNAm score constructed
 279 from the 119 selected CpGs was significantly associated with mortality.

280 **Biological interpretations of 119 CpGs by gene ontology enrichment analysis**

281 The selected panel of 119 CpGs were located within or near 73 genes. Gene
 282 ontology enrichment analysis using these 73 genes resulted in 9 significant pathways
 283 with $FDR < 0.05$ (**Figure 6, Table 4**). These pathways included response to type I
 284 interferon response ($FDR=1.05 \times 10^{-7}$), innate immune response ($FDR=3.45 \times 10^{-5}$),
 285 cytokine response ($FDR=3.69 \times 10^{-3}$) and defense response to virus (**Figure 6**). Our
 286 findings suggested that the selected 119 CpG sites are biologically relevant to HIV
 287 pathogenesis and progression.

Table 4: Gene ontology term enrichment analysis of the selected CpG sites that predict frailty in an HIV infected population

Term	N	%	Genes	Fold Enrichment	P value	FDR
GO:0034340~response to type I interferon	10	11.8	<i>OAS2, IFITM1, IFIT3, MX1, NLRC5, HLA-F, XAF1, HLA-E, IRF7, PSMB8</i>	28.1	6.62E-11	1.05E-07
GO:0045087~innate immune response	19	22.4	<i>TRIM22, MAP4K2, IFITM1, TRIM25, MX1, GFI1, IFIT5, XAF1, LY86, PSMB8, TRIM14, PARP9, PLSCR1, OAS2, IFIT3, NLRC5, HLA-F, HLA-E, IRF7</i>	5.0	2.17E-08	3.45E-05
GO:0051607~defense response to virus	11	12.9	<i>TRIM22, PLSCR1, OAS2, IFITM1, TRIM25, IFIT3, MX1, IFI44L, NLRC5, IFIT5, IRF7</i>	10.5	7.55E-08	1.20E-04
GO:0051707~response to other organism	18	21.2	<i>TRIM22, PDE4B, EPO, IFITM1, NOTCH1, TRIM25, MX1, GFI1, IFIT5, IFI44L, LY86, PLSCR1, OAS2, IFIT3,</i>	4.8	1.17E-07	1.85E-04

							<i>CCDC88B, NLRC5, HLA-E, IRF7</i>
							<i>TRIM22, IFITM1, TFRC, NOTCH1, TRIM25, MX1, PPIB, GFI1, IFIT5, IFI44L, TAP1, PSMB8, TRIM14, PLSCR1, OAS2, SLC1A5, IFIT3, NLRC5, IRF7</i>
GO:0016032~viral process	19	22.4		4.3	1.93E-07	3.07E-04	
							<i>TRIM22, PLSCR1, OAS2, IFITM1, TRIM25, IFIT3, MX1, IFI44L, NLRC5, IFIT5, IRF7</i>
GO:0009615~response to virus	11	12.9		7.7	1.27E-06	2.02E-03	
							<i>TRIM22, PLSCR1, OAS2, IFITM1, TRIM25, IFIT3, MX1, CCDC88B, IFIT5, NLRC5, IFI44L, HLA-E, IRF7</i>
GO:0098542~defense response to other organism	13	15.3		5.8	1.94E-06	3.09E-03	
							<i>TRIM22, EPO, IFITM1, TRIM25, MX1, GFI1, XAF1, PSMB8, PARP9, PLSCR1, OAS2, IFIT3, NLRC5, HLA-F, HLA-E, IRF7</i>
GO:0034097~response to cytokine	16	18.8		4.4	2.32E-06	3.69E-03	
							<i>TRIM22, OAS2, IFITM1, TRIM25, NLRC5, HLA-F, HLA-E, IRF7</i>
GO:0034341~response to interferon-gamma	8	9.4		10.9	7.95E-06	1.26E-02	

288

289 Discussion

290 In this study, we present evidences that DNAm marks in blood are predictive of
 291 frailty and are associated mortality in an HIV-positive population. We identified a panel
 292 of 119 CpG sites that were highly predictive for high frailty and moderately predictive for
 293 low frailty. We also found that the DNAm score constructed by these 119 CpGs was
 294 strongly associated with mortality in the HIV-positive population. More importantly, these

295 clinically informative 119 DNAm features lied in genes involved in HIV pathogenesis and
296 progression. Thus, we discovered a panel of 119 DNAm biomarkers that add knowledge
297 to the epigenetic mechanisms underlying frailty and mortality among people living with
298 HIV.

299 We demonstrated the utility of using DNAm marks to predict frailty in HIV-positive
300 individuals. We took the several steps to avoid overfitting in developing the prediction
301 models: 1) model development and evaluation were conducted separately in training set
302 and validation set. DNAm in two sets were profiled in different time with different
303 platforms; 2) 10-fold cross validation were performed during training in each model. The
304 performances of four machine learning methods were in general consistent except that
305 the SVM did not perform well for low frailty. Our results suggested that this panel of
306 CpG sites was relatively stable and robust although performances of predictive models
307 differed using different methods. Of note, SVM showed the worst performance in
308 predicting high and low frailty that highlights the importance of choosing appropriate
309 machine learning method for model development. Compared to our previously reported
310 a panel of 698 smoking-associated CpGs that are moderately predictive of frailty (15),
311 the panel of 119 CpG sites in the present study includes fewer CpGs (119 CpGs) and
312 shows greater prediction performance on frailty (AUC=0.73 for smoking-associated
313 CpGs and AUC 0.83 for the VACS index-associated CpGs), suggesting that the panel
314 of 119 CpG sites has greater clinical utility for frailty among people living with HIV. The
315 improvement of prediction performance on frailty in the present study may be due to
316 different strategies of CpG selections. In this study, the 119 CpGs are selected from the

317 entire blood methylome and may be more informative than the smoking-associated
318 CpGs in our previous report.

319 We also found that the VACS index-related DNAm score derived from the 119
320 CpGs is strongly associated with mortality. This result provides further knowledge of
321 epigenetic profile into previous report that VACS index is predictive of mortality in
322 people who live with HIV (53). The result is also consistent with previous literature
323 showing that DNAm marks in blood predict mortality in the general population (39).
324 Interestingly, we found no overlapped CpG between our predictive panel and the
325 previously identified CpGs for all-cause mortality that included CpGs associated with a
326 variety of diseases such as diabetes and cancer(39). The discrepancy suggest that
327 methylation-based prediction of mortality may be relevant to causal disease. A recent
328 study reports a significant overlap of mortality-associated CpGs and aging-associated
329 CpGs (Lund et al, 2019). We found no overlapped CpG site between our mortality-
330 related 119 CpG panel and aging-related 353 CpGs in epigenetic clock (7), suggesting
331 that epigenetic mechanisms in HIV-related mortality differ from aging-related mortality.

332 Most importantly, this panel of 119 CpG sites are biologically meaningful and
333 may shed light on our understanding of the biological mechanisms of frailty for HIV-
334 positive individuals. The majority of the 119 CpG sites were located within or near
335 genes that are involved in known HIV pathology and progression. The results from DMR
336 analysis are corresponding to some of the significant genes for frailty. For example,
337 cg26312951, cg22862003 and cg21549285 from the 119 CpG sites are located in the
338 *MX1* gene (Interferon-Induced GTP-Binding Protein Mx1) and were negatively
339 associated with frailty. A DMR on *MX1* also showed a significant with frailty. *MX1*

340 encodes a GTP-metabolizing protein induced by interferon I and II and is involved in
341 interferon gamma signaling and Toll-like signaling pathway. Multiple CpGs near the
342 *HLA-F* and *HLA-E* genes are part of the predictive panel, and these two genes are
343 actively involved in immune response (54). Our previously reported HIV-associated
344 CpG site, cg07839457 in the *NLRC5* gene is also a member of the predictive panel on
345 frailty (25). Other interesting gene regions revealed by DMR analysis include *PARP9* and
346 *IFI44L*, which both genes are involved in *HIV pathogenesis*. The biological relevance of
347 these 119 CpG sites was further supported by the gene ontology enrichment analysis.
348 The top enriched pathways such as type I interferon response and cytokine response
349 may point out important biological pathways that leads to frailty and increased risk of
350 mortality among people living with HIV.

351 We acknowledge several limitations in this study. The generalizability of the
352 selected CpG sites may be limited since our samples were predominantly middle-aged
353 men, which may not generalize to frailty in a different age group. All samples in our
354 study are HIV-positive and the identified CpGs have a limited application to HIV-
355 negative population as we discussed above. Future studies in diverse populations is
356 warranted to validate the selected methylation features. Lastly, the prediction of low
357 frailty is moderate due to imbalanced sample distribution in our present study. We
358 expect that including more samples with low frailty will improve the predictive
359 performance.

360

361 **Conclusions**

362 We identified a panel of 119 predictive DNAm features in blood that are associated with
363 frailty and mortality among people living with HIV. These DNAm features may serve as
364 biomarkers to discriminate high and low frailty groups and may help to prioritize genes
365 to better understand the mechanisms of frailty in HIV-positive population. These DNAm
366 features have potential to serve as biomarkers to monitor HIV progression in future
367 clinical care.
368

369 **Legends**

370 **Figure 1:** Flowchart of analytical procedures on selection of CpG sites in peripheral
371 blood methylome, machine learning prediction models to predict frailty, DNA
372 methylation score to assess mortality, and gene ontology enrichment analysis

373 **Figure 2:** Manhattan and quantile-quantile plots on Veteran Aging Cohort Study (VACS)
374 index ($\lambda=1.038$) in the training set; blue line indicates False Discovery Rate (FDR)=0.05

375 **Figure 3:** Differentially methylation regions (DMRs) are associated with frailty,
376 measured by Veteran Aging Cohort Study (VACS) index. a: *MX1* regional plot showing
377 methylation beta value among subjects with VACS index > 50 and VACS index \leq 50 in
378 the training set. b: *IFI11L* regional plot showing methylation beta value among subjects
379 with VACS index > 50 and VACS index \leq 50 in the training set. c: *PARP9* regional plot
380 showing methylation beta value among subjects with VACS index > 50 and VACS index
381 \leq 50 in the training set

382 **Figure 4:** Receiver operating characteristic curve by XGBoost prediction model on HIV
383 frailty; a) Predicting high HIV frailty (VACS index >50 vs. \leq 50), Area Under Curve
384 (AUC): 0.830 (95% CI: 0.786,0.874); b) Predicting low HIV frailty (VACS index >16 vs.
385 \leq 16), AUC: 0.735 (95% CI: 0.688,0.782)

386 **Figure 5:** Kaplan-Meier curves comparing high and low HIV frailty methylation score
387 groups. Methylation score is calculated by 119 selected probes. High methylation score
388 group has methylation score \geq 0, while low methylation score group <0.

389 **Figure 6:** Gene ontology enrichment analysis of 119 CpG probes (FDR<0.05) for
390 prediction

391 **Figure S1:** Distribution of Veteran Aging Cohort Study index (VACS index) between
392 training set and testing set

393

394 **Table 1:** Study sample characteristics

395 **Table 2:** A panel of 119 CpG sites in blood that predicts high and low frailty in a HIV-
396 positive veteran population

397 **Table 3:** Performance of machine learning model predicting high and low frailty in an
398 HIV-positive population

399 **Table 4:** Gene ontology term enrichment analysis of the selected 119 CpG sites that
400 predict frailty in an HIV-positive population

401 **Table S1:** Differentially methylated regions between high and low frailty in training set

402

Table 2: A panel of 119 CpG sites in blood that predicts high and low frailty in a HIV-infected veteran population

CpG	Chr	Position	Gene	Gene group	Relation to CpG Island	Effect size	SE	P value	FDR*
cg01309328	6	32811253	<i>PSMB8</i>	Body	N_Shore	-5.30	0.888	4.44E-09	9.32E-04
cg08099136	6	32811251	<i>PSMB8</i>	Body	N_Shore	-4.77	0.800	4.63E-09	9.32E-04
cg26312951	21	42797847	<i>MX1</i>	TSS200;5UTR	N_Shore	-2.46	0.417	6.40E-09	9.32E-04
cg08818207	6	32820355	<i>TAP1</i>	Body	N_Shore	-3.40	0.627	9.17E-08	1.00E-02
cg08122652	3	122281939	<i>PARP9</i>	5UTR;TSS1500	N_Shore	-1.79	0.341	2.18E-07	1.45E-02
cg07352001	10	44144445	<i>ZNF32</i>	TSS200;TSS1500	Island	-	2.730	2.31E-07	1.45E-02
cg07325529	6	6613762	<i>LY86</i>	Body		14.30	3.980	2.77E-07	1.45E-02
cg16242615	19	4059988	<i>ZBTB7A</i>	5UTR	Island	-6.27	1.200	2.85E-07	1.45E-02
cg00533183	6	32810742	<i>PSMB8</i>	Body	N_Shore	-5.81	1.120	3.18E-07	1.45E-02
cg12649038	10	116282534	<i>ABLIM1</i>	5UTR		5.03	0.973	3.32E-07	1.45E-02
cg08926253	11	614761	<i>IRF7</i>	Body	Island	-2.91	0.571	4.94E-07	1.96E-02
cg03038262	11	315262	<i>IFITM1</i>	3UTR	N_Shore	-2.86	0.565	5.83E-07	2.12E-02
cg07839457	16	57023022	<i>NLRC5</i>	TSS1500	N_Shore	-1.99	0.396	6.55E-07	2.12E-02
cg22930808	3	122281881	<i>PARP9</i>	5UTR	N_Shore	-1.54	0.306	6.78E-07	2.12E-02
cg03607951	1	79085586	<i>IFI44L</i>	TSS1500		-2.00	0.398	7.41E-07	2.16E-02
cg08260450	6	34993987	<i>ANKS1A</i>	Body		6.35	1.270	8.43E-07	2.30E-02
cg03917473	17	38764244				9.10	1.840	1.04E-06	2.67E-02
cg18234224	1	172917851				6.59	1.350	1.34E-06	3.22E-02
cg06188083	10	91093005	<i>IFIT3</i>	Body		-2.02	0.414	1.40E-06	3.22E-02
cg23570810	11	315102	<i>IFITM1</i>	Body	N_Shore	-2.31	0.475	1.58E-06	3.41E-02
cg03816851	22	18324769	<i>MICAL3</i>	Body	Island	10.30	2.120	1.64E-06	3.41E-02
cg20676542	17	54991456	<i>TRIM25</i>	TSS200	Island	-3.02	0.629	2.10E-06	4.17E-02
cg05019807	9	139410393	<i>NOTCH1</i>	Body	N_Shore	5.97	1.260	2.70E-06	5.07E-02
cg03753191	13	43566902	<i>EPST11</i>	TSS1500	S_Shore	-5.39	1.140	2.88E-06	5.07E-02
cg01971407	11	313624	<i>IFITM1</i>	TSS1500	N_Shelf	-3.53	0.745	2.90E-06	5.07E-02
cg03359362	19	47289611	<i>SLC1A5</i>	TSS1500;Body;5UTR	N_Shore	-	2.670	3.13E-06	5.10E-02
cg22116398	5	196162			S_Shore	5.15	1.090	3.15E-06	5.10E-02
cg26922780	16	88769443	<i>RNF166</i>	Body	N_Shelf	4.27	0.907	3.28E-06	5.12E-02
cg15748006	2	9772375	<i>YWHAQ</i>	TSS1500	S_Shore	4.45	0.949	3.45E-06	5.20E-02

cg22862003	21	42797588	<i>MX1</i>	TSS1500;5UTR	N_Shore	-1.66	0.356	3.88E-06	5.65E-02
cg21366673	6	30459512	<i>HLA-E</i>	Body	S_Shore	-4.19	0.900	4.04E-06	5.69E-02
cg05588757	2	95825608	<i>ZNF514</i>	TSS1500	Island	-19.40	4.220	5.25E-06	6.64E-02
cg06872964	1	79085250	<i>IFI44L</i>	TSS1500		-1.92	0.416	5.26E-06	6.64E-02
cg06376949	10	91173811	<i>IFIT5</i>	TSS1500	N_Shore	-3.21	0.698	5.29E-06	6.64E-02
cg01154505	2	112940409	<i>FBLN7</i>	Body	S_Shore	6.16	1.340	5.32E-06	6.64E-02
cg01871595	13	39336304	<i>FREM2</i>	Body		6.90	1.500	5.59E-06	6.78E-02
cg05696877	1	79088769	<i>IFI44L</i>	5UTR		-1.21	0.265	6.01E-06	6.94E-02
cg24497541	1	55352819	<i>DHCR24</i>	1stExon;5UTR	Island	-9.81	2.150	6.12E-06	6.94E-02
cg11829870	22	50988451	<i>KLHDC7B</i>	3UTR;1stExon	S_Shore	-3.63	0.795	6.28E-06	6.94E-02
cg13720750	5	112309167			N_Shelf	7.28	1.600	6.37E-06	6.94E-02
cg18255813	6	7195966	<i>RREB1</i>	Body		9.38	2.060	6.76E-06	6.94E-02
cg21549285	21	42799141	<i>MX1</i>	5UTR	S_Shore	-1.01	0.224	7.08E-06	6.94E-02
cg14945867	14	54908007	<i>CNIH</i>	1stExon	Island	-29.70	6.550	7.16E-06	6.94E-02
cg09026253	11	313267	<i>IFITM1</i>	TSS1500	S_Shore	-3.48	0.768	7.25E-06	6.94E-02
cg14651616	11	64563992	<i>MAP4K2</i>	Body		-6.98	1.540	7.27E-06	6.94E-02
cg24082730	3	126076366	<i>KLF15</i>	TSS200	Island	-10.90	2.410	7.31E-06	6.94E-02
cg16302816	11	16834800	<i>PLEKH A7</i>	Body		6.50	1.440	7.49E-06	6.96E-02
cg01765174	9	100880960	<i>TRIM14</i>	Body	N_Shore	-4.47	0.993	8.35E-06	7.60E-02
cg10552523	11	313478	<i>IFITM1</i>	TSS1500	N_Shelf	-3.14	0.702	9.42E-06	8.40E-02
cg19025187	3	195808255	<i>TFRC</i>	5UTR	N_Shore	-6.73	1.510	1.07E-05	9.16E-02
cg21081878	21	38334730	<i>HLCS</i>	5UTR	N_Shelf	3.49	0.785	1.08E-05	9.16E-02
cg14299044	10	19972179				5.32	1.200	1.09E-05	9.16E-02
cg24871132	3	149688846	<i>PFN2</i>	TSS200	Island	-21.70	4.900	1.17E-05	9.55E-02
cg27294701	16	88107336	<i>BANP</i>	Body	Island	7.04	1.590	1.18E-05	9.55E-02
cg09296453	6	29692035	<i>HLA-F</i>	Body	Island	-3.24	0.735	1.31E-05	1.04E-01
cg16400434	11	73882363	<i>PPME1 ;C2CD3</i>	TSS200;TSS1500	Island	-25.60	5.850	1.46E-05	1.13E-01
cg16656286	17	4981603	<i>ZFP3</i>	TSS200	Island	22.40	5.130	1.52E-05	1.13E-01
cg26480543	19	55629279	<i>PPP1R12C</i>	TSS1500	S_Shore	-7.33	1.680	1.52E-05	1.13E-01
cg19371652	12	113415883	<i>OAS2</i>	TSS1500		-3.59	0.822	1.52E-05	1.13E-01
cg09597638	17	3907349			N_Shore	-5.48	1.260	1.56E-05	1.14E-01
cg00994629	14	22694547				-7.44	1.710	1.62E-05	1.16E-01

cg02314339	10	91020653				-4.35	1.010	1.86E-05	1.29E-01
cg21468416	9	12701991 4	<i>NEK6</i>	1stExon;TSS1500;5UTR	N_Shore	6.61	1.530	1.89E-05	1.29E-01
cg12149905	3	61547208	<i>PTPRG</i>	TSS200	Island	-	2.820	1.89E-05	1.29E-01
cg20927242	6	29692011	<i>HLA-F</i>	Body	Island	-7.14	1.660	1.97E-05	1.32E-01
cg03332034	16	1823786	<i>EME2;</i> <i>MRPS3</i> 4	TSS1500	Island	-	2.870	2.12E-05	1.36E-01
cg05489271	12	12387169 5	<i>SETD8</i>	Body	N_Shelf	7.83	1.820	2.13E-05	1.36E-01
cg25843003	6	31431312	<i>HCP5</i>	3UTR		-3.29	0.768	2.19E-05	1.36E-01
cg14090510	7	64839023	<i>ZNF92</i>	5UTR;Body	Island	-	3.440	2.21E-05	1.36E-01
cg25467833	1	21233506 1				4.86	1.140	2.28E-05	1.36E-01
cg03725115	6	30458102	<i>HLA-E</i>	Body	Island	-5.09	1.190	2.29E-05	1.36E-01
cg15620384	6	34164405			Island	-6.69	1.570	2.34E-05	1.36E-01
cg08450404	19	58326441	<i>ZNF55</i> 2	TSS200	S_Shore	-8.37	1.960	2.35E-05	1.36E-01
cg06927297	12	11717588 9	<i>RNFT2</i>	TSS200;TSS1500	Island	-	6.460	2.36E-05	1.36E-01
cg21684411	6	31431573	<i>HCP5</i>	3UTR		-6.71	1.570	2.40E-05	1.36E-01
cg20998539	17	62208374	<i>ERN1</i>	TSS1500	S_Shore	-7.32	1.720	2.42E-05	1.36E-01
cg25138854	1	9555557			Island	-7.03	1.650	2.45E-05	1.36E-01
cg13250752	4	13828199 2				-6.09	1.430	2.46E-05	1.36E-01
cg26562691	16	23850404	<i>PRKCB</i>	Body	S_Shelf	6.87	1.620	2.47E-05	1.36E-01
cg23677352	4	14971567 5				-7.92	1.860	2.49E-05	1.36E-01
cg22107533	15	45028083	<i>TRIM6</i> 9	TSS1500		-3.14	0.739	2.53E-05	1.36E-01
cg06412917	12	12485884 4	<i>NCOR2</i>	Body	S_Shelf	6.47	1.520	2.57E-05	1.37E-01
cg01190666	20	62204908	<i>PRIC28</i> 5	5UTR	N_Shore	-3.47	0.818	2.67E-05	1.40E-01
cg11977562	13	11484529 7	<i>RASA3</i>	Body	N_Shelf	7.05	1.660	2.70E-05	1.40E-01
cg13304609	1	79085162	<i>IFI44L</i>	TSS1500		-1.54	0.364	2.72E-05	1.40E-01
cg01537765	19	42914828	<i>LIPE</i>	Body	Island	22.00	5.210	2.75E-05	1.40E-01
cg08275025	11	314493	<i>IFITM1</i>	Body	N_Shore	-4.98	1.180	2.78E-05	1.40E-01
cg05362517	13	37393368	<i>RFXAP</i>	5UTR;1stExon	Island	-	3.260	2.92E-05	1.42E-01
cg21145248	5	17681667 9	<i>SLC34</i> <i>A1</i>	Body		6.04	1.430	2.94E-05	1.42E-01
cg07537655	6	16141561 2	<i>MAP3</i> <i>K4</i>	Body	S_Shelf	1.77	0.419	2.95E-05	1.42E-01
cg26912671	1	66458803	<i>PDE4B</i>	1stExon		10.70	2.540	3.01E-05	1.42E-01

cg11905821	6	31770932	<i>LSM2</i>	Body	N_Shelf	11.80	2.800	3.05E-05	1.42E-01
cg06202053	11	64109021	<i>CCDC8 8B</i>	Body	N_Shore	-5.84	1.390	3.09E-05	1.42E-01
cg13794687	7	12330239 9	<i>LMOD 2</i>	Body		6.73	1.600	3.11E-05	1.42E-01
cg23892836	6	29692085	<i>HLA-F</i>	Body	Island	-3.19	0.759	3.12E-05	1.42E-01
cg15331332	6	29692111	<i>HLA-F</i>	Body	S_Shore	-3.53	0.840	3.12E-05	1.42E-01
cg15804432	15	64454965	<i>PPIB</i>	Body	Island	-4.27	1.020	3.15E-05	1.42E-01
cg09251764	17	6659070	<i>XAF1</i>	TSS200		-5.73	1.370	3.23E-05	1.44E-01
cg20893717	7	10031819 0	<i>EPO</i>	TSS1500	Island	13.80	3.290	3.26E-05	1.44E-01
cg08159663	16	57022486	<i>NLRC5</i>	TSS1500	N_Shore	-3.44	0.821	3.29E-05	1.44E-01
cg10019429	19	32836659	<i>ZNF50 7</i>	1stExon;5UTR	Island	- 17.20	4.120	3.55E-05	1.54E-01
cg11791770	11	611791	<i>PHRF1</i>	3UTR	Island	-6.52	1.560	3.64E-05	1.56E-01
cg11653134	2	66805547			S_Shore	-4.24	1.020	3.81E-05	1.61E-01
cg06981309	3	14626095 4	<i>PLSCR1</i>	5UTR	N_Shore	-1.96	0.471	3.84E-05	1.61E-01
cg07168939	8	14376341 2	<i>PSCA</i>	Body		-5.88	1.420	3.88E-05	1.61E-01
cg04611649	2	15268124 0	<i>ARL5A</i>	5UTR	N_Shelf	7.42	1.790	4.04E-05	1.66E-01
cg16297569	1	92952517	<i>GFI1</i>	TSS1500;TSS200	Island	-8.34	2.020	4.16E-05	1.70E-01
cg01518846	6	26246970	<i>HIST1H 4G</i>	1stExon	Island	- 17.50	4.240	4.41E-05	1.74E-01
cg26852894	4	11157341				6.31	1.530	4.42E-05	1.74E-01
cg06811183	3	48510438	<i>SHISA5</i>	3UTR		-8.23	2.000	4.42E-05	1.74E-01
ch.10.130520 3R	10	63407435				-8.63	2.090	4.43E-05	1.74E-01
cg12461141	11	5710654	<i>TRIM2 2</i>	TSS1500		-3.40	0.829	4.79E-05	1.87E-01
cg13149600	21	43374220	<i>C2CD2</i>	TSS1500	S_Shore	-9.35	2.280	4.83E-05	1.87E-01
cg12660813	1	3192343	<i>PRDM 16</i>	Body	N_Shelf	7.28	1.780	4.96E-05	1.89E-01
cg24447788	19	795310			N_Shore	11.30	2.760	5.02E-05	1.89E-01
cg13861758	9	13814809 9			N_Shelf	5.08	1.240	5.02E-05	1.89E-01
cg09950208	10	13084123 5				-4.99	1.220	5.24E-05	1.95E-01
cg01482620	19	48835971	<i>TMEM 143</i>	3UTR	N_Shore	18.40	4.500	5.28E-05	1.95E-01
cg03634735	7	1992524	<i>MAD1 L1</i>	Body	S_Shore	6.19	1.520	5.37E-05	1.97E-01

*FDR: False discovery rate

403

404

Table S1: Differentially methylated regions between high and low fraity in training set

chr	start	end	value	area	Number of CpG	p.value	fwer	Gene.Symbol	FDR	CpG
3	122281881	122281975	-0.068	0.203	3	2.30E-07	0.001	<i>DTX3L,PARP9</i>	1.03E-04	cg00959259; cg08122652; cg22930808
12	11700321	11700489	-0.054	0.162	3	1.61E-06	0.006	<i>LINC01252</i>	3.60E-04	cg06202470; cg18232235; cg19651115 cg00458211;
1	79088769	79118191	-0.039	0.154	4	1.17E-05	0.049	<i>IFI44,IFI44L</i>	1.75E-03	cg01079652; cg05696877; cg07107453
21	42797588	42797847	-0.046	0.092	2	2.41E-05	0.081	<i>MX1</i>	2.68E-03	cg22862003; cg26312951 cg00855901;
1	79085162	79085765	-0.034	0.168	5	2.99E-05	0.122	<i>IFI44L</i>	2.68E-03	cg03607951; cg06872964; cg13304609; cg17980508
11	319555	319718	-0.036	0.109	3	4.73E-05	0.166	<i>IFITM3</i>	2.80E-03	cg20045320; cg09122035; cg17990365 cg01886988;
11	312518	313624	-0.023	0.230	10	4.89E-05	0.194	<i>IFITM1</i>	2.80E-03	cg01971407; cg04582010; cg05432003; cg09026253; cg10552523; cg11694510; cg20566897; cg22963452; cg27032101 cg12047941; cg15013527; cg25050332; cg27331665; cg03038262;
11	314493	317767	-0.021	0.228	11	5.01E-05	0.195	<i>IFITM1</i>	2.80E-03	cg08000731; cg08275025; cg16379091; cg18434560; cg21686213; cg23570810 cg08924203; cg21549285
21	42798747	42799141	-0.041	0.083	2	5.95E-05	0.197	<i>MX1</i>	2.96E-03	

406 **List of abbreviations**

407 AUC: Area Under Curve

408 CI: Confidence interval

409 DMR: differentially methylated region

410 DNA: Deoxyribonucleic acid

411 DNAm: DNA methylation

412 DAVID: Database for Annotation, Visualization and Integrated Discovery

413 EWAS: epigenome-wide association study

414 FDR: False discovery rate

415 FWER: Family-wise error rate

416 GLMNET: Lasso and Elastic-Net Regularized Generalized Linear Models

417 GO: Gene ontology

418 HIV: Human immunodeficiency virus

419 HM450K: Human Methylation 450K BeadChip

420 NK: Natural killer

421 PC: Principal component

422 QC: Quality control

423 SVM: Support Vector Machines

424 VACS: Veterans Aging Cohort Study

425 VACS index: Veterans Aging Cohort Study index

426 XGBoost: Extreme Gradient Boosting Tree

427

428

429 **Declarations**

430 ***Ethics approval and consent to participate***

431 The study was approved by the committee of the Human Research Subject Protection
432 at Yale University and the Institutional Research Board Committee of the Connecticut
433 Veteran Healthcare System. All subjects provided written consents.

434 ***Availability of data and materials***

435 Demographic and clinical variables and DNAm data for the VACS samples were
436 submitted to GEO dataset (GSE117861) and are available to the public. All codes for
437 analysis are also available upon a request to the corresponding author.

438 ***Competing interests***

439 The authors declare that they have no competing interests.

440 ***Funding***

441 The project was supported by the National Institute on Drug Abuse (R03DA039745,
442 R01DA038632, R01DA047063, R01DA047820) and the National Center for Post-
443 Traumatic Stress Disorder, USA.

444 ***Authors' contributions***

445 CS was responsible for data analysis and manuscript preparation. ACJ provided DNA
446 samples, clinical data, and contributed to manuscript preparation. XZ was responsible
447 for the bioinformatics data processing. VM involved clinical data collection and
448 manuscript preparation. DH and EJ contributed to analytical approach and the
449 manuscript preparation. KX was responsible for the study design, study protocol,
450 sample preparation, data analysis, interpretation of findings, and manuscript
451 preparation.

452 ***Acknowledgements***

453 The authors appreciate the support of the Veteran Aging Study Cohort Biomarker Core
454 and Yale Center of Genomic Analysis.

455

456

457

458 References

- 459 1. Kooij KW, Wit FW, Schouten J, van der Valk M, Godfried MH, Stolte IG, et al. HIV
460 infection is independently associated with frailty in middle-aged HIV type 1-infected individuals
461 compared with similar but uninfected controls. *AIDS*. 2016;30(2):241-50.
- 462 2. Dent E, Kowal P, Hoogendijk EO. Frailty measurement in research and clinical practice: A
463 review. *Eur J Intern Med*. 2016;31:3-10.
- 464 3. Lam K, Pan K, Linnekamp JF, Medema JP, Kandimalla R. DNA methylation based
465 biomarkers in colorectal cancer: a systematic review. *Biochimica et Biophysica Acta (BBA)-*
466 *Reviews on Cancer*. 2016;1866(1):106-20.
- 467 4. Teroganova N, Girshkin L, Suter CM, Green MJ. DNA methylation in peripheral tissue of
468 schizophrenia and bipolar disorder: a systematic review. *BMC genetics*. 2016;17(1):27.
- 469 5. Bakusic J, Schaufeli W, Claes S, Godderis L. Stress, burnout and depression: A systematic
470 review on DNA methylation mechanisms. *Journal of Psychosomatic Research*. 2017;92:34-44.
- 471 6. Li M, D'Arcy C, Li X, Zhang T, Joober R, Meng X. What do DNA methylation studies tell us
472 about depression? A systematic review. *Translational psychiatry*. 2019;9(1):68.
- 473 7. Horvath S. DNA methylation age of human tissues and cell types. *Genome biology*.
474 2013;14(10):3156.
- 475 8. Horvath S, Raj K. DNA methylation-based biomarkers and the epigenetic clock theory of
476 ageing. *Nature Reviews Genetics*. 2018:1.
- 477 9. Hannum G, Guinney J, Zhao L, Zhang L, Hughes G, Sada S, et al. Genome-wide
478 methylation profiles reveal quantitative views of human aging rates. *Molecular cell*.
479 2013;49(2):359-67.
- 480 10. Liu C, Marioni RE, Hedman ÅK, Pfeiffer L, Tsai P-C, Reynolds LM, et al. A DNA
481 methylation biomarker of alcohol consumption. *Molecular psychiatry*. 2018;23(2):422.
- 482 11. Breitling LP, Yang R, Korn B, Burwinkel B, Brenner H. Tobacco-smoking-related
483 differential DNA methylation: 27K discovery and replication. *The American Journal of Human*
484 *Genetics*. 2011;88(4):450-7.
- 485 12. Gao X, Zhang Y, Breitling LP, Brenner H. Relationship of tobacco smoking and smoking-
486 related DNA methylation with epigenetic age acceleration. *Oncotarget*. 2016;7(30):46878-89.
- 487 13. Joubert BR, Håberg SE, Nilsen RM, Wang X, Vollset SE, Murphy SK, et al. 450K
488 epigenome-wide scan identifies differential DNA methylation in newborns related to maternal
489 smoking during pregnancy. *Environmental health perspectives*. 2012;120(10):1425-31.
- 490 14. Lee KW, Pausova Z. Cigarette smoking and DNA methylation. *Frontiers in genetics*.
491 2013;4:132.
- 492 15. Zhang X, Hu Y, Aouizerat BE, Peng G, Marconi VC, Corley MJ, et al. Machine learning
493 selected smoking-associated DNA methylation signatures that predict HIV prognosis and
494 mortality. 2018;10(1):155.
- 495 16. Zhang R, Miao Q, Wang C, Zhao R, Li W, Haile CN, et al. Genome-wide DNA methylation
496 analysis in alcohol dependence. *Addiction biology*. 2013;18(2):392-403.
- 497 17. Kraiczy J, Nayak KM, Howell KJ, Ross A, Forbester J, Salvestrini C, et al. DNA methylation
498 defines regional identity of human intestinal epithelial organoids and undergoes dynamic
499 changes during development. *Gut*. 2019;68(1):49-61.

- 500 18. Nano J, Ghanbari M, Wang W, de Vries PS, Dhana K, Muka T, et al. Epigenome-Wide
501 Association Study Identifies Methylation Sites Associated With Liver Enzymes and Hepatic
502 Steatosis. *Gastroenterology*. 2017;153(4):1096-106.e2.
- 503 19. Delpu Y, Cordelier P, Cho W, Torrisani J. DNA methylation and cancer diagnosis.
504 *International journal of molecular sciences*. 2013;14(7):15029-58.
- 505 20. Figueroa ME, Lugthart S, Li Y, Erpelinck-Verschueren C, Deng X, Christos PJ, et al. DNA
506 methylation signatures identify biologically distinct subtypes in acute myeloid leukemia. *Cancer*
507 *cell*. 2010;17(1):13-27.
- 508 21. Holm K, Hegardt C, Staaf J, Vallon-Christersson J, Jönsson G, Olsson H, et al. Molecular
509 subtypes of breast cancer are associated with characteristic DNA methylation patterns. *Breast*
510 *cancer research*. 2010;12(3):R36.
- 511 22. Berdasco M, Esteller M. Clinical epigenetics: seizing opportunities for translation. *Nature*
512 *Reviews Genetics*. 2019;20(2):109-27.
- 513 23. Mohammad HP, Barbash O, Creasy CL. Targeting epigenetic modifications in cancer
514 therapy: erasing the roadmap to cancer. *Nat Med*. 2019;25(3):403-18.
- 515 24. Erlandson KM, Ng DK, Jacobson LP, Margolick JB, Dobs AS, Palella Jr FJ, et al.
516 Inflammation, immune activation, immunosenescence, and hormonal biomarkers in the frailty-
517 related phenotype of men with or at risk for HIV infection. *The Journal of infectious diseases*.
518 2016;215(2):228-37.
- 519 25. Zhang X, Justice AC, Hu Y, Wang Z, Zhao H, Wang G, et al. Epigenome-wide differential
520 DNA methylation between HIV-infected and uninfected individuals. *Epigenetics*.
521 2016;11(10):750-60.
- 522 26. Chen J, Huang Y, Hui Q, Mathur R, Gwinn M, So-Armah K, et al. Epigenetic Associations
523 with Estimated Glomerular Filtration Rate (eGFR) among Men with HIV Infection. *Clin Infect Dis*.
524 2019.
- 525 27. Mathur R, Hui Q, Huang Y, Gwinn M, So-Armah K, Freiberg MS, et al. DNA Methylation
526 Markers of Type 2 Diabetes Mellitus Among Male Veterans With or Without Human
527 Immunodeficiency Virus Infection. *J Infect Dis*. 2019;219(12):1959-62.
- 528 28. Breitling LP, Saum K-U, Perna L, Schöttker B, Holleczer B, Brenner H. Frailty is associated
529 with the epigenetic clock but not with telomere length in a German cohort. *Clinical Epigenetics*.
530 2016;8(1):21.
- 531 29. Sanchez-Conde M, Rodriguez-Centeno J, Dronda F, Lopez JC, Jimenez Z, Berenguer J, et
532 al. Frailty phenotype: a clinical marker of age acceleration in the older HIV-infected population.
533 *Epigenomics*. 2019;11(5):501-9.
- 534 30. Horvath S, Levine AJ. HIV-1 infection accelerates age according to the epigenetic clock.
535 *The Journal of infectious diseases*. 2015;212(10):1563-73.
- 536 31. Rickabaugh TM, Baxter RM, Sehl M, Sinsheimer JS, Hultin PM, Hultin LE, et al.
537 Acceleration of age-associated methylation patterns in HIV-1-infected adults. *PloS one*.
538 2015;10(3):e0119201.
- 539 32. Gross AM, Jaeger PA, Kreisberg JF, Licon K, Jepsen KL, Khosroheidari M, et al.
540 Methylome-wide analysis of chronic HIV infection reveals five-year increase in biological age
541 and epigenetic targeting of HLA. *Molecular cell*. 2016;62(2):157-68.

- 542 33. Nelson KN, Hui Q, Rimland D, Xu K, Freiberg MS, Justice AC, et al. Identification of HIV
543 infection-related DNA methylation sites and advanced epigenetic aging in HIV-positive,
544 treatment-naive U.S. veterans. *AIDS*. 2017;31(4):571-5.
- 545 34. Marioni RE, Shah S, McRae AF, Chen BH, Colicino E, Harris SE, et al. DNA methylation age
546 of blood predicts all-cause mortality in later life. *Genome biology*. 2015;16(1):25.
- 547 35. Marioni RE, Harris SE, Shah S, McRae AF, von Zglinicki T, Martin-Ruiz C, et al. The
548 epigenetic clock and telomere length are independently associated with chronological age and
549 mortality. *International journal of epidemiology*. 2016;45(2):424-32.
- 550 36. Perna L, Zhang Y, Mons U, Holleczeck B, Saum K-U, Brenner H. Epigenetic age
551 acceleration predicts cancer, cardiovascular, and all-cause mortality in a German case cohort.
552 *Clinical epigenetics*. 2016;8(1):64.
- 553 37. Christiansen L, Lenart A, Tan Q, Vaupel JW, Aviv A, McGue M, et al. DNA methylation
554 age is associated with mortality in a longitudinal Danish twin study. *Aging Cell*. 2016;15(1):149-
555 54.
- 556 38. Fransquet PD, Wrigglesworth J, Woods RL, Ernst ME, Ryan J. The epigenetic clock as a
557 predictor of disease and mortality risk: a systematic review and meta-analysis. *Clin Epigenetics*.
558 2019;11(1):62.
- 559 39. Zhang Y, Wilson R, Heiss J, Breitling LP, Saum K-U, Schöttker B, et al. DNA methylation
560 signatures in peripheral blood strongly predict all-cause mortality. *Nature communications*.
561 2017;8:14617.
- 562 40. Tate JP, Justice AC, Hughes MD, Bonnet F, Reiss P, Mocroft A, et al. An internationally
563 generalizable risk index for mortality after one year of antiretroviral therapy. *Aids*.
564 2013;27(4):563-72.
- 565 41. Justice AC, Dombrowski E, Conigliaro J, Fultz SL, Gibson D, Madenwald T, et al. Veterans
566 aging cohort study (VACS): overview and description. 2006;44(8 Suppl 2):S13.
- 567 42. Lehne B, Drong AW, Loh M, Zhang W, Scott WR, Tan S-T, et al. A coherent approach for
568 analysis of the Illumina HumanMethylation450 BeadChip improves data quality and
569 performance in epigenome-wide association studies. 2015;16(1):37.
- 570 43. Aryee MJ, Jaffe AE, Corrada-Bravo H, Ladd-Acosta C, Feinberg AP, Hansen KD, et al.
571 Minfi: a flexible and comprehensive Bioconductor package for the analysis of Infinium DNA
572 methylation microarrays. 2014;30(10):1363-9.
- 573 44. Jaffe AE, Murakami P, Lee H, Leek JT, Fallin MD, Feinberg AP, et al. Bump hunting to
574 identify differentially methylated regions in epigenetic epidemiology studies. *Int J Epidemiol*.
575 2012;41(1):200-9.
- 576 45. Bebu I, Tate J, Rimland D, Mesner O, Macalino GE, Ganesan A, et al. The VACS index
577 predicts mortality in a young, healthy HIV population starting highly active antiretroviral
578 therapy. *Journal of acquired immune deficiency syndromes (1999)*. 2014;65(2):226.
- 579 46. Kuhn MJ. *Joss. Building predictive models in R using the caret package*. 2008;28(5):1-26.
- 580 47. Hastie T, Tibshirani R, Friedman J, Franklin J. The elements of statistical learning: data
581 mining, inference and prediction. *The Mathematical Intelligencer*. 2005;27(2):83-5.
- 582 48. Kotsiantis SB, Zaharakis I, Pintelas P. Supervised machine learning: A review of
583 classification techniques. *Emerging artificial intelligence applications in computer engineering*.
584 2007;160:3-24.

- 585 49. Chen T, Guestrin C, editors. Xgboost: A scalable tree boosting system. Proceedings of the
586 22nd acm sigkdd international conference on knowledge discovery and data mining; 2016:
587 ACM.
- 588 50. Ogutu JO, Piepho H-P, Schulz-Streeck T. A comparison of random forests, boosting and
589 support vector machines for genomic selection. BMC Proc. 2011;5 Suppl 3(Suppl 3):S11-S.
- 590 51. Brodersen KH, Ong CS, Stephan KE, Buhmann JM, editors. The balanced accuracy and its
591 posterior distribution. 2010 20th International Conference on Pattern Recognition; 2010: IEEE.
- 592 52. Fresno C, Fernández EAJB. RDAVIDWebService: a versatile R interface to DAVID.
593 2013;29(21):2810-1.
- 594 53. Justice AC, Modur SP, Tate JP, Althoff KN, Jacobson LP, Gebo KA, et al. Predictive
595 accuracy of the Veterans Aging Cohort Study index for mortality with HIV infection: a North
596 American cross cohort analysis. Journal of acquired immune deficiency syndromes (1999).
597 2013;62(2):149-63.
- 598 54. O'Callaghan CA, Bell JIJr. Structure and function of the human MHC class Ib molecules
599 HLA-E, HLA-F and HLA-G. 1998;163(1):129-38.

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

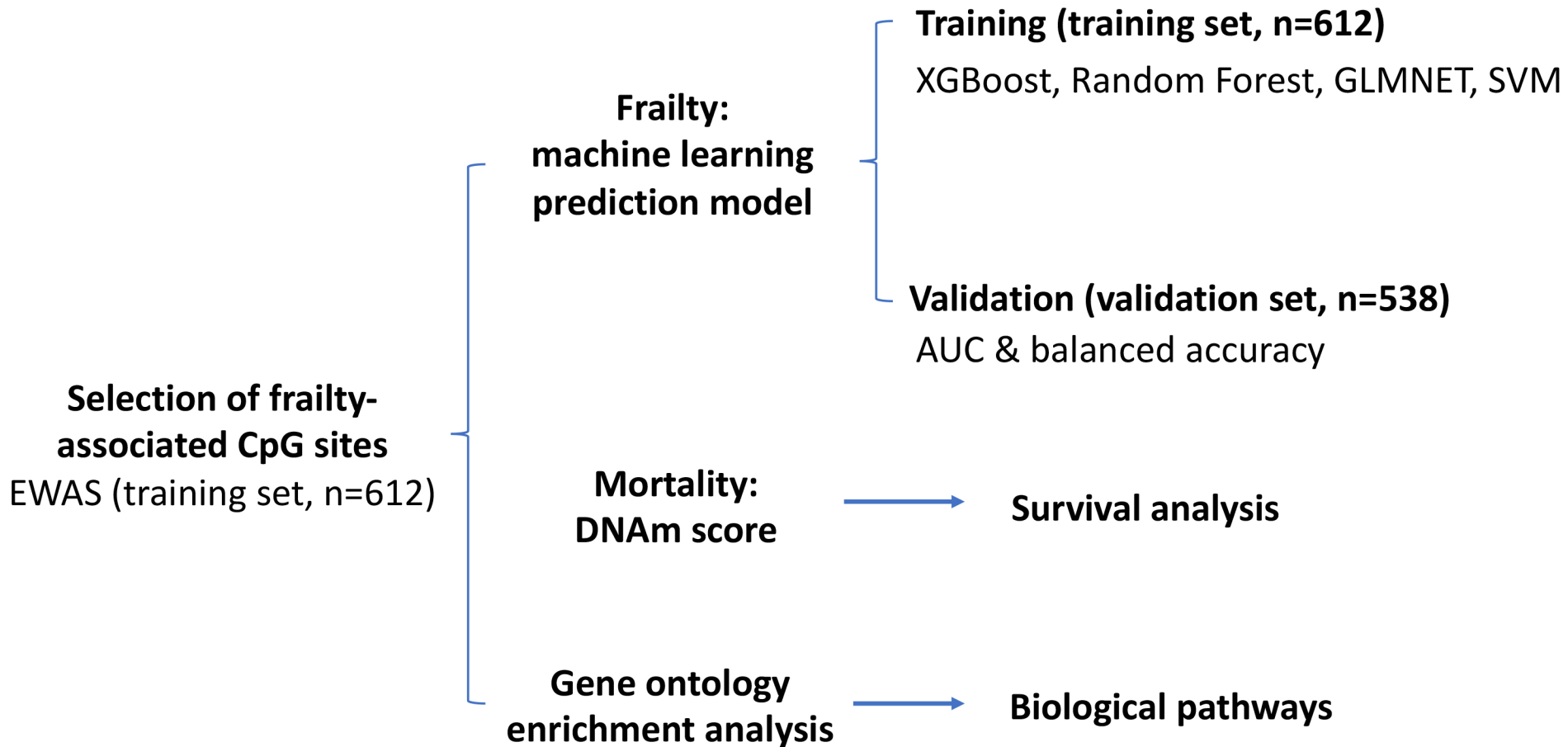


Figure 1

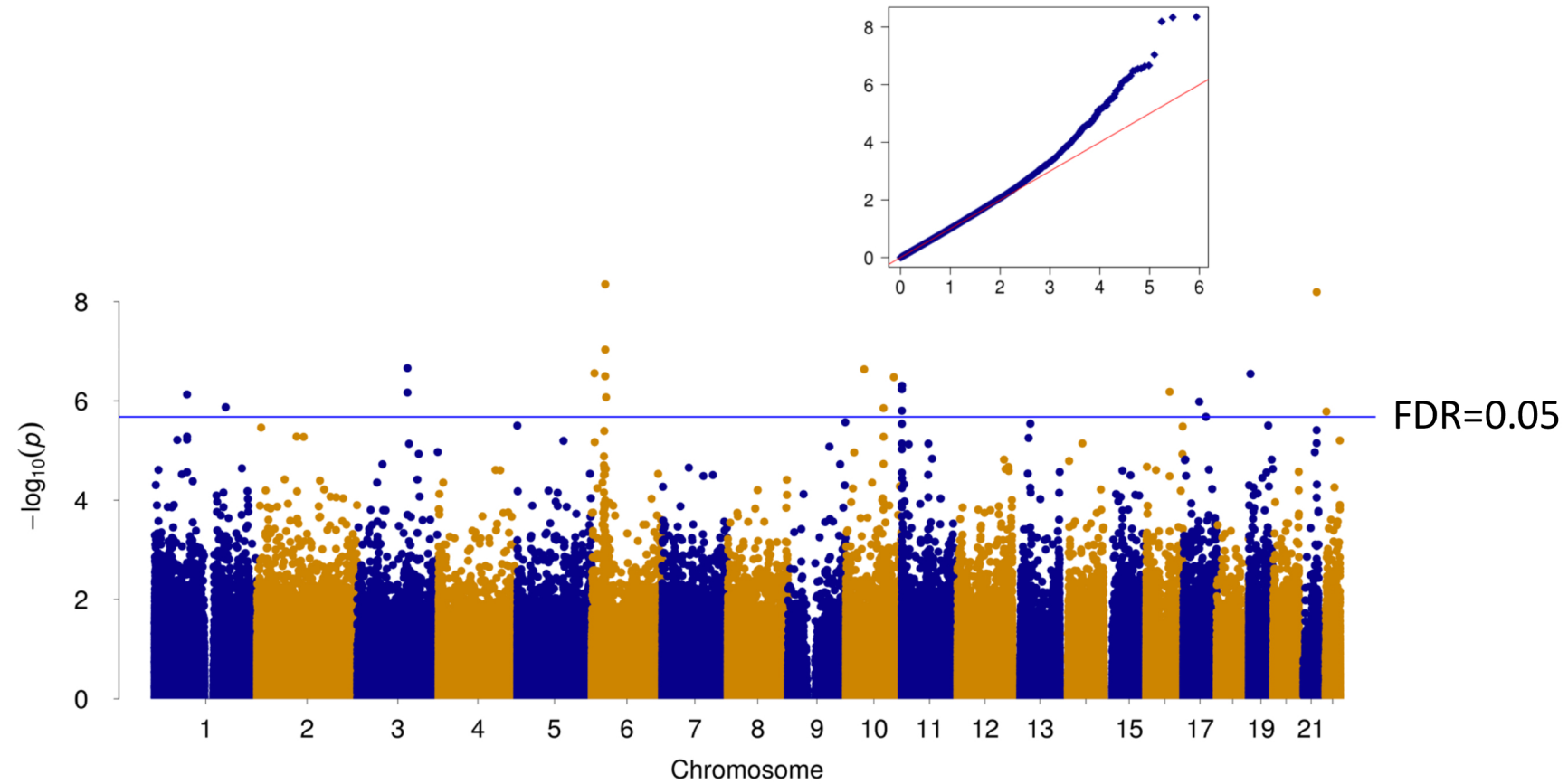


Figure 2

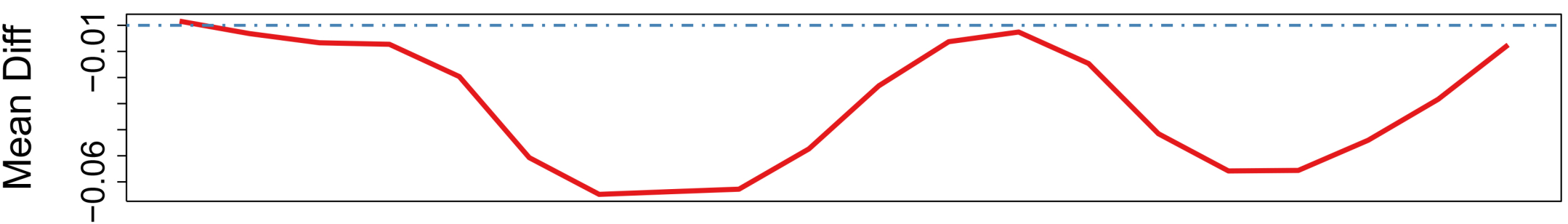
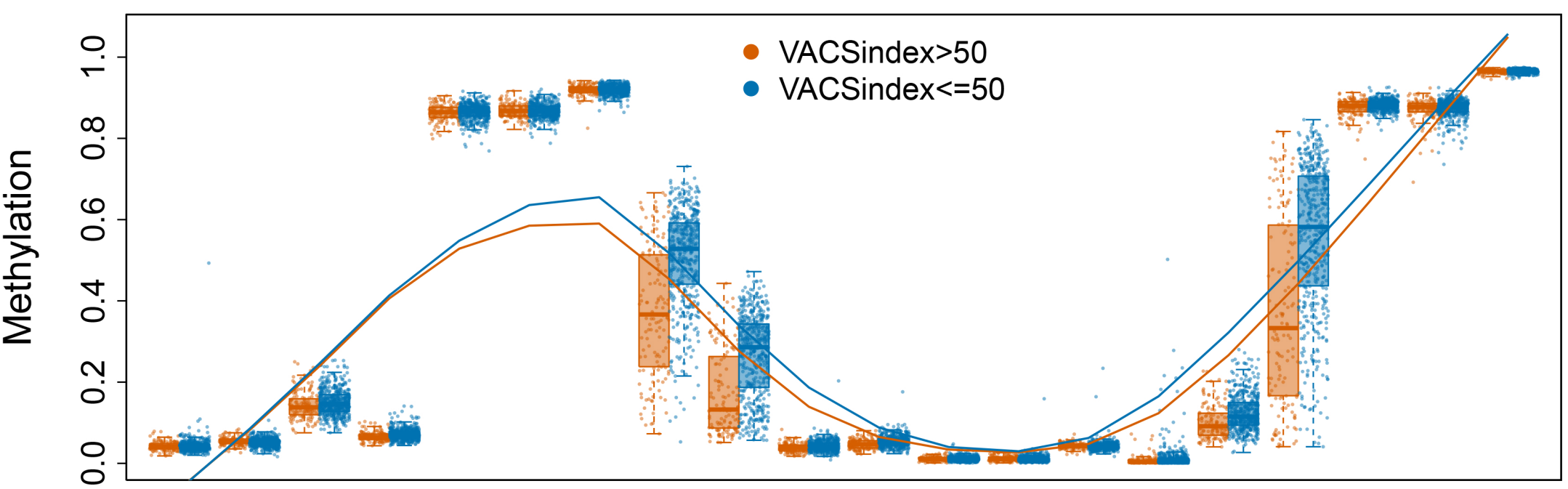
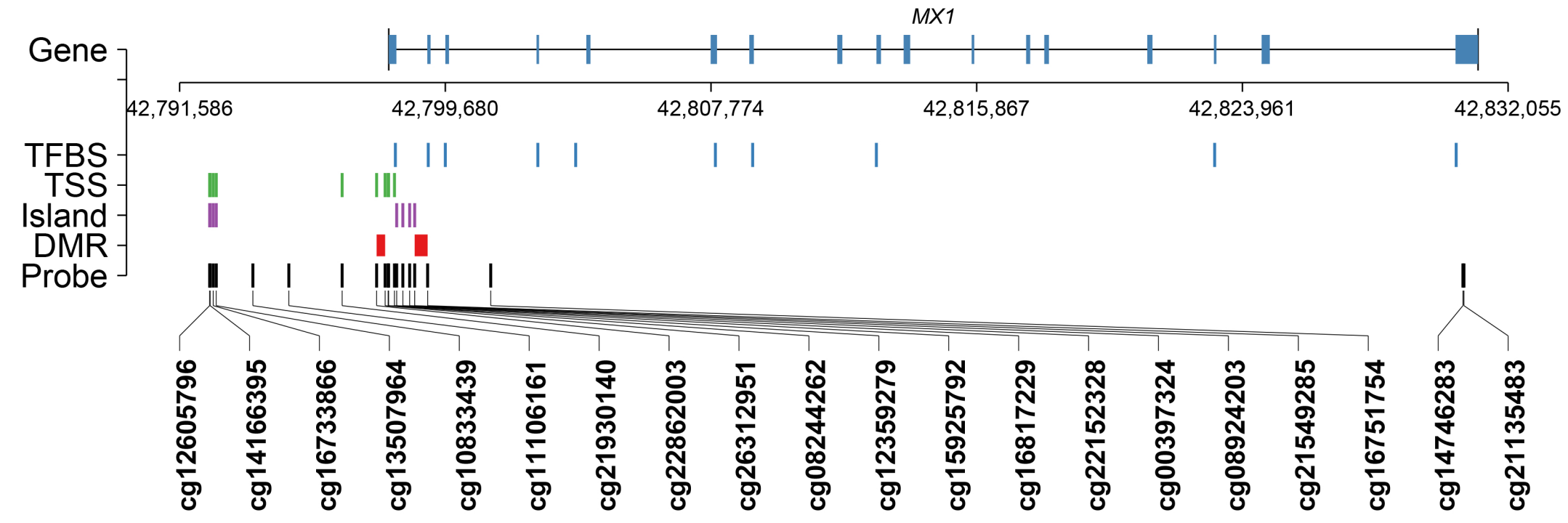


Figure 3a

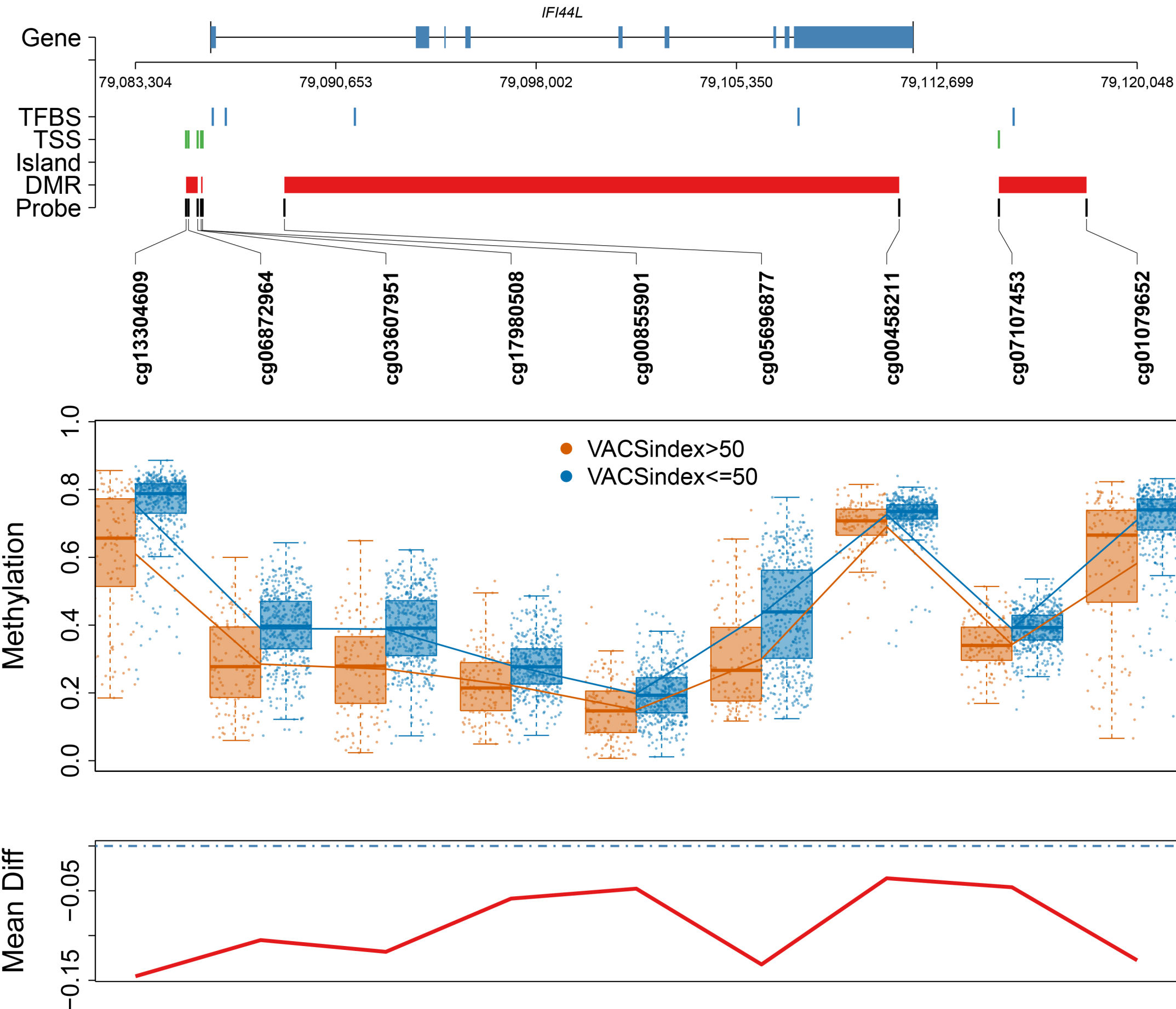


Figure 3b

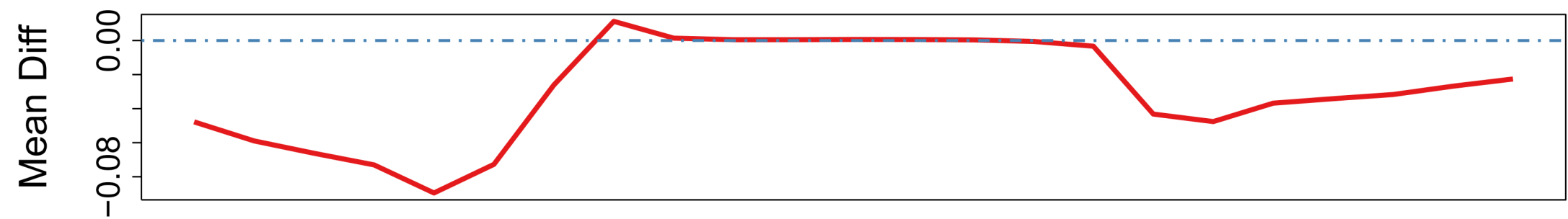
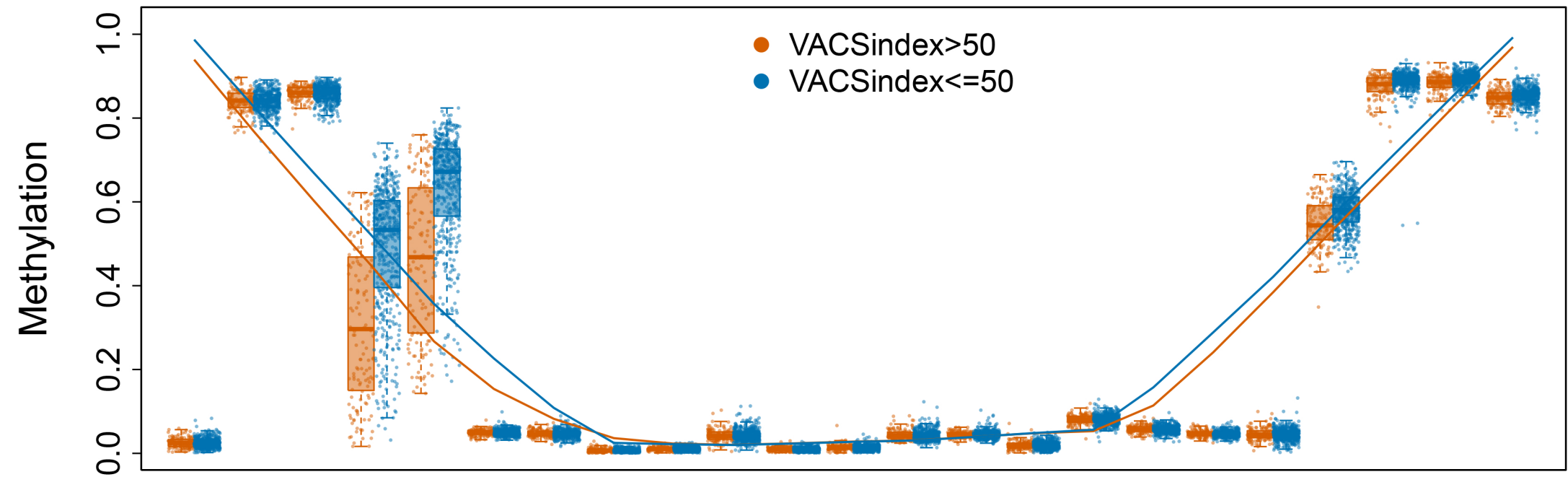
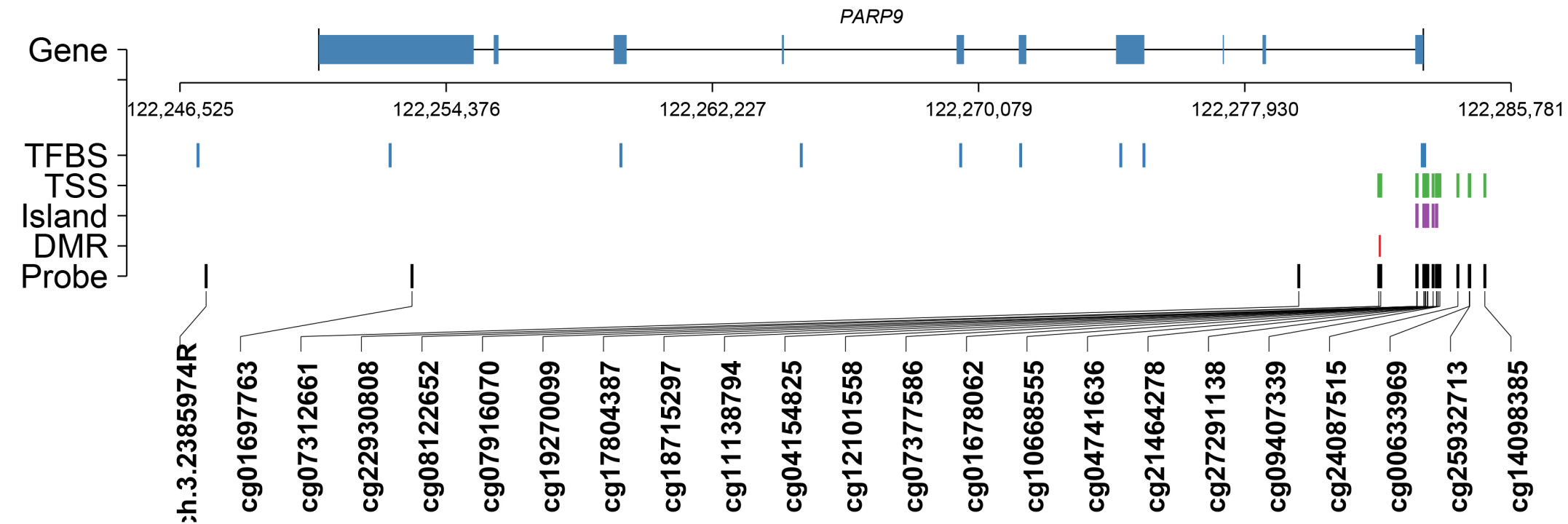


Figure 3c

High frailty (VACS index > 50)

Low frailty (VACS index ≤ 16)

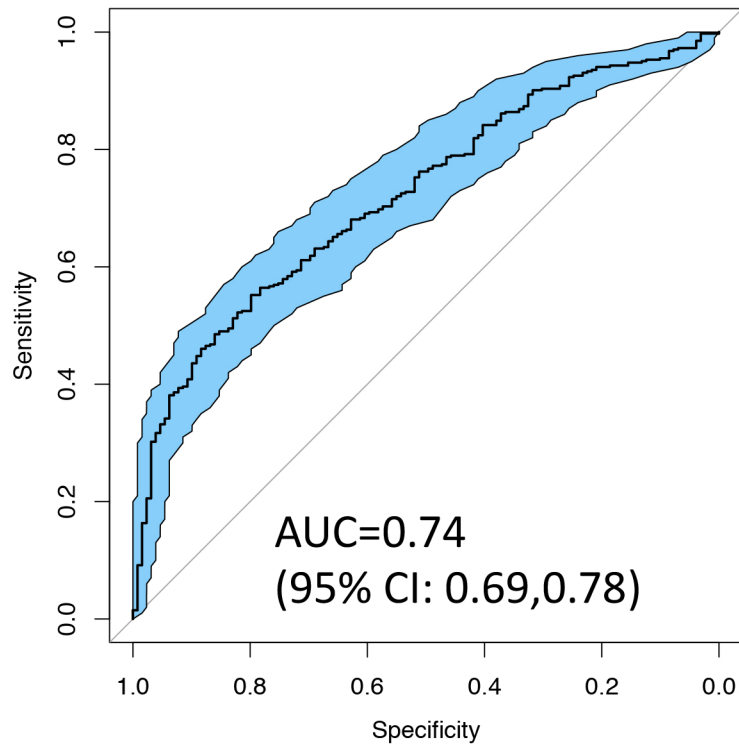
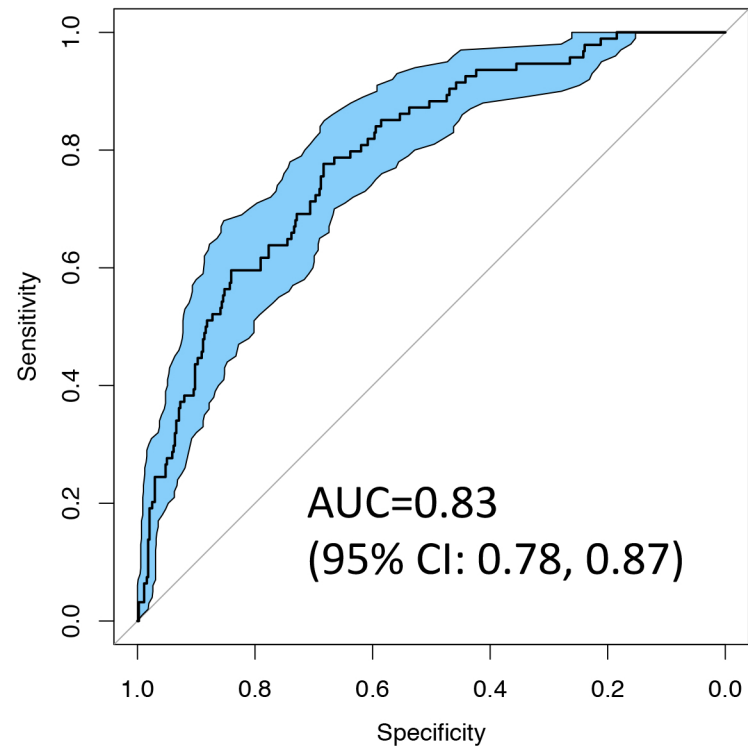


Figure 4

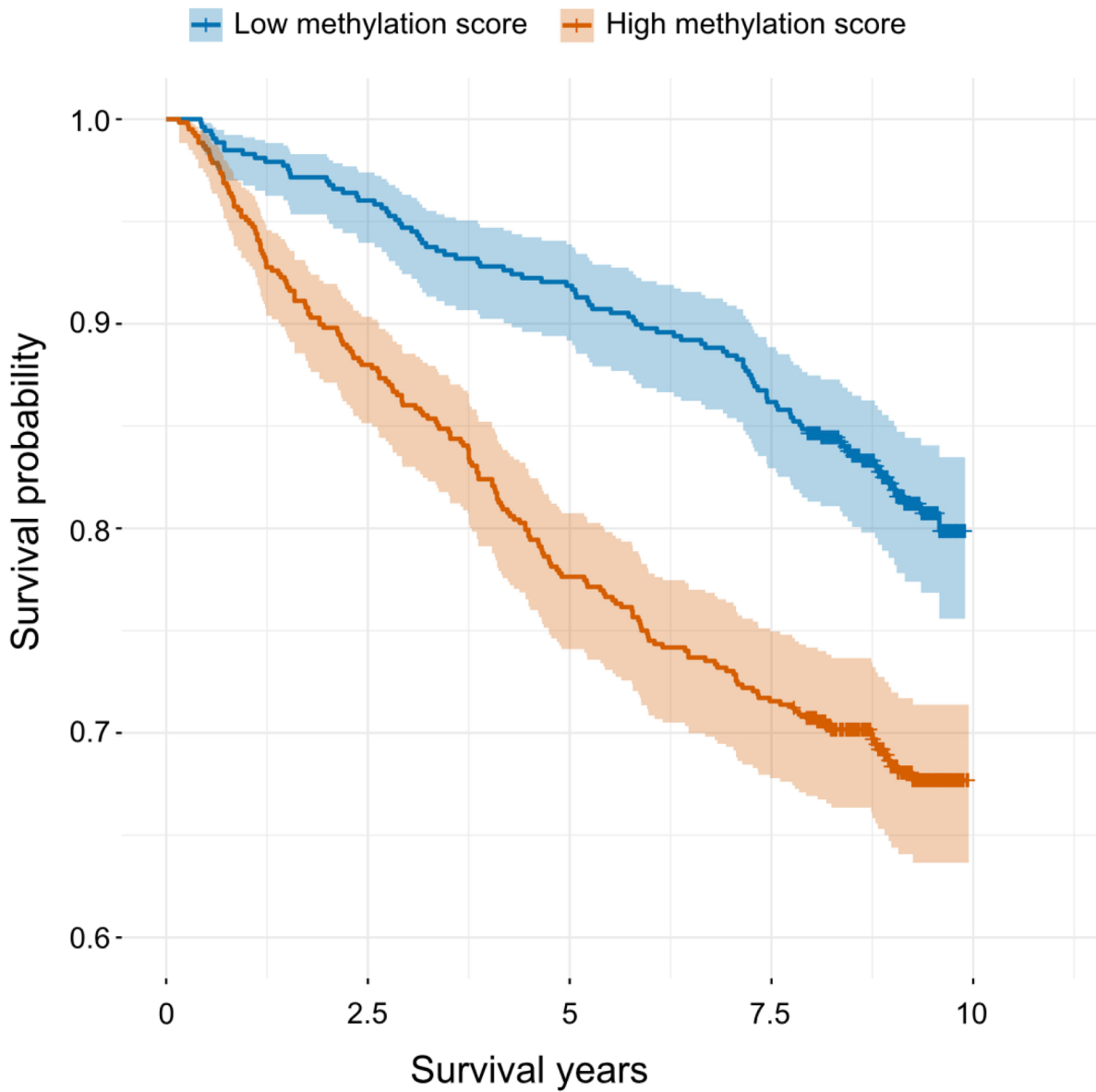


Figure 5

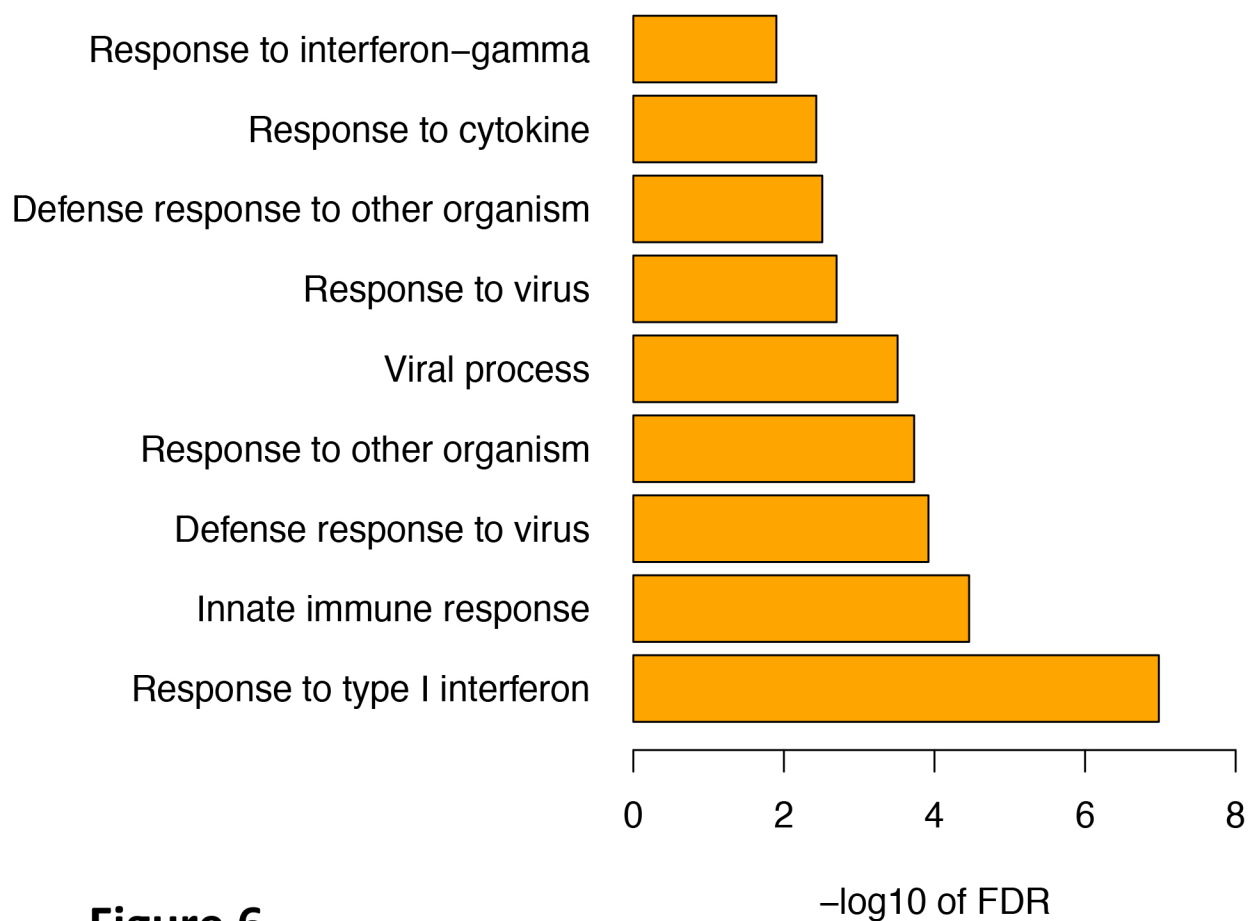


Figure 6