

## Histone mark age of human tissues and cells

Lucas Paulo de Lima Camillo<sup>1,2,\*</sup>, Muhammad Haider Asif<sup>3</sup>, Steve Horvath<sup>4</sup>, Erica Larschan<sup>5,6</sup>, Ritambhara Singh<sup>3,5,\*</sup>

<sup>1</sup>School of Biological Sciences, University of Cambridge, UK

<sup>2</sup>School of Clinical Medicine, University of Cambridge, UK

<sup>3</sup>Department of Computer Science, Brown University, USA

<sup>4</sup>Altos Labs, Cambridge, UK

<sup>5</sup>Center for Computational Molecular Biology, Brown University, USA

<sup>6</sup>Department of Molecular Biology, Cell Biology and Biochemistry, Brown University, USA

\*Corresponding authors

*E-mail:* [lpd27@cam.ac.uk](mailto:lpd27@cam.ac.uk), [ritambhara@brown.edu](mailto:ritambhara@brown.edu)

**Background:** Aging involves intricate epigenetic changes, with histone modifications playing a pivotal role in dynamically regulating gene expression. Our research comprehensively analyzes seven key histone modifications across various tissues to understand their behavior during human aging and formulate age prediction models.

**Results:** These histone-centric prediction models exhibit remarkable accuracy and resilience against experimental and artificial noise. They showcase comparable efficacy when compared with DNA methylation age predictors through simulation experiments. Intriguingly, our gene set enrichment analysis pinpoints vital developmental pathways crucial for age prediction. Unlike in DNA methylation age predictors, genes previously recognized in animal studies as integral to aging are amongst the most important features of our models. We also introduce a pan-histone-mark, pan-tissue age predictor that operates across multiple tissues and histone marks, reinforcing that age-related epigenetic markers are not restricted to particular histone modifications.

**Conclusion:** Our findings underscore the potential of histone marks in crafting robust age predictors and shed light on the intricate tapestry of epigenetic alterations in aging.

**Keywords:** aging, prediction, histone modifications

## 1 Background

2 Aging is marked by noticeable changes mainly at cellular and organismal levels, encompassing phenomena like  
3 epigenetic disturbances, genomic instability, proteostasis loss, nutrient-sensing deregulation, and dysbiosis [1, 2].  
4 This understanding has spawned a variety of omics age predictors in fields such as epigenetics, transcriptomics,  
5 proteomics, metabolomics, and microbiotics [3, 4]. Most studies focus on on blood chemistry, transcriptomics,  
6 and DNA methylation, revealing several aging biomarkers, including those based on DNA methylation, telomere  
7 length, and proteomics [3, 5]. Blood tests, facilitated by deep neural networks, offer notable accuracy with a  
8 median absolute error of about five to six years [6, 7]. RNA sequencing contributes rich transcriptomic data  
9 for age predictors, which tend to be cell type or tissue-specific [8–12]. Cytosine methylation have emerged as  
10 the most favored molecular measurement for crafting age predictors that apply to all human tissues, with pan-  
11 tissue predictors achieving a median absolute error nearing four years [13–16]. Single-cell pan-tissue predictors  
12 utilizing DNA methylation have also been realized [17]. Newly, pan-mammalian epigenetic clocks applicable to  
13 all mammalian species have been presented [18]. Reflecting upon the achievements of age predictors centered  
14 on cytosine methylation, it beckons whether other epigenetic shifts, notably those anchored in histone levels,  
15 could engender mammalian aging clocks of similar precision. Clocks informed by histone marks carry potential  
16 relevance, resonating well with the histone code. Despite studies delineating the nuanced relationship between  
17 aging and histone marks [19, 20], a multifaceted mammalian age predictor rooted in histone mark data remains  
18 to be formulated.

19 To bridge this gap, we harness ENCODE data [21, 22] to analyze seven histone marks in human tissues  
20 and cells. We pinpoint a discernible shift from heterochromatin-linked modifications to those linked with  
21 euchromatin, corroborated by other studies [19, 23]. Furthermore, the variance in these histone modifications  
22 across genes escalates with age, hinting at an epigenetic regulatory decline. We introduce the first age predictors  
23 grounded in histone modification ChIP-Seq data. Impressively, their performance rivals that of DNA methylation  
24 age predictors, adjusted for training sample size. Our explorations divulge pivotal pathways and genes for  
25 age prediction. Developmental pathways and micro RNAs conspicuously dominate most histone modification  
26 age predictors. We also unearth that histone modifications can be broadly categorized as either activating  
27 or repressive for age predictor construction, irrespective of their unique roles. Interestingly, specific genes  
28 manifest consistent age-related trends across both these categories. Capitalizing on these insights, we present  
29 the inaugural pan-histone mark, pan-tissue age predictor.

30 To encapsulate, our research underscores the potential of histone modification data in age prediction. We  
31 emphasize the pivotal role of epigenetic regulation in aging and spotlight the prospective utility of histone

32 modifications as aging biomarkers. Concurrently, we illuminate key pathways and genes pivotal to epigenetic  
33 aging.

## 34 Results

### 35 Dynamics of histone mark during aging

36 In this study, we explore publicly available chromatin-immunoprecipitation sequencing (ChIP-Seq) human data  
37 from the ENCODE project [21, 22] (Figure 1). We focused on seven key histone modifications. Three, H3K4me3,  
38 H3K27ac, and H3K9ac, are broadly associated with euchromatin; another two, H3K9me3 and H3K27me3, are  
39 broadly associated with heterochromatin; one, H3K36me3, is associated with transcription elongation and  
40 heterochromatin; and one, H3K4me1, is associated with enhancers [24, 25].

41 Before any attempt to create age predictors based on histone modifications, we first analyzed data derived  
42 from human tissue from the ENCODE project to understand age-related dynamics [21, 22]. We obtained a  
43 total of 1814 samples ( $n$ ) from ChIP-Seq data for H3K4me3 ( $n = 359$ ), H3K27ac ( $n = 359$ ), H3K27me3 ( $n$   
44  $= 291$ ), H3K4me1 ( $n = 264$ ), H3K36me3 ( $n = 257$ ), H3K9me3 ( $n = 248$ ), and H3K9ac ( $n = 36$ ). The  
45 samples represent 82 tissues with ages ranging from embryonic to 90-plus years (Supplementary Figure 2a) with  
46 a roughly equal number of males and females (48.2% males, 50.6% females, 1.2% not available). Cancerous  
47 tissue constituted 0.9% of observations ( $n = 17$ ). The sequencing was performed with seven different Illumina  
48 instruments at four labs in universities across the United States. A summary of the relevant statistics can be  
49 found in Supplementary Figure 2 and attest to the breadth and diversity of the data collected.

50 We used the processed ChIP-Seq data files that display the probability that a genomic region is enriched  
51 for a specific histone mark compared to the control of sequencing DNA without the immunoprecipitation step.  
52 Effectively, each sample has a p-value for each single nucleotide in the genome. The lower the p-value, the  
53 higher the confidence that the locus contains the histone mark of interest. Given the high dimensionality of the  
54 data (high number of features  $p$  compared to the number of observations  $n$ ), with 3 billion nucleotides in the  
55 human genome, we decided to reduce the number of features by summarizing the values across genomic regions.  
56 We opted for averaging over the gene bodies of protein-coding and noncoding genes to facilitate interpretation  
57 unless stated otherwise (minimum, median, and maximum bin sizes of 7, 3897, and 2473538 bps respectively).  
58 The Homo Sapiens Ensembl annotation 105 provided the genomic locations [26]. To summarize the values  
59 in each gene as a single feature, we averaged the negative log<sub>10</sub> of p-values and then arcsinh-transformed to  
60 stabilize the variance (see Methods). In the end, the number of features was reduced by nearly 50 thousand,

61 from roughly 3 billion to 62241 per sample.

62 We plotted our data using uniform manifold approximation and projection (UMAP) to determine whether age  
63 was a major differentiating factor in our samples. As expected, each histone modification is generally separated  
64 from the others (Figure 2a). Interestingly, however, there are two main clusters, one with the activating marks  
65 H3K4me3, H3K27ac, and H3K9ac, and another with the repressive modifications H3K9me3, H3K27me3, and  
66 the elongation mark H3K36me3. H3K4me1, typically enriched in enhancers and neither clearly activating nor  
67 repressive, is located between the two clusters. Next, we colored the UMAP plot with age rather than histone  
68 marks (Figure 2b). While there is no apparent separation between young and middle-aged samples, old samples  
69 (>70 years) are relatively separated from the rest of the data. We confirmed this observation by replotting UMAP  
70 stratified by modification (Figure 2c). While it is much easier to differentiate a sample based on the type of  
71 histone mark, age is relevant enough — at least in old age — to contribute towards the UMAP projection. We  
72 also plotted the data with principal component analysis (PCA) to rule out any potential artifacts from UMAP,  
73 yielding similar results (Supplementary Figures 2d-f).

74 After broadly analyzing the data through low-dimensional projections, we focused on uncovering age-related  
75 trends. It has been widely reported that aging is accompanied by loss of heterochromatin and activation of  
76 constitutively repressed genes [19, 23]. Indeed, the mean signal of all three repressive histone modifications has  
77 a negative Pearson's correlation ( $r$ ) with age (Figure 2d), with H3K9me3 (Pearson's  $r = -0.35$ ,  $p$ -value =  $9.7e-9$ )  
78 and H3K27me3 (Pearson's  $r = -0.24$ ,  $p$ -value =  $4.9e-5$ ) reaching significance. Likewise, the mean signal of all  
79 three activating histone marks has a positive Pearson's correlation with age, with H3K4me3 (Pearson's  $r = 0.2$ ,  
80  $p$ -value =  $1.7e-4$ ) reaching significance. The mean signal of H3K4me1 barely changes with age, with Pearson's  
81  $r$  of only  $-0.01$ . Our results add to the evidence from several previous studies indicating loss of heterochromatin  
82 with aging.

83 Another interesting metric to track across aging is how variable the histone modification signal becomes,  
84 as previous studies have reported an increase in entropy during aging in DNA methylation [27–29]. While the  
85 entropy calculation for an unbound histone mark enrichment is not as straightforward as for DNA methylation,  
86 the signal variance normalized by its mean can give insights into the increased variability during aging. Interest-  
87 ingly, Pearson's correlation is indeed positive for all seven histone marks (Figure 2e), with H3K4me3 (Pearson's  
88  $r = 0.19$ ,  $p$ -value =  $2.4e-4$ ), H3K4me1 (Pearson's  $r = 0.18$ ,  $p$ -value =  $2.7e-3$ ), H3K9me3 (Pearson's  $r = 0.17$ ,  
89  $p$ -value =  $6.1e-3$ ), and H3K27me3 (Pearson's  $r = 0.19$ ,  $p$ -value =  $9.1e-4$ ) reaching significance. The broad  
90 increase in normalized signal variance with age suggests that any tight regulation to maintain histone marks to  
91 specific genomic regions becomes less effective with spillover to other loci.

92 While the data so far point towards robust age-related trends, we wanted to determine whether the genes'  
93 signals correlate with age. For each of the 62241 genes, we calculated the correlation coefficient with respect  
94 to age and plotted it on a histogram (Figure 2f). We chose Spearman's  $r$  over Pearson's  $r$  given that multiple  
95 age-related changes are non-linear, so a non-parametric coefficient is best suited to detect such correlations.  
96 As shown by the red shade, a surprisingly large proportion of genes significantly correlate with age ( $p$ -value  $<$   
97  $8.0 \times 10^{-7}$ , i.e.,  $p$ -value  $< 0.05$  with Bonferroni correction for the 62241 genes). As expected for the repressive  
98 histone marks, many more genes have a negative rather than positive coefficient. Surprisingly, however, the  
99 same is true for the three activating histone marks. Overall, given a large number of genes whose histone marks  
100 significantly correlated with age, it was likely that we could develop a histone mark age predictor from our data.

### 101 **Performance of pan-tissue histone mark age predictors**

102 Given the dynamics of histone modifications we observed during aging, we set out to create age predictors.  
103 Given the relatively small number of samples, we opted for a 10-fold nested cross-validation setup. Nine folds  
104 are used for internal 9-fold cross-validation to select the appropriate hyperparameters. Then, a predictor is  
105 trained on these nine folds with the best hyperparameter and is tested in the remaining external fold. This  
106 process is repeated ten times. It is worth emphasizing that samples originating from the same biosample were  
107 not split into different folds, as this may have artificially inflated the performance (Supplementary Figure 1).

108 To create an apt age predictor, we had three requirements: (1) the approach can suitably handle data  
109 with high dimensionality ( $p$  features  $\gg n$  samples), (2) the approach is robust to technical variation and  
110 experimental noise, and (3) the approach is easily interpretable. We introduced a feature-reduction step to  
111 fulfill the first requirement by training an ElasticNet model and selecting features with non-zero coefficients ( $p'$ )  
112 [30]. This is the only step in the overall age predictor with a hyperparameter ( $\lambda$ ), representing the strength  
113 of regularization. Moreover, having an initial set of reduced features allow us to easily interpret the most  
114 important genes and pathways through gene set enrichment analysis. For the second requirement, it has been  
115 shown that PCA can vastly improve the reliability of DNA methylation epigenetic age predictors by removing  
116 technical noise [31]. Therefore, we transformed the data using PCA calculated with a truncated support vector  
117 decomposition, generating  $(p' - 1)$  principal components. Finally, we used an automatic relevance determination  
118 regression (ARD), a form of regularized Bayesian regression, for the last requirement. It can easily be interpreted  
119 as, similarly to linear regression, each feature has a coefficient representing the model's weight. In addition, it  
120 provides an uncertainty value for each prediction. Choosing a different model to ElasticNet after feature selection  
121 and noise reduction also avoids the issue of double dipping. While we and others have previously shown that

122 deep learning can improve the performance and interpretation of pan-tissue DNA methylation epigenetic age  
123 predictors [14, 32], we opted for the aforementioned machine learning approach given the low number of samples.  
124 For more details of our modeling approach, see Methods.

125 Since chromatin immunoprecipitation is noisy and highly dependent upon the quality of the antibody [33],  
126 we expected a good but not impressive performance. We measured the performance of the age predictors using  
127 the following metrics: Pearson's correlation coefficient ( $r$ ), median absolute error (MAE), and root-mean-square  
128 error (RMSE). Surprisingly, all performed exceedingly well (Figure 3a), except for the H3K9ac age predictor -  
129 likely given the small sample size of 36. The H3K4me3 age predictor, in particular, was the best performer,  
130 with  $r=0.94$ , MAE= 4.31, and RMSE=8.74. Though the setup is not directly comparable, it is remarkable  
131 that this pan-tissue histone mark age predictor has similar reported performance compared to some of the most  
132 used DNA methylation age predictors [13, 27]. In summary, histone mark age predictors have remarkably low  
133 prediction errors.

134 It seems that the sample size is a significant determinant of the performance (Figure 3c), as it is highly  
135 negatively correlated with RMSE ( $r = -0.74$ ,  $p = 0.059$ ). It is possible that a larger sample size would make  
136 histone modification age predictors match or even exceed the performance of DNA methylation age predictors.  
137 Therefore, we downloaded all human tissue ChIP-seq samples imputed with Avocado [34]. The original 1814  
138 samples plus 1379 imputed samples were added to 3193 samples. We reran the nested cross-validation adding  
139 the imputed samples to the training sets. It has been suggested that the imputed signals contain enough  
140 biological information to be useful in several downstream analyses [35]. However, the performance was overall  
141 very similar for our age prediction tasks (Supplementary Figure 3g-l). This might suggest that imputed samples  
142 are unlikely to help our age predictors; perhaps the performance might be already saturated or the age-related  
143 changes are too subtle for Avocado to reconstruct.

144 In addition to testing the performance on data using features based on the average signal value over  
145 gene bodies, we also explored binning the ChIP-Seq data into (1) solely intergenic regions, (2) genes and  
146 intergenic regions (whole genome), (3) 20318 CpG dinucleotides common to the Illumina Methylation arrays  
147 27k, 450k, and EPIC, and (4) Horvath's 353 CpG sites from his pan-tissue DNA methylation age predictor [13].  
148 Despite the different lengths of the genomic loci, that should have not biased our data transformation since  
149 we average the signal over the entire bin. Given that heterochromatin is present mainly in noncoding genomic  
150 regions, we expected a better performance for the histone modification age predictors based on repressive histone  
151 modifications. This is observed — albeit with a minor improvement — for the H3K9me3 age predictor ( $r =$   
152  $0.78$  vs.  $0.74$ , MAE =  $8.67$  vs.  $8.20$ , RMSE =  $14.52$  vs.  $15.44$ ) and the H3K36me3 age predictor ( $r = 0.84$  vs.

153 0.80, MAE = 7.17 vs. 7.86, RMSE = 12.32 vs. 13.48) (Supplementary Figure 3a). The performance is virtually  
154 identical for the whole genome setting (Supplementary Figure 3a). The performance of the 20138 CpG sites  
155 is similar but slightly worse (Supplementary Figure 3c). Interestingly, it is particularly so for the age predictors  
156 of the repressive marks H3K9me3 ( $r = 0.66$  vs.  $0.74$ , MAE = 11.02 vs. 8.20, RMSE = 17.06 vs. 15.44),  
157 H3K27me3 ( $r = 0.90$  vs.  $0.92$ , MAE = 7.49 vs. 6.20, RMSE = 12.22 vs. 11.10), and H3K36me3 ( $r = 0.69$   
158 vs.  $0.80$ , MAE = 9.48 vs. 7.86, RMSE = 16.46 vs. 13.48). A similar trend is observed for the age predictors  
159 using the histone modification signal from Horvath's 353 CpG sites, though with overall poorer performance  
160 (Supplementary Figure 3d). The methylation status of these CpG sites may interfere with how much epigenetic  
161 information can be gained from that particular locus — for instance, if there is high histone acetylation, the  
162 gene is almost certainly active. Still, lack of histone mark repression does not mean the gene is active as it  
163 might be methylated (though CpG methylation and histone repression typically are well correlated). By binning  
164 at different places in the genome, we can see that age-related information is degenerate throughout the genome,  
165 as age predictors using inputs from completely different loci perform well. Nonetheless, there are specificities  
166 for the performance of histone mark age predictors given the function of the modifications, i.e., some marks  
167 perform slightly better or worse than others depending on the loci of the bins.

168 While it is impossible to directly compare the performance of DNA methylation epigenetic age predictor  
169 versus histone mark age predictors without paired data from the same sample, we attempted to make a rough  
170 comparison. For such, we ran 100 simulations by randomly drawing the same number of samples of each  
171 histone mark from the pan-tissue DNA methylation data set used to create AltumAge [14]. We subjected  
172 this random pool to the same nested cross-validation setup using the same machine-learning approach. In the  
173 end, we had performance metrics for 100 simulations of each sample size for DNA methylation age predictors.  
174 With these results, we compared how well each histone mark age predictor fell into the distribution of DNA  
175 methylation age predictors (Figure 3b). If a histone mark age predictor performs well, say over 90th percentile,  
176 it means that the reported metric was better than 90 out of the 100 simulations of a DNA methylation age  
177 predictor with the same number of samples. Overall, while the H3K9me3 and H3K36me3 age predictors are  
178 in the 0th percentile for RMSE, the H3K4me3, H3K9ac, and H3K27ac ones are in the 93rd, 92nd, and 67th  
179 percentile, respectively. Similar results were found for Pearson's correlation and MAE (Supplementary Figure  
180 3e,f). It is worth emphasizing that the AltumAge data was highly skewed towards younger ages, in which DNA  
181 methylation age predictors perform better, in contrast to our histone mark data's more uniform age distribution  
182 (Supplementary Figure 2a). Overall, the performance of the histone mark age predictors was approximately  
183 in line with the DNA methylation age predictors, with activating histone marks outperforming and repressive

184 histone marks underperforming.

185 Moreover, we assessed how robust and reliable our histone mark age predictors are. Our data contained  
186 samples with biological triplicates, which allowed us to analyze how reliable each histone mark age predictor is  
187 to experimental noise. Except for H3K4me1 and H3K36me3, all other histone mark age predictors showed an  
188 intraclass correlation coefficient above 0.9 (Figure 3e). In addition to assessing the reliability of the models to  
189 experimental variation, we wanted to test how they performed under the addition of artificial noise. For each  
190 test fold in the nested cross-validation, we added random Gaussian noise to the test data with up to 1.5 standard  
191 deviations in 0.3 standard deviation increments. Even in the most extreme scenario with 1.5 standard deviations,  
192 most models' performance remained similar, except for the H3K4me3 age predictor (Figure 3f). These results  
193 show that the age predictors are robust and reliable to experimental and artificial noise.

194 Lastly, we hypothesized that our choice of the uncertainty-aware ARD regression would give us insights into  
195 the epigenetic drift that occurs over time (Supplementary Figure 3g). We expected to see an increased model  
196 uncertainty with the sample's age. Though relatively weak correlations, we did notice a statistically significant  
197 relationship for the following age predictors: H3K36me3 ( $r = 0.14$ ,  $p = 0.025$ ), H3K4me3 ( $r = 0.11$ ,  $p =$   
198  $0.033$ ), and H3K9me3 ( $r = 0.14$ ,  $p = 0.029$ ). Despite the weak correlations, these findings show that the age  
199 predictors might be learning the well-described phenomenon of epigenetic drift.

## 200 **Inference of pan-tissue histone mark age predictors**

201 While DNA methylation age predictors have been the most used tools to measure age, the insights gained from  
202 them into what constitutes epigenetic aging are limited. The most important genes based on the location of  
203 the relevant CpG sites are often difficult to relate to the rest of the aging literature. Therefore, we sought to  
204 carefully analyze the genes that comprise our histone mark age predictors, i.e., the genes selected after the first  
205 step with ElasticNet. In the previous nested cross-validation setup, we trained 10 models in total. However,  
206 to be able to interpret the findings more clearly, we ran a single 10-fold cross-validation to choose the best  
207 hyperparameter  $\lambda$  and trained a single histone modification age predictor with the entirety of the data for each  
208 histone.

209 First, we began by visualizing an upset plot with all histone mark age predictors except for H3K9ac, given its  
210 low sample size and poor performance (Figure 4a). The models selected a subset from 341 genes for H3K9me3  
211 up to 1275 for H3K27ac. As expected, the selected genes were often shared across similar histone marks.  
212 For instance, the two marks with the most genes in common were H3K27ac and H3K4me3, and the three  
213 marks with the most genes in common were H3K27ac, H3K4me3, and H3K4me1. Surprisingly, though few,



214 some selected genes were common across both activating and repressive histone mark age predictors. For the  
215 principal components derived from the age predictor genes, the ARD regression decreased the coefficients of  
216 only a small subset of genes (about 5%) to zero (Supplementary Figure 4b). This indicates that only a tiny  
217 fraction of the 62241 features is sufficient to predict age.

218 Next, we investigated which pathways were important for the histone mark age prediction. Gene set enrich-  
219 ment analysis (GSEA) can reveal important gene ontology processes which are either over or underrepresented  
220 in a set of genes. Using Panther DB [36], we ran seven GSEAs, one for each set of selected genes from the  
221 histone mark age predictors (Figure 4b-h). Several developmental pathways are overrepresented in the top 10  
222 GO biological processes. Some examples are regionalization (H3K4me3, H3K27ac, H3K9ac, H3K27me3), pat-  
223 tern specification process (H3K4me3, H3K9ac, H3K27me3), anterior/posterior pattern specification (H3K4me3,  
224 H3K27ac, H3K27me3), anatomical structure morphogenesis (H3K27ac, H3K36me3). This brings further ev-  
225 idence to the fact that aging is simply a maladaptive continuation of development, which are fundamental  
226 for fitness early in life but whose continuation result in organismal decay. For H3K4me1, processes related to  
227 telomere organization, muscle regeneration, and immune response are heavily overrepresented. For H3K9me3,  
228 processes related to H3K27me3 regulation and fat proliferation are the most important.

229 Complementary to the GSEA, we also looked into which ENSEMBL gene biotypes were over or underrep-  
230 resented in each model (Supplementary Figure 4a). For such, we used Fisher's exact test with a Bonferroni  
231 correction for the number of gene biotypes ( $n = 39$ ). Some general trends emerge, with an overrepresenta-  
232 tion of micro RNAs (all histone marks), small nuclear RNAs (H3K9me3, H3K27me3, H3K4me1, H3K36me3,  
233 H3K27ac), small nucleolar RNAs (H3K9me3, H3K4me3, H3K4me1, H3K36me3, H3K27ac), and miscellaneous  
234 RNAs (H3K9me3, H3K4me1, H3K36me3, and H3K27me3). Micro, small nucleolar, and small nuclear RNAs  
235 have been linked to several age-related phenomena [19, 23]. Protein-coding genes are vastly underrepresented  
236 in the age predictors for all histone modifications besides H3K9ac.

237 Following the gene set analysis, we looked into the importance of individual genes. We focused on the top  
238 three protein-coding genes with the highest positive and negative contributions toward the final age prediction for  
239 each age predictor (Figure 4i). Calculating the individual gene importance is possible by inverse-transforming  
240 the coefficients of principal components from the ARD regression back into coefficients for the genes. The  
241 overarching theme was the importance of histone-coding genes. Amongst those are H1-1, H2AC15, H2BC8,  
242 H3C7, H3C11. For activating histone marks, histone genes had a negative coefficient (more histone enrichment  
243 translates to lower predicted age), whereas, for the repressive histone modifications, the opposite was true.  
244 These genes represent the components of the nucleosome histones H2A, H2B, and H3, and linker histone H1,

245 all of which have been linked to aging [19, 23, 37]. Other relevant age-related genes which contribute towards  
246 our histone mark age predictors are NOG, which plays an important role in early development in all germ layers;  
247 HOXD8, important in body patterning; TXNIP, whose inhibition can protect against age-related Alzheimer's  
248 disease in mice and whose upregulation causes oxidative stress [38, 39]; PER1 and PER3, circadian-clock genes  
249 that are known to be involved in aging [40, 41]; TBX3, which is highly expressed in embryonic stem cells  
250 and facilitates cellular reprogramming [42]; TNFSF9, which skews hematopoiesis during aging [43]; BTG2,  
251 which drives senescence [44]; CH25H, whose upregulation contributes to the development of osteoarthritis and  
252 inflammation in obesity and diabetes [45, 46]. In contrast to several DNA methylation age predictors, histone  
253 mark age predictors are enriched in factors that are known to play a role in aging.

### 254 **Creation of a pan-histone-mark pan-tissue age predictor**

255 With the results, we had some indications that age predictors trained with similar histone marks behave alike.  
256 UMAP clusters activating and repressive marks together (Figure 2a); there is a remarkable correlation between  
257 the performance of the age predictors with the number of training samples (Figure 3c); some genes are used by  
258 the same models (Figure 4b), and so are similar pathways (Figure 4c-h); the age predictors are generally enriched  
259 in similar gene biotypes (Supplementary Figure 4a). Thus, we set out to test whether a pan-histone-mark age  
260 predictor with reasonable performance was viable.

261 First, we made a grid plot contrasting the distribution of Spearman's correlation between age and each gene  
262 for every two histone marks (Figure 5a). As expected, there is usually a positive correlation between activating  
263 histone marks and, similarly, between repressive ones. The opposite is true when comparing an activating to  
264 a repressive mark. Nevertheless, looking at the density plots, some genes appear to have similar Spearman's  
265 correlation even when contrasting activating and repressive histone modifications. Second, we sorted the protein-  
266 coding genes by the highest positive and negative overall Spearman's correlation with age (Supplementary Figure  
267 5b). Indeed, several histone marks display the same age-related trends. This is because of the generality of the  
268 trend in some genes, despite the mainstream assumption that activating and repressive histone marks change  
269 in opposite directions. This further supports the hypothesis that a histone mark age predictor trained on one  
270 type of histone modification could plausibly predict age with another histone mark as input.

271 Next, we reran our nested cross-validation but rather than only predicting the histone mark's test fold  
272 of interest, we instead predicted the test fold of all histone marks. As expected, each histone mark age  
273 predictor performs the best when predicting the age of the histone mark with which it was trained (Figure  
274 5b). Interestingly, however, several histone mark age predictors can use other similar histone marks as input

275 with decent performance. For instance, the age predictor trained on H3K27ac performs well ( $r > 0.75$ ) when  
276 presented with H3K4me3 and H3K9ac data. The age predictor trained on H3K4me3 performs similarly well  
277 when presented with H3K27ac and H3K9ac data. Conversely, the age predictor trained on H3K27me3 has highly  
278 negative correlations when predicting the age with the three activating marks.

279 To test whether the models simply capture “on” or “off” information of a gene rather than a histone-specific  
280 signature, we reran the nested cross-validation, flipping the input sign to the ARD regressor. If that was the case,  
281 then an age predictor trained on an activating histone mark could theoretically predict the age for a repressive  
282 histone mark with the negative of the input signal and vice-versa. This is true despite less significant correlation  
283 values across the board (Supplementary Figure 5a). For instance, the flipped H3K4me3 age predictor can predict  
284 on H3K27me3 data ( $r > 0.6$ ), and the flipped H3K27me3 age predictor can predict on H3K27ac ( $r > 0.6$ ) and  
285 H3K4me3 ( $r > 0.35$ ). This shows that the ChIP-seq signal can be viewed as a sliding scale from repression to  
286 activation for some age predictors. This led us to rerun the nested cross-validation, using for training either all  
287 activating or all repressive marks plus H3K36me3 — as it is also associated with heterochromatin [47]. Indeed,  
288 activating and repressive histone age predictors perform exceedingly well ( $r > 0.8$ ) on activating and repressive  
289 histone marks, respectively.

290 However, one of the most striking observations is that some histone mark age predictors — H3K4me1 and  
291 H3K9ac — have positive correlations for six of the seven histone marks. This led us to believe that while  
292 some genes contribute towards age prediction based on information akin to a sliding scale of activation and  
293 repression, others must supply information differently. Likely, it would be based on age-related cross-histone  
294 epigenetic information, which would make creating a pan-histone pan-tissue histone modification age predictor  
295 viable. Therefore, we reran the nested cross-validation, using all seven types of histone mark for training or  
296 testing, again separating the folds by biosample. The performance for all seven histone marks roughly matches  
297 the one from the histone mark-specific age predictor, if not slightly better (Figure 5c). The pan-histone age  
298 predictor has a Pearson’s  $r$  of 0.87, MAE of 6.65 years, and RMSE of 12.09 years). To further test the  
299 generalizability of our pan-histone age predictor, we tested it in untouched data thus far from 568 primary cells  
300 spanning 12 histone marks taken from the ENCODE project (Figure 5d). The performance of an age predictor  
301 trained on tissues is expected to drop given that the *in vitro* milieu and passaging can induce changes akin to  
302 aging [48]. The model’s performance on these cultured cells is still significant ( $r = 0.53$ , MAE = 12.57, RMSE  
303 = 19.01), as the performance of well-known DNA methylation age predictors is modified by several *in vitro*  
304 variables [49].

305 Next, we looked into some critical pathways and genes for the pan-histone age predictor. Moreover, using

306 Panther DB [36], we analyzed the gene ontology biological processes which were either over or underrepresented  
307 in the set of genes selected with the ElasticNet step (Figure 5f). Most gene ontology terms are related to de-  
308 velopmental and transcriptional processes. Like histone-specific age predictors, the gene set is underrepresented  
309 in protein-coding genes and overrepresented in micro, miscellaneous, small nuclear, and small nucleolar RNA  
310 ( $p$ -adjusted  $< 0.0001$ , Figure 5g). Similar genes appear in the histone-specific age predictors when looking at  
311 the protein-coding genes with the highest positive and negative contribution to the pan-histone age predictor  
312 (Figure 5h). Amongst those are H4C7, a component of nucleosomes; HOXD4, important in body patterning;  
313 NR1D1 and PER3, involved in regulating circadian rhythms. Curiously, several olfactory receptor genes also  
314 appeared important for age prediction. Recently, this sense has been implicated in lifespan regulation in worms  
315 [50].

## 316 Discussion

317 The allure of DNA methylation age predictors, given their remarkable precision, has significantly redirected  
318 the focus in aging research towards epigenetics. Previously, the foundation for epigenetic age predictors pre-  
319 dominantly rested on DNA methylation. However, recent advancements have led to the inception of a model  
320 grounded in chromatin accessibility, which has shown promising results [51]. Until this study, a predictor based  
321 on histone marks was a void in the research domain. Our work elucidates that ChIP-Seq data, derived from seven  
322 histone marks, can underpin the formulation of highly precise age predictors. A simulated analysis indicates that  
323 histone mark age predictors, given comparable sample sizes, could potentially surpass their DNA methylation  
324 counterparts in terms of accuracy, especially concerning activating histone marks. Fundamentally, our findings  
325 champion the age-associated dynamics of histone modifications as potent aging biomarkers, integral for devising  
326 resilient age predictors.

327 While DNA methylation age predictors are commendable in performance, their interpretability often remains  
328 obscure. Crucial CpG sites integral to these models often pertain to genes of elusive function or ones whose  
329 aging impact is dubious. For instance, a recent endeavor to construct a lifespan DNA methylation predictor  
330 highlights a CpG site near BCL11B as vital [52]. However, its attenuation scarcely affects the epigenetic age  
331 in mice. Another CpG site related to the ELOV2 gene, crucial in fatty acid metabolism, while having potential  
332 health implications, bears a tenuous link with aging. In contrast, our histone mark age predictors teem with well-  
333 established aging biology. These models underscore pivotal genes associated with development, inflammation,  
334 senescence, and stem cell sustenance, among others. Interestingly, some pathways enriched in our predictors  
335 mirror those in a recent accelerated aging mouse model [53].

336 Moreover, our results pave the way for formulating intriguing longevity interventions. Although the signifi-  
337 cance of genes in our models is inherently correlative, it would be fascinating to explore interventions targeting  
338 these gene modifications. For example, yeast lifespan has been shown to extend upon overexpressing certain his-  
339 tone proteins [54]. Additionally, medications influencing histone modifications have demonstrated their potential  
340 in governing health and lifespan in specific organisms. Examples include NAD<sup>+</sup> precursor supplementation im-  
341 pacting sirtuin enzymes [55, 56] and the effects of alpha-ketoglutarate on certain histone demethylases [57–59].  
342 Recent revelations suggest that compounds influencing histone modifications can rejuvenate cells *in vitro* [60,  
343 61]. Thus, our models' top hits and interventions affecting histone modifications might chart innovative paths  
344 in aging research.

345 One of this paper's pivotal contributions is the demonstration that a unified age predictor can adeptly  
346 decipher age from various epigenetic data types. Merging different data modalities for a unified model has  
347 precedent [51]. Yet, our results suggest distinct histone marks encapsulate consistent age-related data. A case  
348 in point is the effectiveness of a predictor trained on H3K4me3 data when applied to H3K27ac data. We have  
349 crafted a pan-histone, pan-tissue age predictor compatible across numerous histone marks both in tissue and *in*  
350 *vitro* primary cells. Remarkably, this predictor can seamlessly assimilate multiple data types without necessitating  
351 model weight adjustments, underlining the significance of age-related patterns over specific genomic locus  
352 changes.

## 353 Conclusion

354 In summary, this study presents a holistic inspection of age-related changes in histone marks. We introduce  
355 innovative age predictors, underscoring the potential to utilize a solitary model capable of handling varied  
356 data modalities for age prediction. A limitation to consider is the dearth of quality ChIP-Seq data beyond  
357 the ENCODE project. High-quality ChIP-Seq data is scarce compared to DNA methylation and is needed for  
358 the appropriate processing of the bigWig files necessary for our models. Optimally, our predictors would be  
359 cross-validated using samples adhering to ENCODE guidelines. Nonetheless, the prohibitive cost of ChIP-Seq,  
360 compared to the prevalent methylation arrays, has constrained the creation of expansive datasets at this study's  
361 juncture. We are optimistic about future experiments building upon the insights this paper offers.

## 362 Methods

### 363 Data

364 To generate and interpret the age predictors, we collected 1814 human tissue ChIP-Seq samples from the  
365 ENCODE project in the bigWig format [21, 22] for the histone modifications H3K4me3 (n = 359), H3K27ac (n =  
366 359), H3K27me3 (n = 291), H3K4me1 (n = 264), H3K36me3 (n = 257), H3K9me3 (n = 248), and H3K9ac (n  
367 = 36). The data was generated from Bradley Bernstein's lab at the Broad Institute, John Stamatoyannopoulos's  
368 lab at UW, Joseph Costello's lab at UCSF, and Bing Reng's lab at UCSD.

369 We generated a feature matrix by averaging the negative log<sub>10</sub> of p-values signal across all nucleotides per  
370 gene body according to the Homo Sapiens Ensembl annotation 105 [26]. Of note, the p-values are already  
371 -log<sub>10</sub>-transformed in the ENCODE bigWigs. Then, these averages were arcsinh-transformed. The summary of  
372 the transformation is as follows:

$$\text{Histone mark enrichment at locus A} = \sinh^{-1} \left( \frac{1}{\text{len}(\text{locus A})} \sum_{j=i}^{i+\text{len}(\text{locus A})} -\log_{10}(\text{p-values}_j) \right)$$

373 Samples in which more than 10% of the features were unavailable were discarded. Then, missing values for  
374 features were encoded as 0. In Section 2 specifically, we also tested averaging the negative log of p-values signal  
375 across intergenic regions, genes, and intergenic regions (whole genome), 20318 CpG dinucleotides common to  
376 the Illumina Methylation arrays 27k, 450k, and EPIC, and Horvath's 353 CpG sites from his pan-tissue DNA  
377 methylation age predictor [13].

378 Embryonic samples had their age encoded as  $\frac{7 \times (w-40)}{365}$ , where  $w$  is the gestational week, and 40 is the  
379 number of weeks of a normal gestation [62]. Therefore, some samples had negative ages. Some samples whose  
380 age was described as 90-plus by the ENCODE project, likely for anonymity reasons, were encoded as 90 years.

381 The distribution of the data is represented in Supplementary Figure 2.

382 Moreover, to test the performance *in vitro* of our pan-histone pan-tissue age predictor, we gathered another  
383 568 samples of primary cells spanning 12 histone marks (H2AFZ, H3K4me1, H3K4me2, H3K4me3, H3K9ac,  
384 H3K9me1, H3K9me3, H3K27ac, H3K27me3, H3K36me3, H3K79me2, H4K20me1) taken from the ENCODE  
385 project [21, 22].

386 Lastly, to assess the tentative increase in performance that could arise from increasing the sample size  
387 through imputation, we further downloaded all available Avocado-imputed samples from ENCODE for the  
388 seven prominent histone marks we analyzed in the paper, 197 of each, totaling 1379 samples.

## 389 Age predictor performance evaluation

390 For all experiments assessing the performance of the age predictors, the setup consisted of 10-fold nested cross-  
391 validation with observations from the same biosample remaining in the same folds to not artificially boost the  
392 performance. Nine folds are used for internal 9-fold cross-validation to select the appropriate hyperparameter.  
393 The nine folds are used for training the model, which is tested in the remaining external fold. Cancer samples  
394 were always removed.

395 Our age predictors consist of three steps. The feature selection method employs an ElasticNet model with a  
396 0.9 L1 to L2 proportion to choose  $p$  features whose absolute coefficient is above zero. The only hyperparameter  
397  $\lambda$  in the age predictors is the strength of regularization of the feature-selection ElasticNet, which controls  
398 how many variables are chosen. The tested hyperparameter values were  $\lambda \in \{0.1, 0.05, 0.01, 0.005, 0.001\}$ .  
399 Secondly, to improve the robustness of the age predictor, we applied PCA calculated with a truncated support  
400 vector decomposition to generate  $(p' - 1)$  principal components. Finally, we used an automatic relevance  
401 determination regression [63], a form of regularized Bayesian regression, to predict age. All of the steps above  
402 were created with `sklearn` library in Python. If not mentioned otherwise, all other hyperparameters were the  
403 standard ones in the package.

404 To roughly compare the performance of using histone mark data to DNA methylation data [14], we used  
405 AltumAge's pan-tissue dataset and pooled 100 random samples of the same size as the sample size for each  
406 histone modification. Then, we ran the same three steps we used to create our histone modification age  
407 predictors and ended up with 100 values for each performance metric. With this information, we compared  
408 where the performance of the histone mark age predictors fell within the distribution for the DNA methylation  
409 age predictors taking sample size into account. However, it must be emphasized that the ENCODE ChIP-seq  
410 dataset does not have the same age and tissue distribution as AltumAge's DNA methylation dataset.

411 In the setup to assess the performance of our models with the addition of the Avocado-imputed samples,  
412 whenever a model was trained on a particular histone mark, the imputed samples were added to the training  
413 set (with tumors removed) for both each inner fold of the nested cross-validation and when training the model  
414 to predict each test set of the cross-validation. The imputed samples were only added to try to improve the  
415 performance of predicting the age of actual samples rather than attempting to predict the age of imputed  
416 samples.

## 417 **Age predictor interpretation**

418 To interpret a particular age predictor, we ran a 10-fold cross-validation to select the best regularization hyper-  
419 parameter in the ElasticNet feature selection step. Then, we trained the age predictor using the entirety of the  
420 data in each setting. For the gene set enrichment analysis, we selected the genes that passed feature selection:  
421 H3K4me3 ( $p' = 1240$ ), H3K27ac ( $p' = 1275$ ), H3K27me3 ( $p' = 922$ ), H3K4me1 ( $p' = 892$ ), H3K36me3 ( $p'$   
422  $= 870$ ), H3K9me3 ( $p' = 341$ ), H3K9ac ( $p' = 102$ ), and pan-histone ( $p' = 3739$ ). To determine the individual  
423 contribution of each feature towards the final age prediction, we inverse-transformed the coefficients of principal  
424 components from the automatic relevance determination regression back into coefficients for the genes.

## 425 **Statistical Analysis**

426 The p-values associated with Pearson's and Spearman's correlation coefficients were obtained using the functions  
427 `pearsonr` and `spearmanr` from the python package `scipy`.

428 To create the UMAP and PCA plots in Figure 2 and Supplementary Figure 2, we ran Python's `dynamo`  
429 package function `dyn.tl.reduceDimension` with either UMAP or PCA as the basis with standard parameters.

430 The intraclass correlation coefficient measures how well multiple measurements agree with one another and  
431 is used to assess model reliability. To calculate the intraclass correlation coefficient, we used Python's `pingouin`  
432 package with a single-rater, absolute-agreement, two-way random-effects model per [64] guidelines and similarly  
433 to [31].

434 Notebooks with the analyses and a complete list of the python package versions are also available on our  
435 GitHub ([URLXXXX](#)).

## 436 **Declarations**

### 437 **Ethics approval and consent to participate**

438 Not applicable.

### 439 **Consent for publication**

440 Not applicable.



## 441 **Availability of data and materials**

442 After publication, the code to rerun our results will be available on our GitHub (URLXXXX). It takes about three  
443 weeks for all scripts to run on an Amazon AWS ml.t3.2xlarge instance.

444 The data will be available after the publication on Zenodo, for review we have provided the datasets on  
445 Google Drive (URLXXXX).

## 446 **Competing Interests**

447 L.P.D.L.C. was a part-time employee and a share-option holder of Shift Bioscience Ltd during part of the  
448 development of this manuscript.

## 449 **Author Contributions**

450 L.P.D.L.C. conceived of the presented idea, devised the methodology, ran experiments, and wrote the first draft  
451 of the manuscript. M.H.A. ran experiments and collected the datasets. S.H. assisted with the analysis and  
452 biological interpretation of the results and contributed to the final manuscript. E.L. assisted with the analysis  
453 and biological interpretation of the results and contributed to the final manuscript. R.S. supervised the project,  
454 devised the methodology, and contributed to the final manuscript.

## 455 **Funding**

456 No funding sources are reported for this work.

## 457 **Acknowledgements**

458 The title of this paper was inspired by [13].

## 459 Figure Legends

Figure 1: Schematic showing the main steps of creating the histone mark age predictors. First, bigWig files containing the enrichment p-values per nucleotide of chromatin immunoprecipitation samples were gathered from the ENCODE project. Then, the dimensionality of the data is reduced by summarizing the p-values into bins that represent the signal. Then, modeling is done through the application of an ElasticNet for feature selection followed by Automatic Relevance Determination regression to predict age. Image created with BioRender.

Figure 2: Age-related changes in seven histone modifications across human tissues. (a, b) Uniform manifold approximation and projection (UMAP) of the dataset containing the average signal of 62241 genes for 1814 samples grouped by histone modification (a) and age (b). (c) UMAP colored by age for each histone modification. (d, e) Linear regression plot with 95% confidence interval based on 1000 bootstraps for the signal mean (d) and signal variance (e) normalized by signal mean over age for each histone modification. (f) Histogram of Spearman's correlation for each 62241 features across age per histone modification. Bins shaded in red represent statistically significant correlations ( $p < 0.05$  with Bonferroni's correction).

Figure 3: Performance of pan-tissue histone mark age predictors. (a) Scatter plot of the predicted age versus real age of each histone mark age predictor using genes as features. Each of the 10 test folds of the nested cross-validation is shown in a different color. A dotted black line representing  $x=y$  is shown alongside a colored, solid regression line with its 95% confidence interval based on 1000 bootstraps. (b) Histogram of the root-mean-square error (RMSE) for age predictors trained on 100 random samples pooled from AltumAge's DNA methylation dataset with the same number of samples as each histone mark. A colored, vertical line shows where in the RMSE distribution the age predictor trained with the histone mark data would lie. (c) Scatter plot of the RMSE of each histone mark age predictor against sample size with a regression line with its 95% confidence interval based on 1000 bootstraps. (d) Three-dimensional scatter plots for samples that were done in triplicates. A dotted black line representing  $x=y=z$  is shown alongside a colored, solid regression line with its 95% confidence interval based on 1000 bootstraps. (e) Bar plot showing the intraclass correlation coefficient with error bars representing 95% confidence interval for each histone mark. (f) Point plot with 95% confidence interval based on 1000 bootstraps of the mean absolute error of each histone mark age predictor under added artificial Gaussian noise.

Figure 4: Inference of pan-tissue histone mark age predictors. (a) Upset plot for the subset of genes selected for six of the seven histone modification age predictors. Top 10 gene ontology biological processes from gene set enrichment analysis from Panther DB of H3K4me3 (b), H3K27ac (c), H3K9ac (d), H3K4me1 (e), H3K36me3 (f), H3K27me3 (g), and H3K9me3 (h). (i) Top 3 protein-coding genes for each histone mark age predictor with positive and negative coefficients.

Figure 5: Creation of a pan-histone-mark pan-tissue age predictor. (a) Grid plot with Spearman's correlation of the histone modification signal for a particular gene over age. Density plots, histograms, and regression lines show the direction of the correlation. (b) Bubble plot of Pearson's correlation coefficient when age predictors are trained on different histone marks from the ones they aim to predict. Scatter plot of the predicted age of the pan-histone-mark age predictor versus real age using genes as features stratified by histone modification (c) or grouped together (d). Each of the 10 test folds of the nested cross-validation is shown in a different color. A dotted black line representing  $x=y$  is shown alongside a colored, solid regression line with its 95% confidence interval based on 1000 bootstraps. (e) Similarly, a scatter plot of the pan-histone-mark age predictor trained using all of the tissue sample data to predict the age of primary cells from 11 different histone marks. Each color represents a distinct histone modification. (f) Top 20 gene ontology biological processes from gene set enrichment analysis from Panther DB for the selected genes from the pan-histone-mark predictor. (g) Bar plot with the proportion of ENSEMBL's gene biotype for the selected genes in each histone mark age predictor. P-values were rectified with Bonferroni's correction (\*,  $p < 0.01$ ; \*\*,  $p < 0.001$ ; \*\*\*,  $p < 0.0001$ );. (h) Top 10 protein-coding genes for the pan-histone-mark age predictor with positive and negative coefficients.

Supplementary Figure 1: Setup for the nested cross-validation of the different histone mark age predictors. First, the data is split into ten folds. Each fold is divided into training and validation (9/10 of the data) and test (1/10 of the data), always maintaining observations from the same biosample. There is another 9-fold cross-validation using the training and validation part of the data for hyperparameter tuning. When the best hyperparameter is found, the age predictor is trained on the entire training and validation data and used to predict the remaining test set. After this is done for all ten folds, there is a prediction for each observation in the entire data set. Therefore, performance metrics can be calculated. Image created with BioRender.

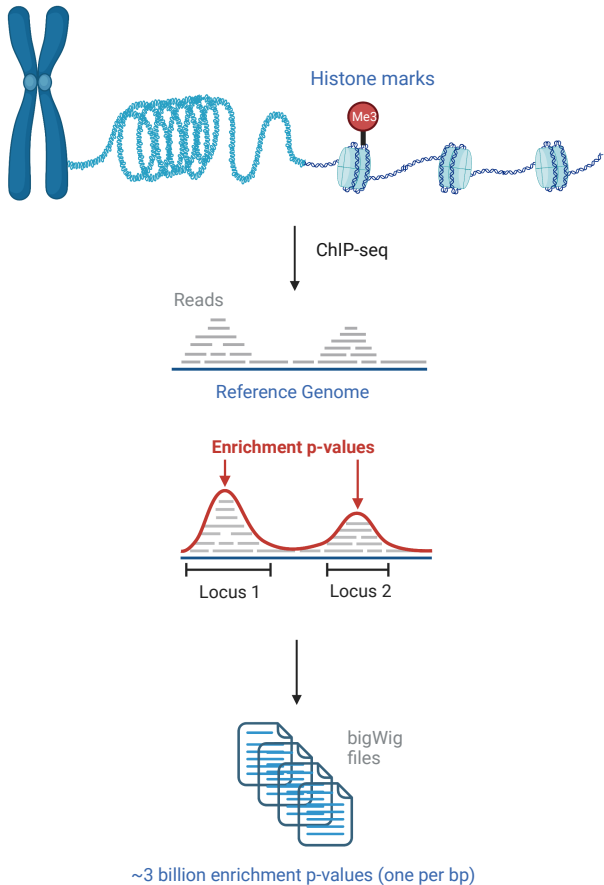
Supplementary Figure 2: (a) Histogram of the number of samples per histone modification over age. (b) Histogram of the density of transformed signal values per gene. (c) Pie plots with the distribution of the 1814 ChIP-Seq samples by platform, run end, tissue, histone, mapped read length, sex, lab, project, disease status, cancer status, fragmentation method, and library size. (d, e) Principal component analysis (PCA) plot of the dataset grouped by histone modification (d) and age (e) showing the first and second principal components (9.3% and 6.1% of the explained variance, respectively). (f) PCA colored by age for each histone modification.

Supplementary Figure 3: Scatter plot of the predicted age versus real age of each histone mark age predictor using (a) intergenic regions, (b) genes and intergenic regions (whole genome), (c) 20318 CpG dinucleotides common to the Illumina Methylation arrays 27k, 450k, and EPIC, and (d) Horvath's 353 CpG sites from his pan-tissue DNA methylation age predictor [13] as features. Each of the 10 test folds of the nested cross-validation is shown in a different color. A dotted black line representing  $x=y$  is shown alongside a colored, solid regression line with its 95% confidence interval based on 1000 bootstraps. Histogram of Pearson's correlation (e) and median absolute error (MAE) (f) for age predictors trained on 100 random samples pooled from AltumAge's DNA methylation dataset with the same number of samples as of each histone mark. A colored, vertical line shows where in the RMSE distribution the age predictor trained with the histone mark data would lie. (g) Scatter plot of the predicted standard deviation versus real age of each histone mark age predictor using genes as features. Each of the 10 test folds of the nested cross-validation is shown in a different color. A colored, solid regression line with its 95% confidence interval based on 1000 bootstraps. Scatter plot of the predicted age versus real age of each histone mark age predictor trained in addition with Avocado-imputed samples [34] using (h) gene bodies, (i) intergenic regions, (j) genes and intergenic regions (whole genome), (k) 20318 CpG dinucleotides common to the Illumina Methylation arrays 27k, 450k, and EPIC, and (l) Horvath's 353 CpG sites from his pan-tissue DNA methylation age predictor [13] as features. Each of the 10 test folds of the nested cross-validation is shown in a different color. A dotted black line representing  $x=y$  is shown alongside a colored, solid regression line with its 95% confidence interval based on 1000 bootstraps.

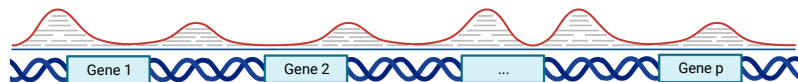
Supplementary Figure 4: (a) Bar plot with the proportion of ENSEMBL's gene biotype for the selected genes in each histone mark age predictor. P-values were rectified with Bonferroni's correction (\*,  $p < 0.01$ ; \*\*,  $p < 0.001$ ; \*\*\*,  $p < 0.0001$ ;). (b) Doughnut plots for each histone modification age predictor. On the left, the proportion of principal components whose coefficients were positive (yellow), zero (gray), or negative (blue) is displayed; on the right, the weight of each gene to the total prediction is shown, with positive genes with positive coefficient in yellow and negative in blue.

Supplementary Figure 5: (a) Bubble plot of Pearson's correlation coefficient when age predictors are trained on certain histone marks and attempt to predict others but with the negative value of the input to the ARD regression part of the age predictor model. (b-g) Regression plots of age versus histone mark signal values for six genes generally go up or down with age. Shaded is the 95% regression confidence interval based on 1000 bootstraps.

## STEP ① Chromatin Immunoprecipitation data from the ENCODE project



## STEP ② Dimensionality reduction

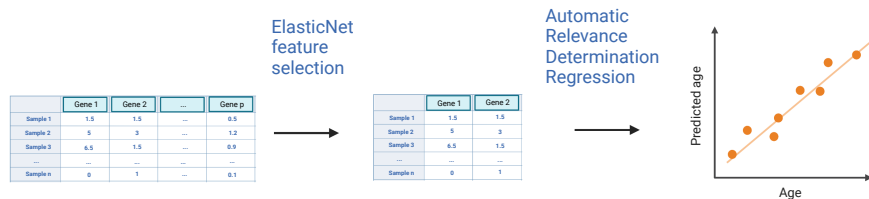


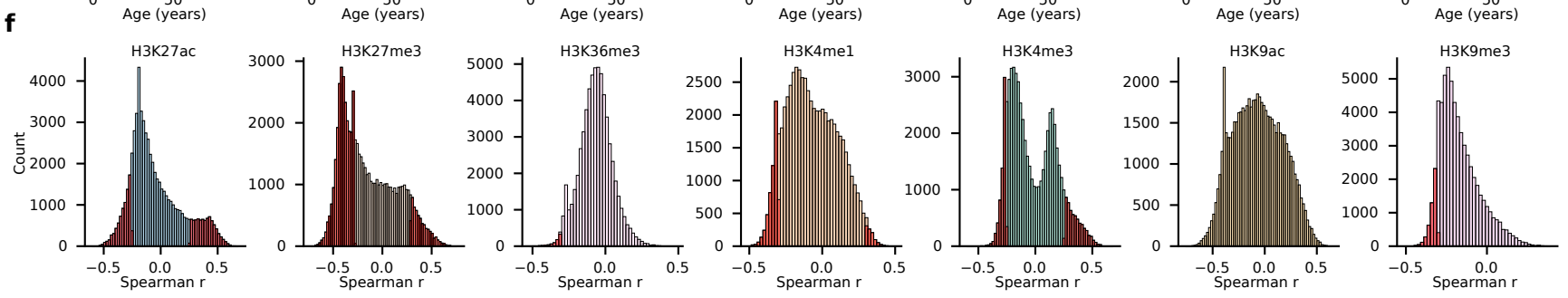
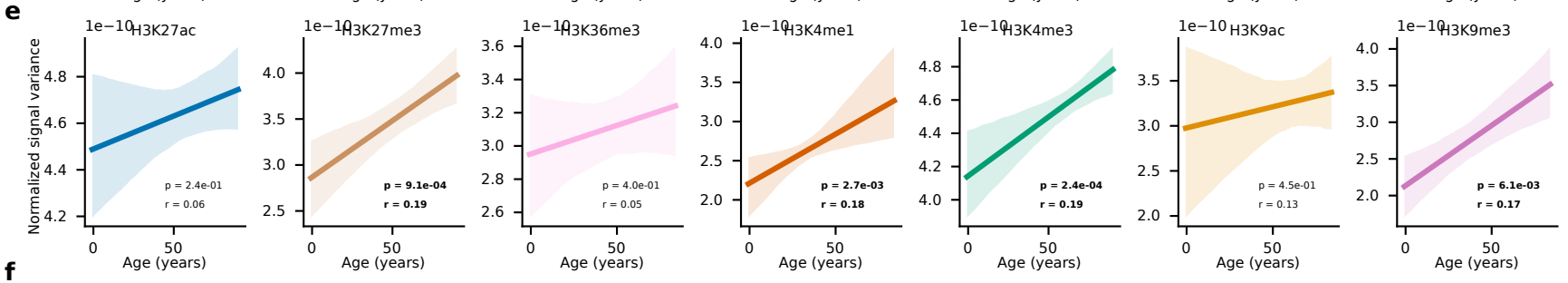
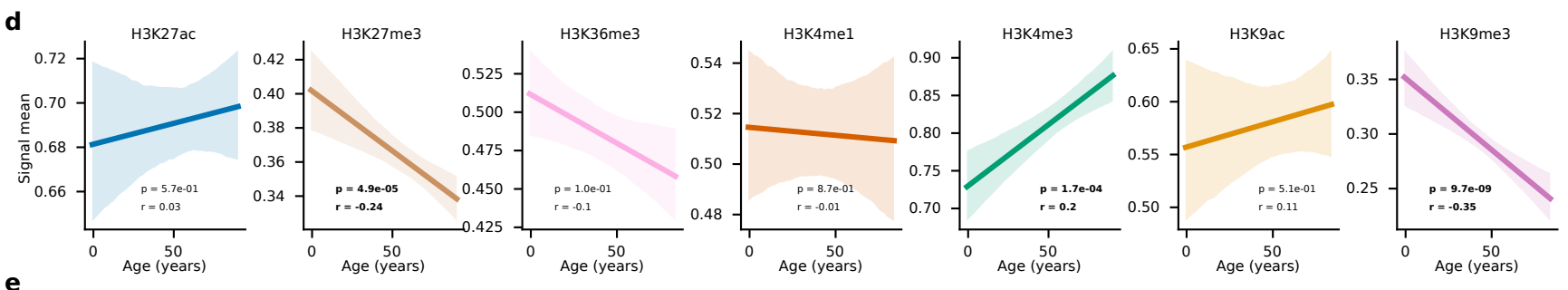
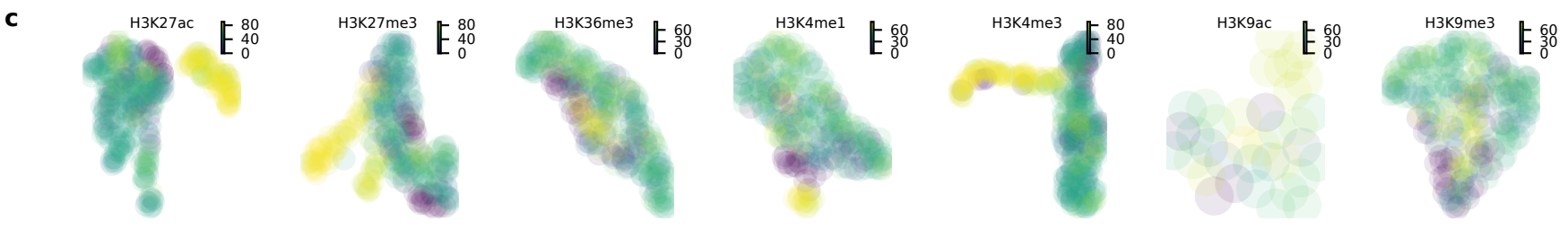
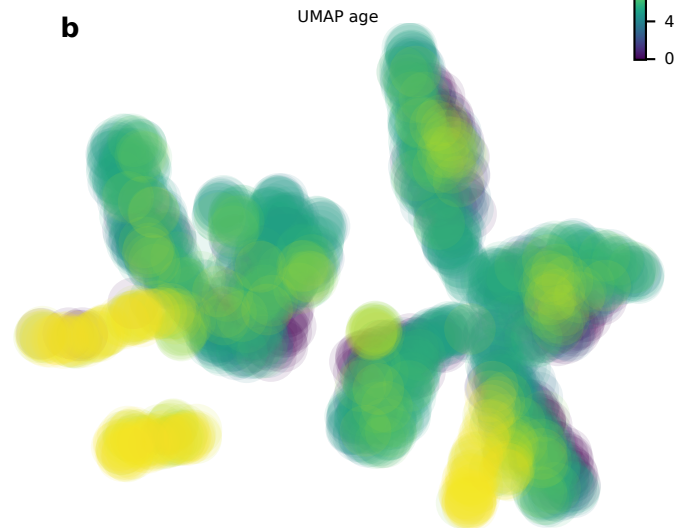
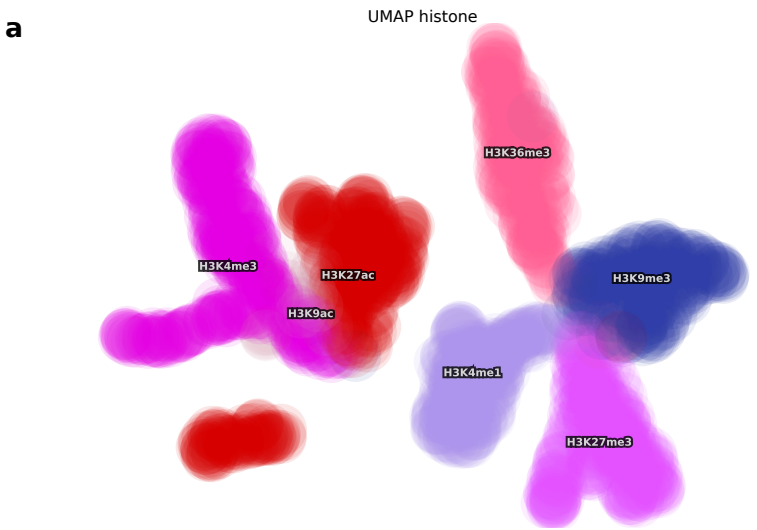
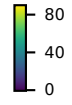
Summarize signal by gene bodies, CpG sites, etc.

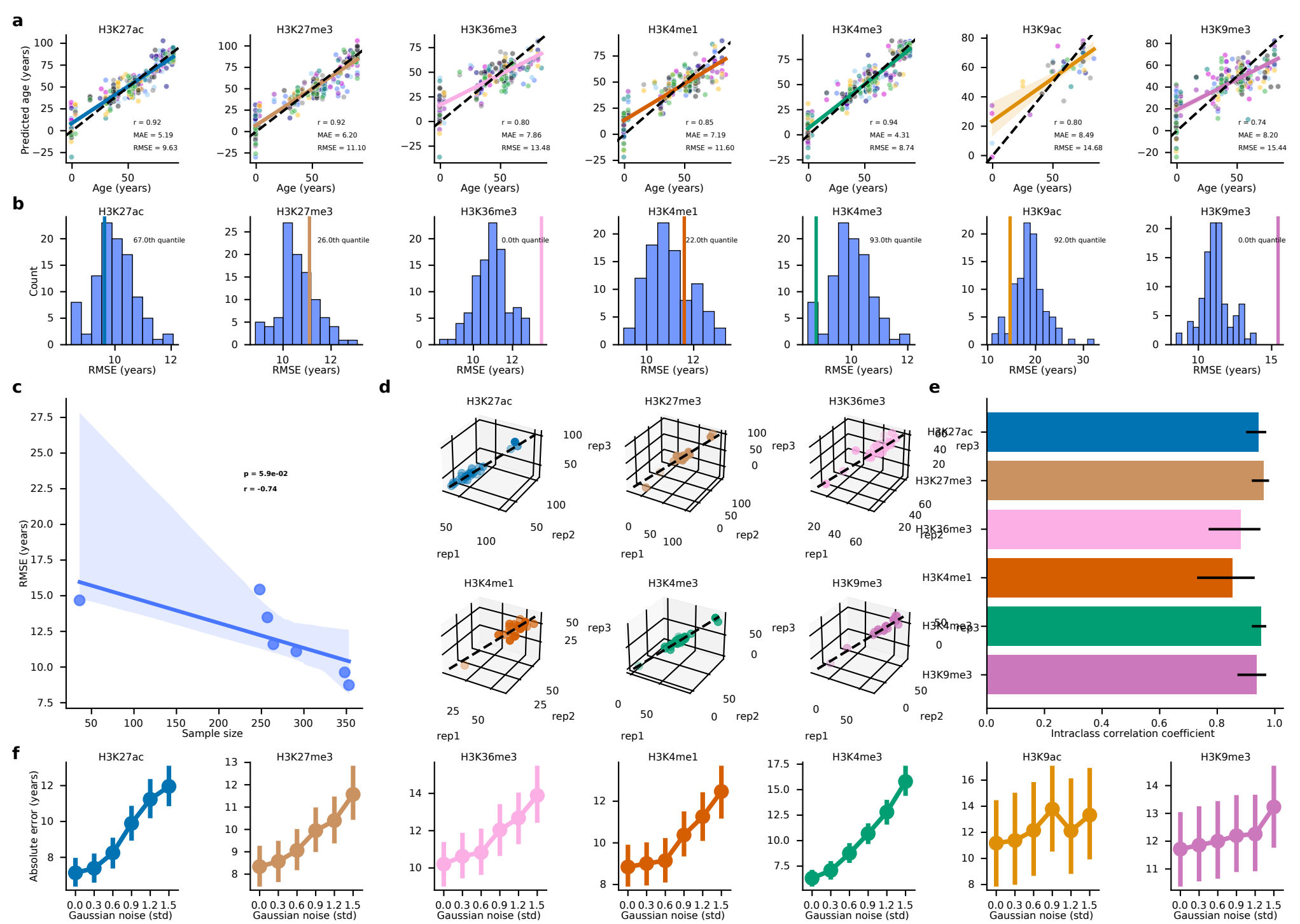
$$\text{Histone mark enrichment at locus A} = \sinh^{-1} \left( \frac{1}{\text{len}(\text{locus A})} \sum_{j=i}^{i+\text{len}(\text{locus A})} -\log(\text{p-values}_j) \right)$$

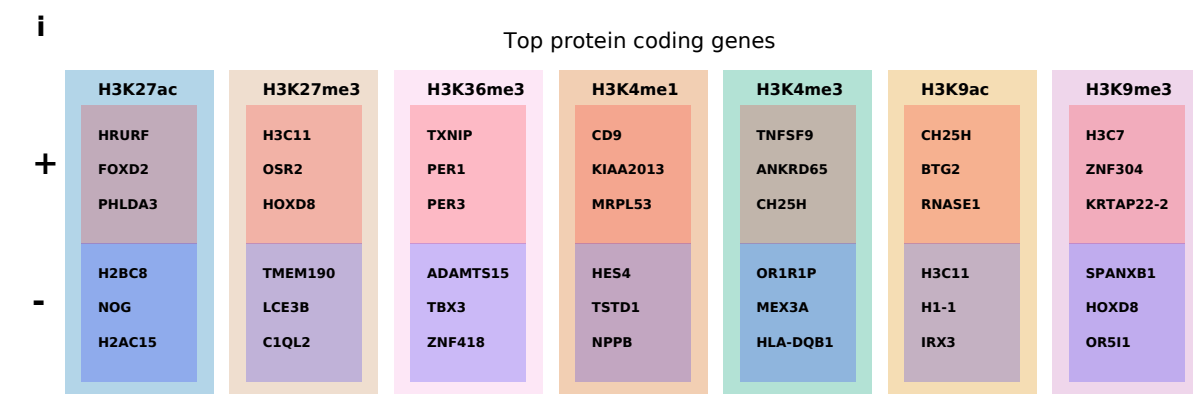
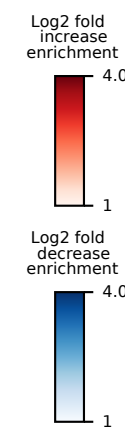
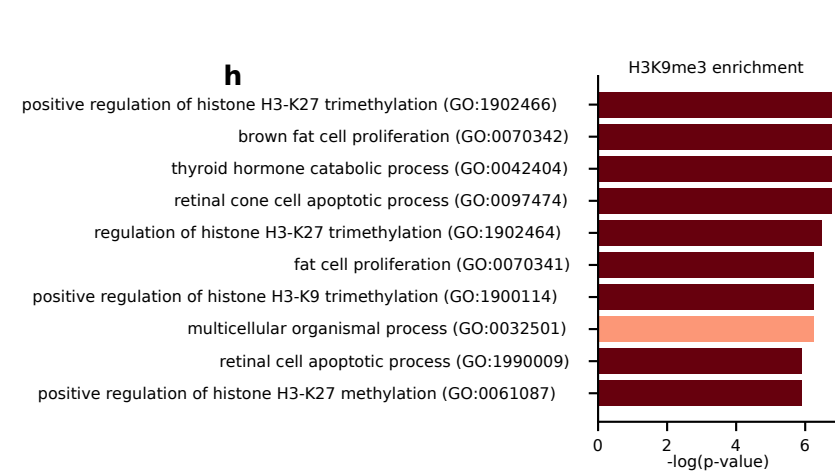
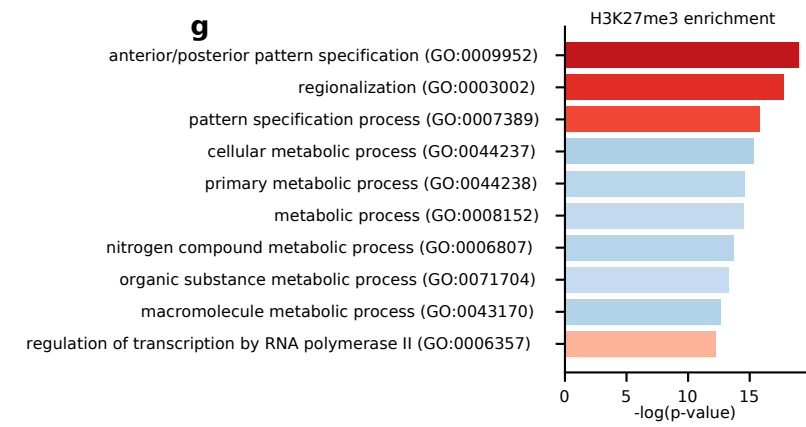
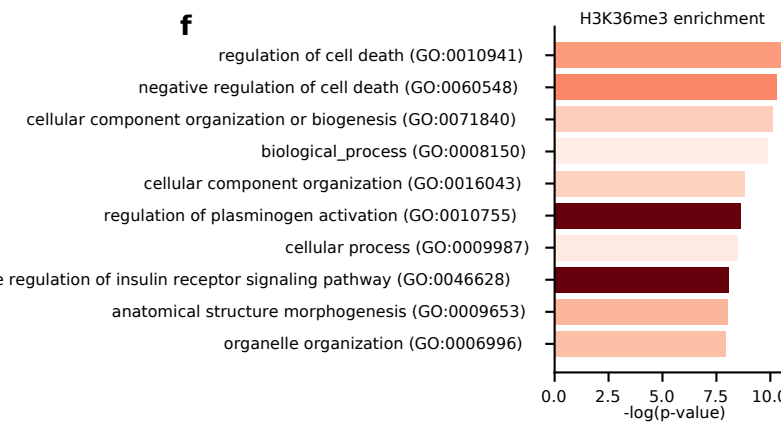
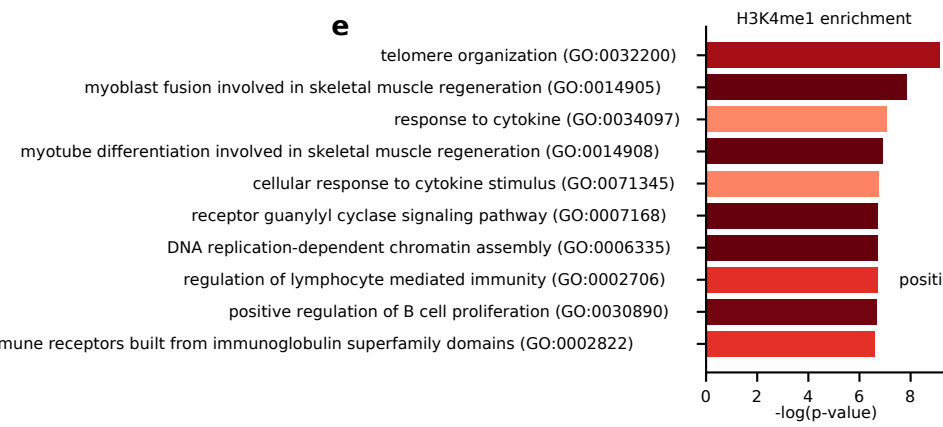
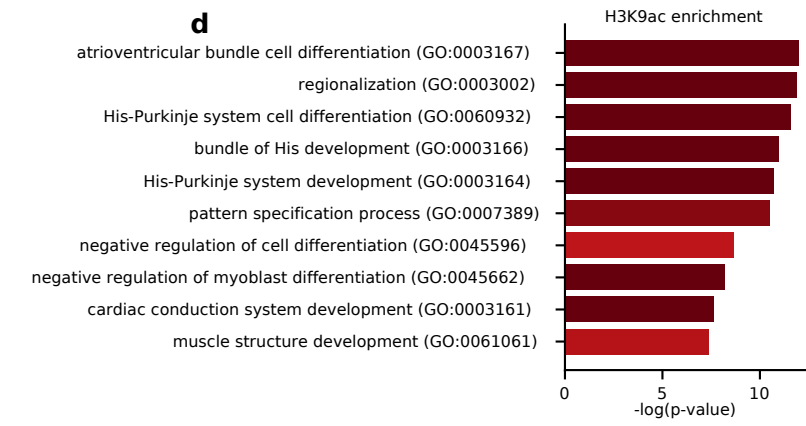
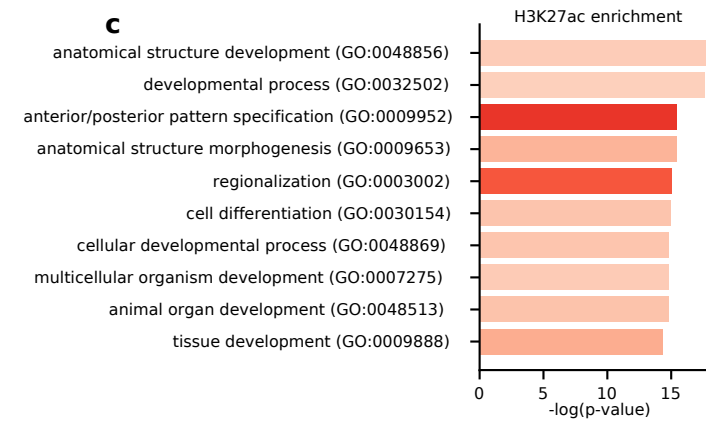
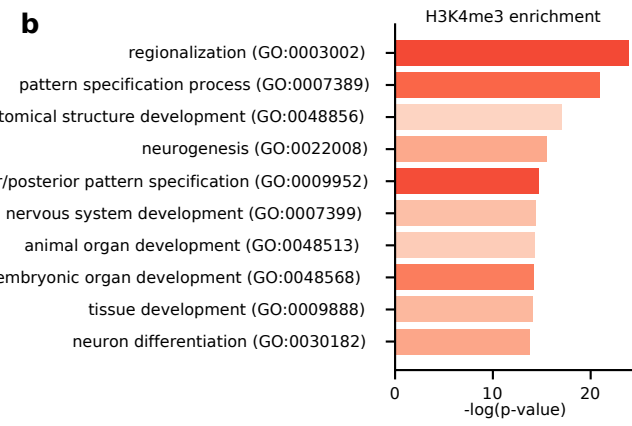
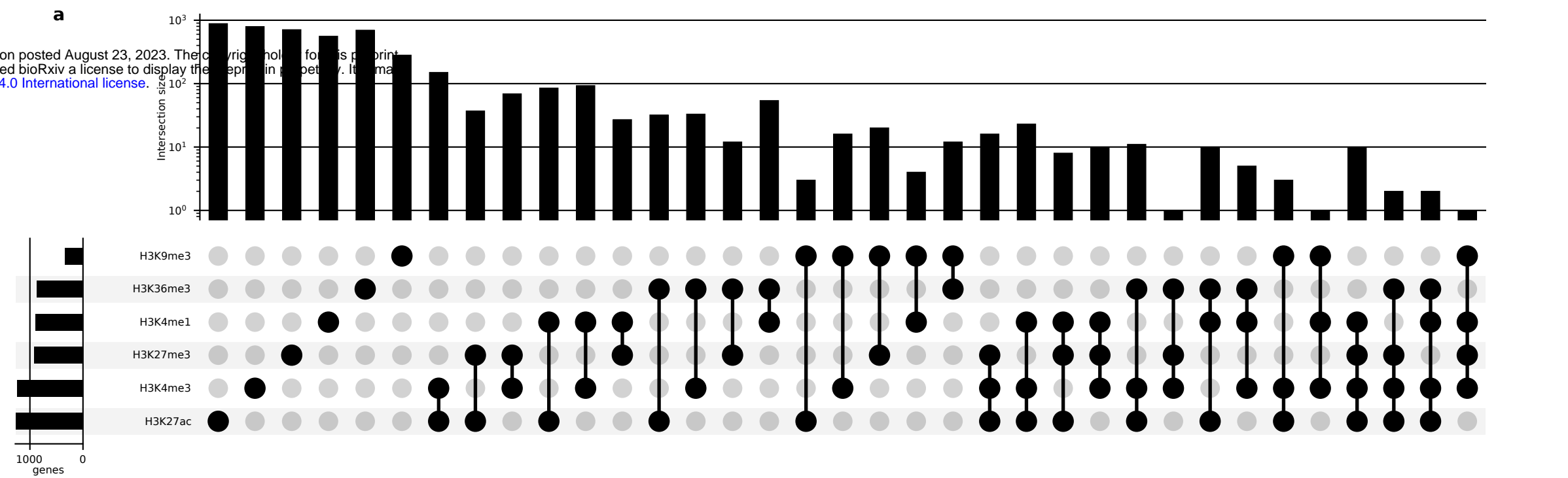
	Gene 1	Gene 2	...	Gene p
Sample 1	1.5	1.5	...	0.5
Sample 2	5	3	...	1.2
Sample 3	6.5	1.5	...	0.9
...	...	...	...	...
Sample n	0	1	...	0.1

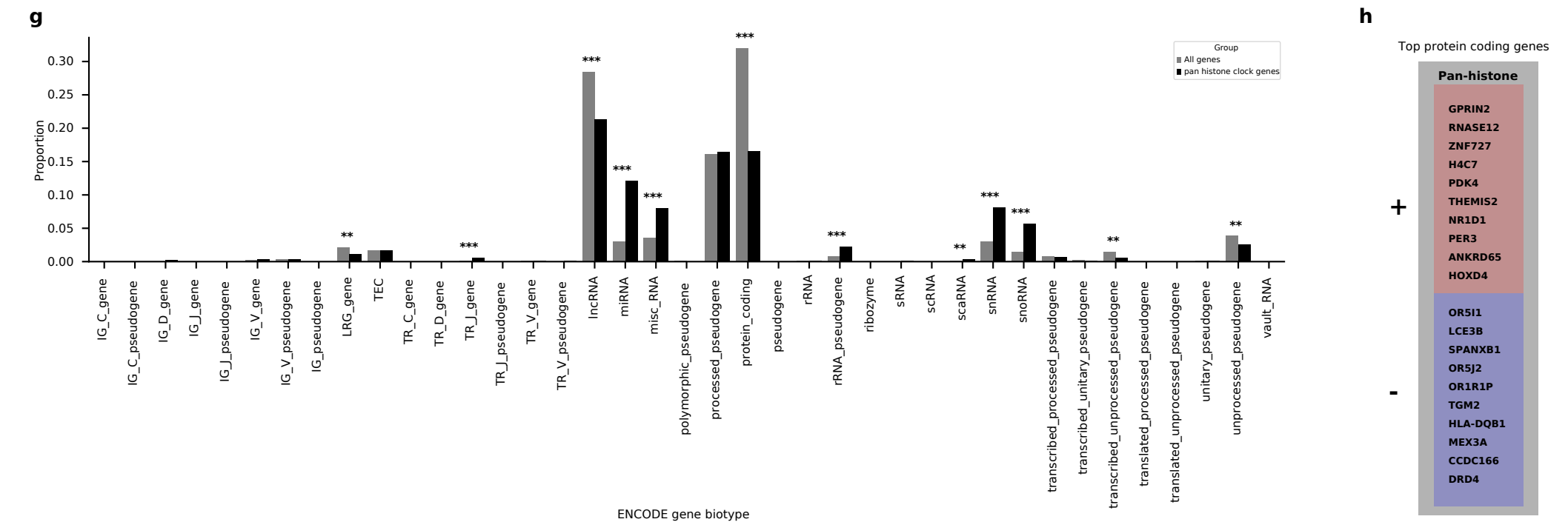
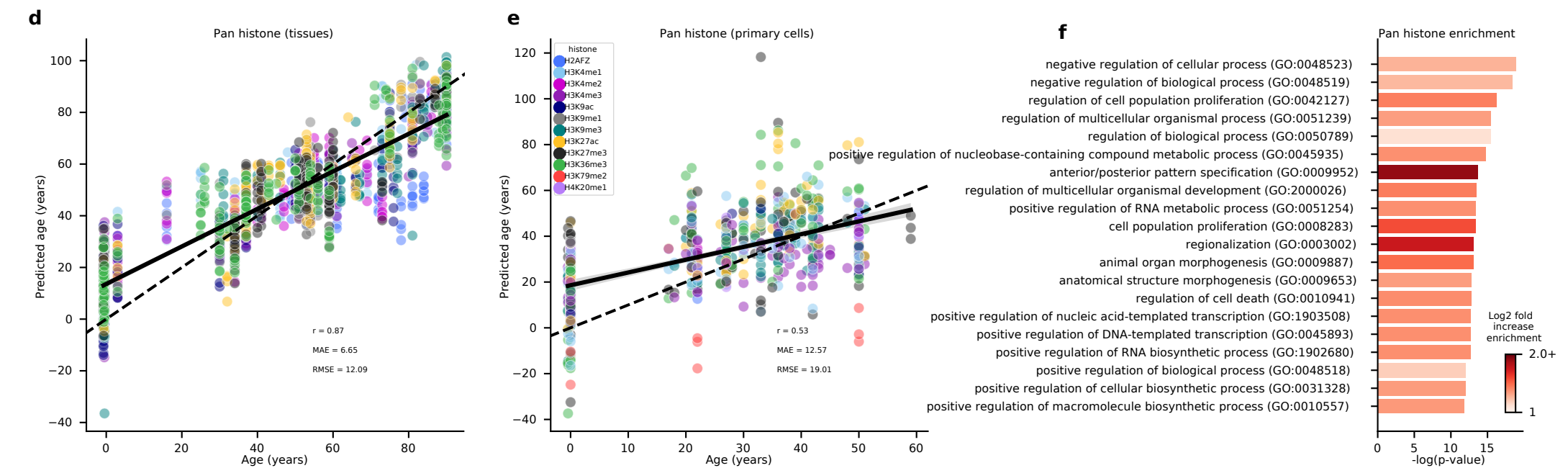
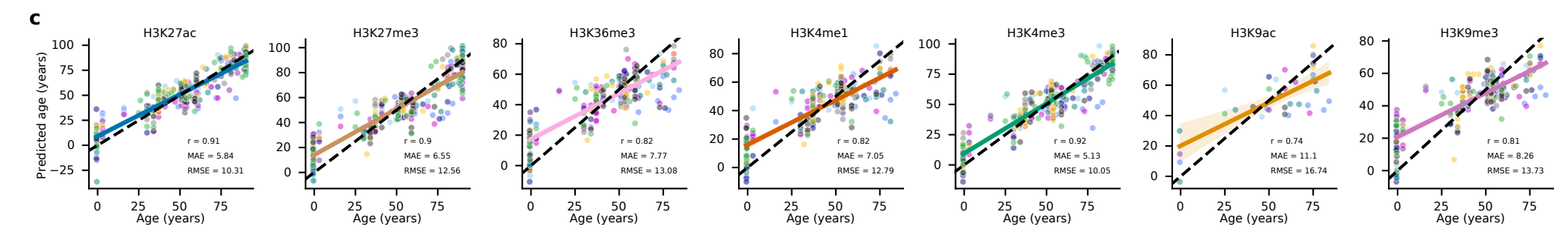
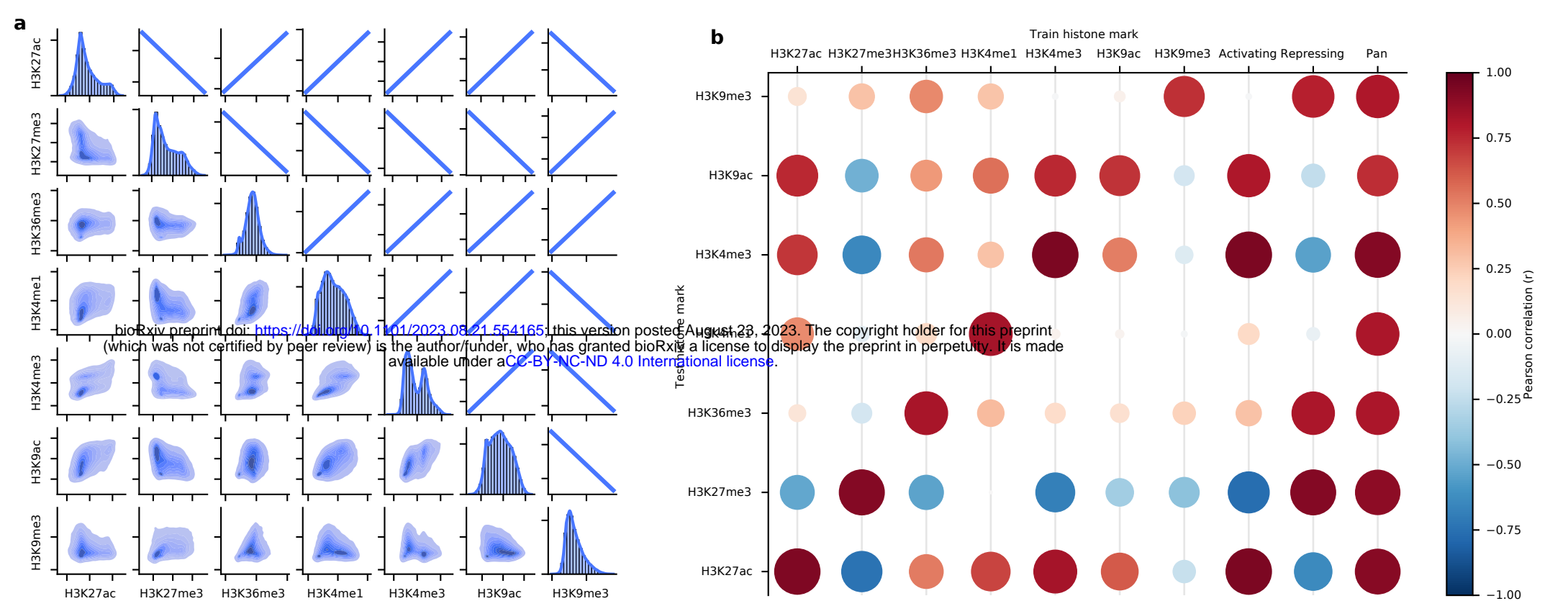
## STEP ③ Modelling



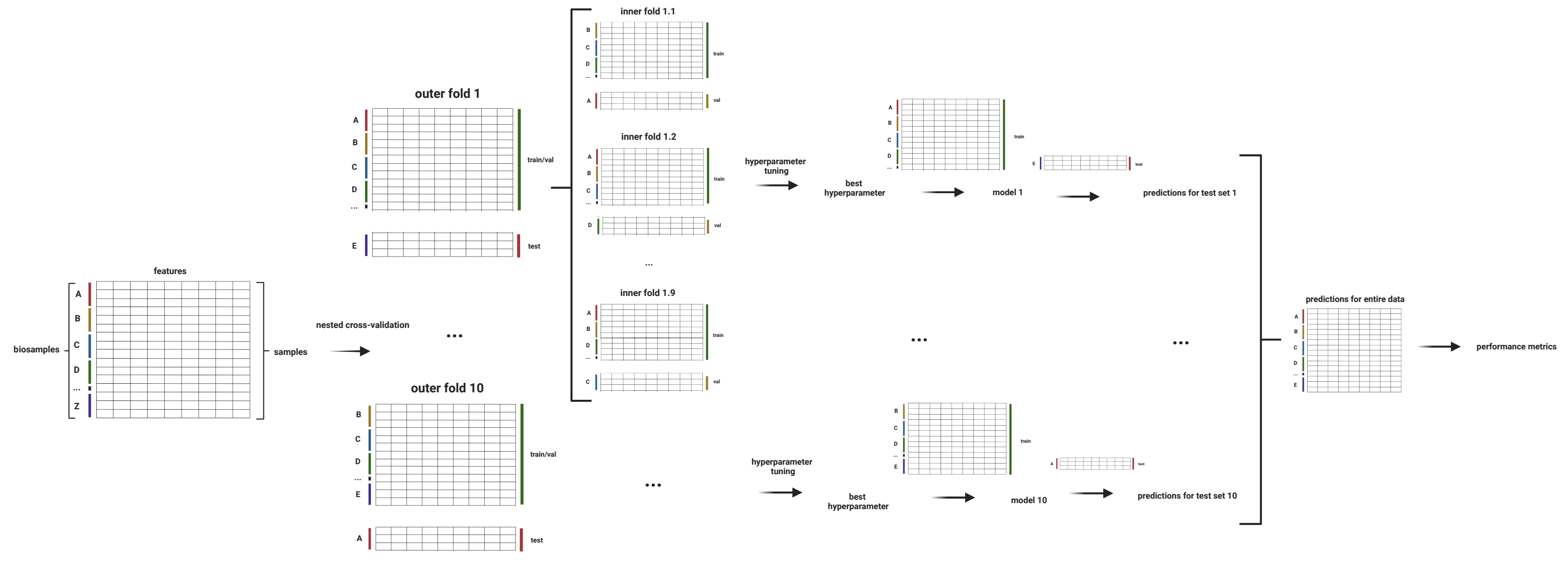




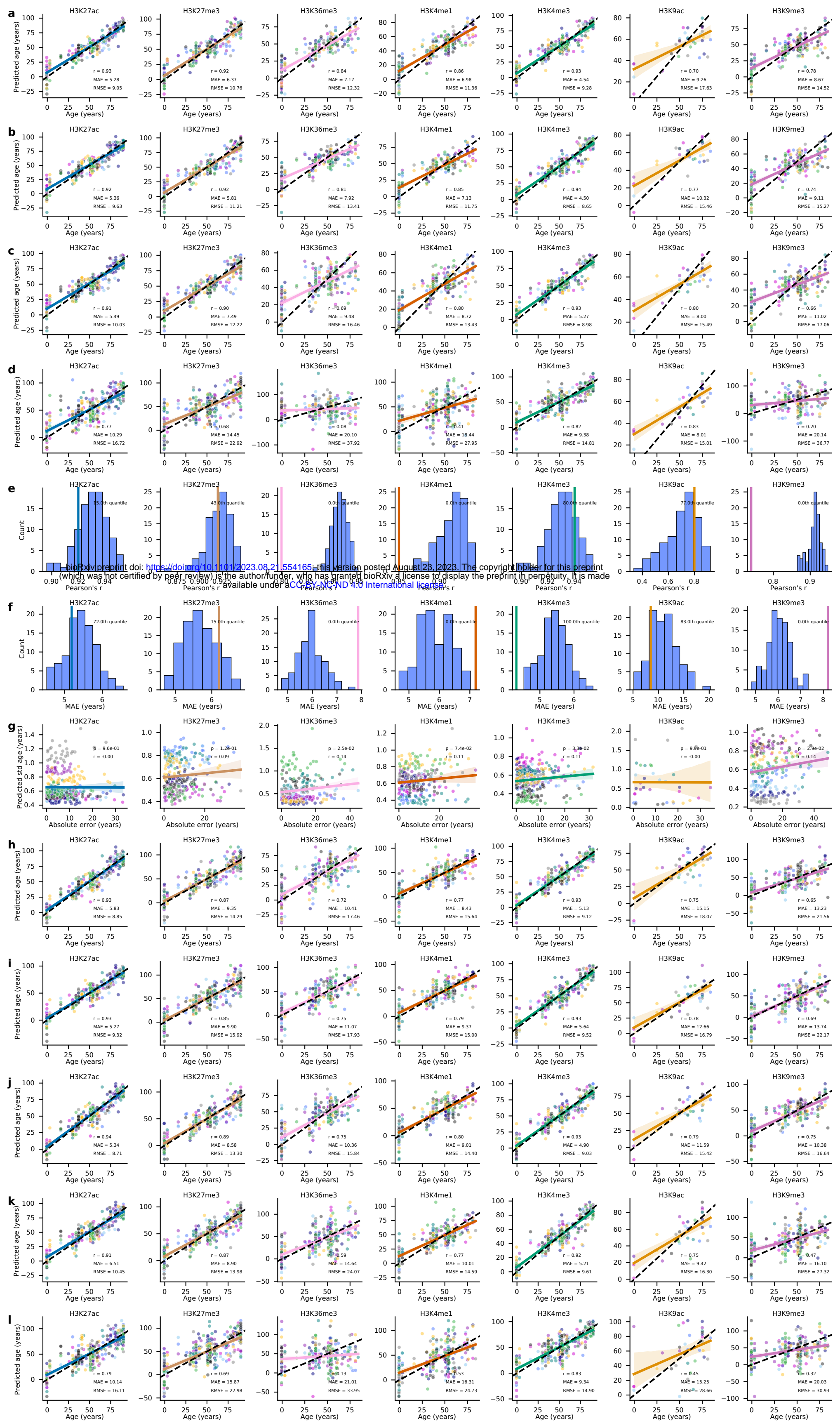




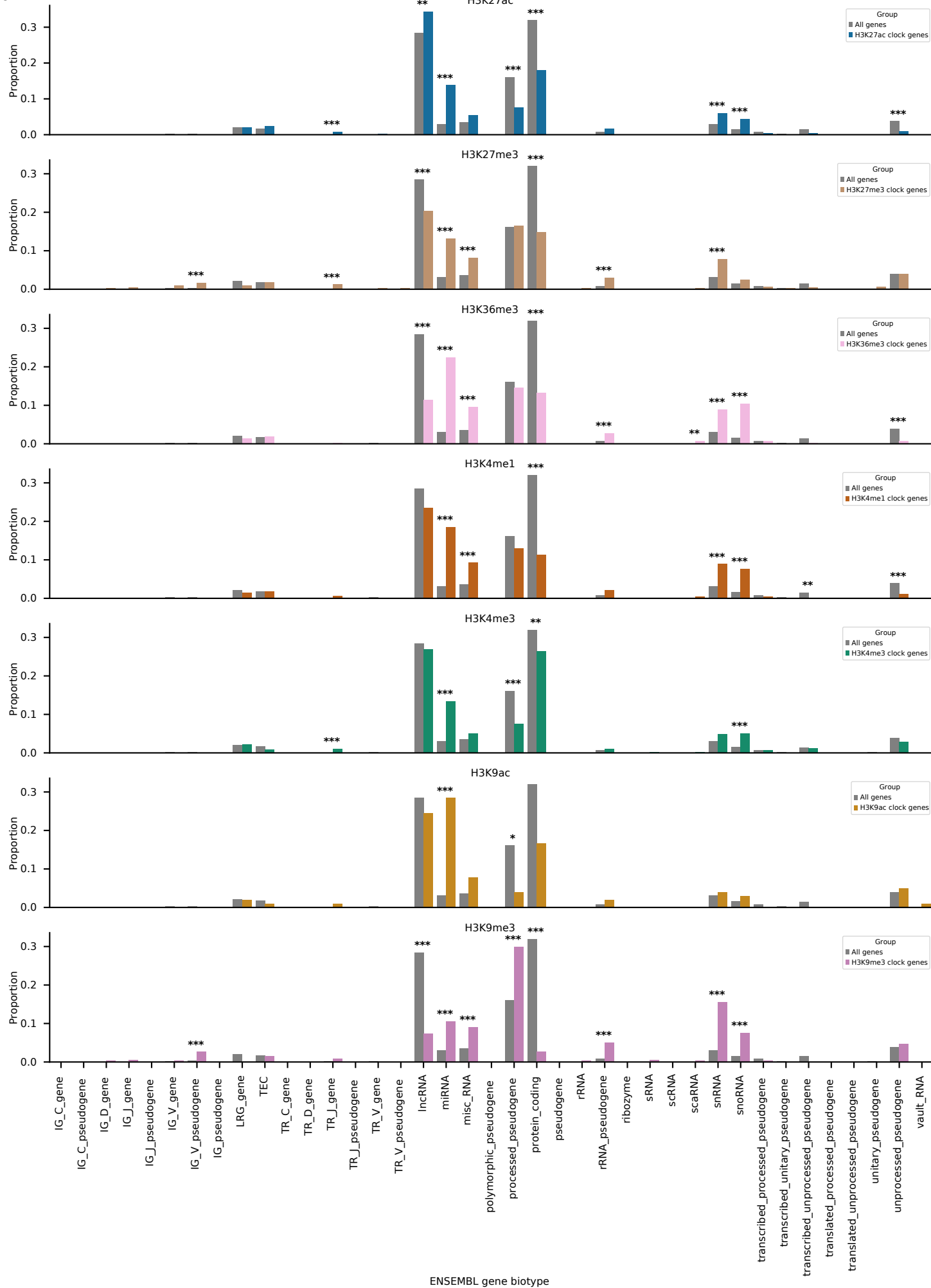








**a**



**b**

