# EpiMix: an integrative tool for epigenomic subtyping using DNA methylation

Yuanning Zheng[1], John Jun[1], Kevin Brennan[1] and Olivier Gevaert[1]

[1]Stanford Center for Biomedical Informatics Research (BMIR), Department of

Medicine & Department of Biomedical Data Science, Stanford University, Stanford,

CA 94305, USA

**Correspondence Information:**

Olivier Gevaert, Ph.D., Associate Professor

Stanford Center for Biomedical Informatics Research (BMIR)

Department of Medicine, Stanford University

1265 Welch Rd, Stanford, CA 94305-5479

Office: 650-721-2378

Email: ogevaert@stanford.edu

# Abstract

DNA methylation (DNAme) is a major epigenetic factor influencing gene expression with alterations leading to cancer, immunological, and cardiovascular diseases. Recent technological advances enable genome-wide quantification of DNAme in large human cohorts. So far, existing methods have not been evaluated to identify differential DNAme present in large and heterogeneous patient cohorts. We developed an end-to-end analytical framework named "EpiMix" for population-level analysis of DNAme and gene expression. Compared to existing methods, EpiMix showed higher sensitivity in detecting abnormal DNAme that was present in only small patient subsets. We extended the model-based analyses of EpiMix to cis-regulatory elements within protein-coding genes, distal enhancers, and genes encoding microRNAs and lncRNAs. Using cell-type specific data from two separate studies, we discovered novel epigenetic mechanisms underlying childhood food allergy and survival-associated, methylation-driven non-coding RNAs in non-small cell lung cancer.

# Main text

DNA methylation (DNAme) is one of the major epigenetic marks in humans. It is defined as the addition of a methyl ($CH_3$) group to DNA that occurs primarily at the cytosine of cytosine-guanine dinucleotide (CpG) sequence. DNAme regulates various biological processes by affecting gene expression, and aberrant DNAme plays a critical role in the development and progression of many human diseases[1–3]. Recent experimental methods based on microarrays or next-generation sequencing have enabled genome-wide quantification of DNAme at single-nucleotide resolution. Due to its quantitative and cost-effective nature, microarray-based technology has emerged as the method of choice for profiling DNAme in large human cohorts. For example, The Cancer Genome Atlas (TCGA) project has used the microarray technology to generate DNAme profiles in over 10,000 specimens representing 33 cancer types. The Gene Expression Omnibus database (GEO) and other public repositories also host a large number of DNAme datasets across cancers and other complex diseases.

Over the last decade, a number of computational approaches have been developed to identify genes that are abnormally methylated in human diseases. Some methods are tailored to the analysis of DNAme data from bisulfite sequencing[4–7], while others are designed for array-based data or can be adapted to both data platforms[8–12]. Many existing methods identify differentially methylated loci by comparing all samples from an experimental group versus samples in a control group. This type of comparison works well when the experimental population is assumed to be homogenous. However, when the study population is large, abnormal DNAme may be present in only a subset

71 of the patients, and this intra-population variation has been observed in cancers and

72 many other diseases[13–15]. In cases where abnormal DNAme occurred in only a small

73 subset of the patients, existing methods are not capable of capturing the signals of

74 differential methylation. Therefore, there is a critical need to use a statistical approach

75 to model the distribution of DNAme in large patient cohorts, and to identify the patient

76 subsets with differential DNAme profiles. This epigenetic subtyping can be essential

77 to improve personalized diagnosis, treatment and drug discovery.

78

79 Furthermore, gene expression in mammalian cells is a result of a complex process

80 coordinated by a broad range of genomic regulatory elements[16,17]. In many studies,

81 CpG sites were mapped to genes based on linear genomic proximity. This mapping

82 logic assumes that the transcriptional activity can be affected only when the genes are

83 overlapped or close to the differentially methylated sites. However, emerging evidence

84 has shown that distal enhancers, which may locate at a great linear genomic distance

85 from their target genes, play a critical role in orchestrating spatiotemporal gene

86 expression programs[18]. Abnormal DNAme at enhancers was frequently reported in

87 cancers and many other diseases[19,20]. Therefore, the analysis of enhancer

88 methylation can improve our understanding of how gene expression is regulated

89 across physiological and pathological conditions.

90

91 Existing computational tools focus on the DNAme analysis of protein-coding genes.

92 Besides protein-coding genes, non-coding RNAs, such as microRNAs (miRNAs) and

93 long non-coding RNAs (lncRNAs), play an important role in regulating cell

94 functions[21,22]. Recent studies have shown that DNAme is a major epigenetic

95  mechanism regulating non-coding RNA expression[23,24]. With existing methods, it is

96  challenging to decipher how DNAme regulates non-coding RNA expression.

97

98  Here, we present EpiMix, a comprehensive analytical framework for population-level

99  analysis of DNAme and gene expression. EpiMix utilizes a model-based

100  computational approach to identify abnormal DNAme at diverse genomic elements,

101  including cis-regulatory elements within or surrounding protein-coding genes, distal

102  enhancers, and genes encoding miRNAs and lncRNAs. In two separate studies, we

103  showed that EpiMix identified novel methylation-driven pathways in T cells from

104  childhood food allergy and methylation-driven non-coding RNAs in non-small cell lung

105  cancer patients. To improve usability, we disseminated EpiMix's algorithms in

106  Bioconductor[25], enabling end-to-end DNAme analysis. Furthermore, we developed a

107  web tool for interactive exploration and visualization of EpiMix's results

108  (https://epimix.stanford.edu). Overall, EpiMix can be used to discover novel epigenetic

109  biomarkers for disease subtypes and therapeutic targets for personalized medicine.

110

111  # Results

112

113  **Overview of EpiMix Workflow**

114

115  EpiMix is an end-to-end analytical framework for modeling DNAme at diverse genomic

116  elements and for identifications of differential DNAme associated with gene

117  expression. The EpiMix framework consisted of four functional modules: (1) data

118  downloading, (2) preprocessing, (3) DNAme modeling and (4) functional analysis

119 (**Fig.1**). To analyze DNAme at functionally diverse genomic elements, we

120 implemented four alternative analytic modes: "Regular," "Enhancer", "miRNA" and

121 "lncRNA." Both the Regular and Enhancer modes aimed to detect differential DNAme

122 associated with the expression of protein-coding genes. The Regular mode analyzed

123 DNAme sites within or immediately surrounding the genes, while the Enhancer mode

124 specifically analyzed DNAme at distal enhancers. The miRNA and lncRNA modes

125 were built for the detection of DNAme affecting the expression of miRNAs and

126 lncRNAs. After the methylation-driven genes were identified, users could perform

127 comprehensive exploratory analyses using the functional analysis module. The

128 functional analysis module was built with both in-house developed methods and

129 integrating existing computational tools to enable diverse functional analyses and

130 visualization of the differential DNAme.

Fig.1 Overview of EpiMix workflow. EpiMix includes four modules: Downloading, Preprocessing, Methylation modeling and Functional analysis. Data from public repositories (i.e., TCGA and GEO) can be automatically downloaded and preprocessed by EpiMix. Alternatively, users can input their own custom datasets. The preprocessing module includes functions for quality control, batch effect normalization, and missing value imputation. To model DNAme, EpiMix enables four alternative analytic modes: Regular, Enhancer, miRNA and lncRNA. Each mode uses a custom algorithm to analyze DNAme at a specific type of genomic element. One major output from the methylation modeling is a matrix of functional CpG-gene pairs, illustrating the differentially methylated CpGs whose DNAme states were associated with gene expression. After the differentially methylated genes have been identified, users can perform diverse analytical tasks with EpiMix's functional analysis module. This includes pathway enrichment analysis, genome-browser style visualization, gene regulatory network analysis, epigenetic biomarker discovery and identification of methylation-associated disease subtypes.

.

## Identifications of abnormal DNAme present in small sample

## subsets

145

146  To assess the sensitivity of EpiMix in identifications of differential DNAme that was

147  present in only specific patient subsets, we performed simulation experiments. We

148  used a dataset that jointly profiled DNAme data and messenger RNA abundance in

149  human naïve CD4+ T cells[26]. The dataset contains quiescent T cells and antigen-

150  activated T cells from 103 human subjects. The DNAme data were obtained from

151  Infinium MethylationEPIC array, and the messenger RNA expression data were

152  obtained from RNA-Seq.  We randomly sampled a subset of CpGs (n = 300) from the

153  quiescent group as baselines, such that the average beta values of the selected CpGs

154  ranged from 0.1 to 0.9. Then, for each CpG, we randomly selected a subset of samples

155  from the activation group and combined them with the baseline group (**Fig.2a** and

156  **Methods**), such that the final proportions of samples from the activation group in the

157  combined dataset ranged from 3% to 50%, and the mean differences in beta values

158  between the activated and the baseline samples ranged from 0.1 to 0.7. We then

159  compared the DNAme of the synthetic populations to the baseline population (**Fig.2a**).

8

160

Fig.2 **a,** Design of the simulation study. The dataset contained experimentally purified naïve CD4+ T cells from 103 human subjects. Cells from each subject were divided into half and either activated with the T-cell antigen or left resting in the media. The baseline group contained quiescent samples from all 103 subjects. The experimental group contained quiescent samples from all subjects and the antigen-activated samples from $N$ subjects, where $N$ ranged from 3 to 103. We compared the DNAme of the experimental group to the baseline group and tested whether EpiMix can detect the signals of differential methylation. **b,** Correlation between the delta beta values and the minimum detection threshold for the prevalence (left axis) and actual count (right axis) of the activated samples in the experimental group. The simulation was repeated 300 times using a different CpG site at each time, and the mean detection threshold was shown. **c,** Density plots showing the mixture models when delta beta was 0.1 and the differential methylation was present in 3%, 5% and 25% of the experimental group. **d,** Density plots showing the mixture models when delta beta was 0.3 and the differential methylation was present in 3%, 5%, and 25% of the experimental group. **e)** Number of differentially methylated CpGs detected by different methods when the differential methylation was present in from 3% to 25% of the population. For all methods, the same set of CpGs were used, and the total number of CpGs at each prevalence was 2,700.

9

175

176    We found that the sensitivity of EpiMix was determined by the magnitude of differences

177    in DNAme between the quiescent and the activated subjects. When the delta beta was

178    0.1, EpiMix detected differential DNAme that was present in 3% to 25% of the synthetic

179    population, with a mean minimum detection threshold of 11.0% (absolute sample

180    count = 13) (**Fig.2b, c**). When the delta beta was 0.2 or higher, the minimum detection

181    threshold ranged from 3% to 10%, with a mean threshold of 3.4% (absolute sample

182    count = 4) (**Fig.2b, d**). These results indicated that EpiMix was able to detect abnormal

183    DNAme that was present in only small subsets of a tested population, and the

184    sensitivity was positively correlated with the magnitude of differences in DNAme.

185
186    Next, we compared the performance of EpiMix with other existing methods in

187    identifications of differential DNAme, including Minfi[10], iEVORA[27] and RnBeads[12,28].

188    When the differential DNAme was present in 3% of the population, EpiMix detected

189    the differential methylation signals at 1,747 CpG sites, whereas the other methods did

190    not capture any differential DNAme (**Fig.2e)**. When the differential DNAme was

191    present in 5% of the population, EpiMix identified 3.1 times more differentially

192    methylated CpGs than iEVORA, and 3.6 times more CpGs than Minfi and RnBeads.

193    Minfi and RnBeads only detected CpGs with high magnitude differences in DNAme,

194    with an average delta beta of 0.6. In contrast, EpiMix detected CpGs with delta beta

195    ranging from 0.1 to 0.7, with an average threshold of 0.3. When the prevalence of

196    differential DNAme was 15% or higher, EpiMix detected similar numbers of CpGs to

197    the other three methods. These results indicated that EpiMix had higher sensitivity to

198    detect differential DNAme that was present in only small sample subsets.

199

## Modeling of DNA methylation at *cis*-regulatory elements within protein-coding genes

To test the Regular mode of EpiMix, we used the complete, real dataset from antigen-activated T cells and quiescent T cells (n = 103 subjects per group)[26]. In the activated T cells, 1,090 CpGs were differentially methylated compared to the quiescent cells. Integrative analysis with RNA-seq data showed that the differentially methylated CpGs were functionally associated with the expression of 748 protein-coding genes (**Supplementary Table 1**). Of the differentially methylated CpGs, 746 (68.4%) CpGs associated with 504 genes were hypomethylated and 327 (30.0%) CpGs associated with 238 genes were hypermethylated (**Fig.3a**). This result indicated that antigens induced a widespread loss of DNAme. Gene ontology (GO) analysis showed that the hypomethylated genes were associated with lymphocyte proliferation (e.g., *CCND2*, *CCND3*, *CDK6*, *CDK14*), T cell activation (e.g., *BCL2*, *CCL5*, *HLA-DPA1*, *HLA-DRB1*), glycoprotein biosynthesis (e.g., *AGO2*, *ALG9*, *B3GNT5*, *B4GALT5*) and cytokine receptor activity (*IL1R1*, *IL1R2*, *IL21R*, *IL23R*) (**Supplementary Table 2**). This result confirmed that EpiMix identified differential DNAme associated with T cell activation.

Many of the CpGs were differentially methylated in only a subset of the patients. For instance**,** the *Human Leukocyte Antigen DRB1* (*HLA-DRB1*) gene was hypomethylated in the antigen-activated T cells from 25% of the subjects, whereas the majority (75%) of the subjects had a normal methylation state similar to the quiescent T cells (**Fig.3b**). As expected, gene expression levels of *HLA-DRB1* were significantly increased in the hypomethylated compared to the normally methylated subjects (**Fig.3c**). Overall, the prevalence of hypomethylation ranged from 5.9% - 100%, with

11

225     a mean prevalence of 69.6% (**Fig.3d**). The prevalence of hypermethylation ranged

226     from 5.8% - 100%, with a mean prevalence of 47.3% (**Fig.3e**). These results indicated

227     that the antigen-induced response in T cells varied between different individuals.

228

229     We next investigated the genomic distribution of the differentially methylated CpGs.

230     Thirty-nine percent (39.5%) of the CpGs were located at the promoters, and 56.4%

231     were located at introns (**Supplementary Fig.1a**). Using publicly available chromatin

232     immunoprecipitation-sequencing (ChIP-seq) data of human naïve CD4+ T cells, we

233     found that the abnormal DNAme was significantly enriched at active promoters

234     marked by H3K4me3 and H3K27ac, active enhancers marked by H3K4me1, and to a

235     lesser extent, actively transcribed gene bodies marked by H3K36me3

236     (**Supplementary Fig.1b**). These results demonstrated that EpiMix was able to identify

237     aberrant DNAme at lineage-defining *cis*-regulatory elements.

238

239     To allow users to investigate the genomic locations and chromatin states associated

240     with the differentially methylated sites, EpiMix enables genome browser-style

241     visualization. We illustrated this functionality with hypomethylation in two regions of

242     the interleukin-receptor gene *IL21R* (**Fig.3f**). The first region was located at the

243     promoter, which overlapped with DNase I hypersensitivity sites and activating histone

244     modifications (i.e., H3K4me1, H3K4me3 and H3K27ac). The second region was

245     located at the three-prime untranslated region, enriched with histone modifications

246     marking for active enhancers (i.e., H3K4me1 and H3K27ac). In concordance with this

247     DNA hypomethylation, *IL21R* expression levels were significantly increased

248     (**Supplementary Table 1,** Wilcoxon rank-sum test, $P < 3.19E-08$).

249

12

**Fig.3 Identifications of differential DNAme resulting from antigen-induced T cell activation. a,** Proportions of the hypo-, hyper- and dual methylated CpGs in antigen-activated T cells. The dual methylated CpGs refer to the CpGs that were hypomethylated in some individuals, while hypermethylated in some other individuals. **b,** Mixture model of a CpG associated with the *HLA-DRB1* gene, and **c,** *HLA-DRB1* gene expression levels in different mixtures. Red indicates hypomethylation (n = 26), while blue indicates normal methylation (n = 77). Gene expression levels were compared with Wilcoxon rank-sum test. **d-e,** Density plots showing the prevalence distribution of the d) hypo- and e) hyper-methylated CpGs **f,** Genome-browser style visualization of the chromatin state, DM values, and transcript structure of the *IL21R* gene. The hypomethylated CpGs were labeled in red. The differential methylation (DM) value represents the mean difference in beta values between the hypomethylated subjects versus the normally methylated subjects. DM = 0: normal methylation; DM < 0: hypomethylation.

13

261 **Identification of functional DNA methylation at distal enhancers in food allergy**

262

263 To demonstrate the Enhancer mode of EpiMix, we used the same CD4+ T cell

264 dataset[26]. In this dataset, 82 human subjects were diagnosed with food allergy and 21

265 subjects were non-allergic controls. The differential response of T cells to antigen-

266 induced activation between different individuals may be associated with the allergic

267 status. We then characterized allergy-associated changes in DNAme by comparing

268 antigen-activated T cells from the allergic patients to those from the non-allergic

269 controls. Using a permutation approach (**Supplementary Fig.2** and **Methods**), we

270 identified 107 differentially methylated enhancers that were functionally linked to the

271 expression of 119 genes. The number of target genes of each enhancer ranged from

272 1 to 3, resulting in 131 significant enhancer-gene pairs (**Supplementary Table 3**).

273 This result is consistent with the previous studies showing that enhancers typically

274 loop to and are associated with the activation of 1 to 3 promoters[29,30]. Of the functional

275 enhancers, 21/107 (19.6%) enhancers associated with 24 genes were

276 hypomethylated, 82/107 (76.7%) enhancers associated with 92 genes were

277 hypermethylated (**Fig.4a**). This result indicated that there was a global gain of DNAme

278 at enhancers in food allergy.

279

280 The genomic distance between enhancers and their target genes ranged from 4.5 kb

281 to 1.7 Mb, with a median distance of 148 kb (**Fig.4b**). In a previous study, Jin et al.

282 used high-throughput chromosome conformation capture (Hi-C) assay to investigate

283 promoter-enhancer interactions and demonstrated that approximately 25% of the

284 enhancer-promoter pairs are within a 50 kb range and approximately 57% spans 100

285 kb or greater genomic distance, with a median distance of 124 kb[31]. Another study by

14

286  Rao et al. showed that the distance between enhancers and promoters spans from 40

287  kb to 3 MB, with a median distance of 185 kb[32]. Our data agree with these

288  experimentally generated results. To further characterize the enhancer-gene linkage,

289  we investigated how often did the functional enhancers associate with the nearest

290  gene promoter. We ranked the 20 adjacent genes of each enhancer by their genomic

291  distance to the enhancer. **Fig.4c** showed that only 6.1% of the times did the enhancer

292  associate with the nearest promoter, whereas the majority of the enhancers skipped

293  one or more intervening genes to associate with promoters farther away. In line with



294  **Fig. 4 Identifications of differentially methylated enhancers associated with food allergy. a,** Proportions of
295  the hypo-, hyper- and dual methylated enhancers in children with food allergy. **b,** Distribution of the linear genomic
296  distance between enhancers and their gene targets. **c,** For each functional enhancer, the 20 adjacent genes were
297  ranked by genomic distance. Bars show the proportions of the functionally linked genes in each rank. **d**, Mixture
298  model of the *LDLR* gene (top panel) and *LDLR* gene expression levels in different mixtures (bottom panel). Red
299  indicates normal methylation (n = 72), while blue indicates hypermethylation (n = 10). Gene expression levels were
300  compared by Wilcoxon rank-sum test. **e,** Integrative visualization of the chromatin states and the adjacent genes
301  of the hypermethylated enhancer shown in panel d. The genes in the functional CpG-gene pairs are shown in red,
302  while the others are shown in black. **f,** Enriched TF motifs and odds ratios for the differentially methylated enhancers.
303  To find significantly enriched motifs, we used all the distal CpGs as the background and the functional enhancers
304  as the targets.

15

305　this result, a previous study using the chromosome 5C assay showed that only ~7%

306　of the time did the distal elements loop to the promoter of the nearest gene, whereas

307　the majority of enhancers bypass the nearest promoter and loop to promoters farther

308　away[33]. These results confirmed that EpiMix identified true distal *cis*-regulatory events.

309

310　The genes linked to the differentially methylated enhancers were related to the lipid

311　metabolism (*LDLR*, *CAT*, *LPIN2*, *SREBF1, PIK3C2B*) and T cell activation (*CASP3*,

312　*MALT*, *PRKCZ*, *SMAD3*). **Fig.4d** showed that the enhancer linked to the *LDLR* gene

313　was hypermethylated in 12.2% of the allergic patients, and the gene expression of

314　*LDLR* was significantly decreased in the hypermethylated patients. Integrative

315　visualization (**Fig.4e**) showed that the hypermethylated enhancer overlapped with the

316　Dnase I hypersensitivity site and was enriched with histone modifications marking for

317　active enhancers, including H3K4me1 and H3K27ac, and to a lesser extent, H3K4me3

318　and H3K9ac. The *LDLR* gene encodes a low-density lipoprotein receptor that

319　transports cholesterol from the blood into the cell, which plays a critical role in

320　regulating T cell lipid metabolism[34]. Our results suggested that T cells from a small

321　subset of the allergic patients may have an abnormal lipid metabolic profile due to

322　enhancer hypermethylation.

323

324　Enhancers are enriched for sequences bound by site-specific transcription factors

325　(TFs). Hypermethylation of enhancers suppresses gene transcription by decreasing

326　the binding affinity of TFs[35,36]. We then carried out motif enrichment analysis of the

327　differentially methylated enhancers. We identified significant enrichment of binding

328　sites for Jun-related factors (JUN, JUND), Fos-related factors (FOS, FOSL1, FOSL2,

329　FOSB), BATF-related factors (BATF, BATF3), and Interferon-regulatory factors (IRF2,

330     IRF5, IRF7) (**Fig.4f** and **Supplementary Table 4**). These results agree with the

331     evidence showing that Jun-related factors, BATF-related factors and Interferon-

332     regulatory factors play a critical role in regulating the immune gene activation in T cells,

333     and dysregulation of their activity causes aberrant immune response[37,38]. Our results

334     demonstrated that the abnormal DNAme at enhancers affected the target gene

335     response of these TFs and increased the subsequent risk for developing food allergy.

336

337     **Identification of methylation-driven miRNAs in human lung cancer**

338

339     Similar to protein-coding genes, miRNA-coding genes are transcriptionally regulated

340     by DNAme[39,40]. To demonstrate the miRNA mode of EpiMix, we used a lung

341     adenocarcinoma dataset containing DNAme and miRNA expression profiles of 457

342     tumors and 32 adjacent normal tissues[41]. The DNAme data were acquired from the

343     HM450 array, and the gene expression data were obtained from high-throughput

344     microRNA sequencing (miRNA-Seq).

345

346     Both tumors and normal tissues from the lung are composed of multiple cell types,

347     majorly including epithelial cells, fibroblasts, hematopoietic cells and endothelial cells.

348     Studies have shown that DNAme profiles are cell-type specific[42,43]. When using data

349     collected at the tissue ("bulk") level for DNAme analysis, the differential DNAme may

350     result from variations in cell-type proportions between different individuals. To resolve

351     the confounding effects from intra-tumoral heterogeneity, we used previously

352     validated computational methods to decompose tissue compositions and to infer cell-

353     type-specific methylomes and transcriptomes (**Supplementary Fig. 3** and

354     **Methods**)[44,45]. We then applied EpiMix to the deconvoluted data of each individual cell

355    type. In epithelial cells, we identified 272 differentially methylated CpGs functionally

356    associated with the expression of 92 miRNA genes (**Fig.5a** and **Supplementary**

357    **Table 5**). In fibroblasts, we found 12 hypomethylated CpGs functionally associated

358    with the expression of 3 miRNA genes (**Supplementary Fig. 4a-b**). Although we

359    discovered 9 differentially methylated CpGs in hematopoietic cells and 6 CpGs in

360    endothelial cells, none of the differential DNAme were functionally correlated with

361    gene expression. We further compared the differentially methylated gene lists

362    identified using data from bulk tissues versus the ones using individual cell types. Over

363    80% of the differentially methylated genes identified in epithelial cells could also be

364    identified using data from bulk tissues (**Supplementary Fig. 4a-b**). These results

365    demonstrated that, although tumors are composed of multiple cell types, the majority

366    of differential methylation events occurred in epithelial cells.

367

368    We next focused our analysis on the deconvoluted data of epithelial cells. Of the 272

369    differentially methylated CpGs, 138 (50.8%) CpGs associated with 66 genes were

370    hypomethylated and 55 (20.2%) CpGs associated with 37 genes were

371    hypermethylated. Sixty-five percent (63.6%) of the functional CpGs were located at

372    the promoters, and this proportion was significantly higher than randomly selected

373    CpGs (**Supplementary Fig.1c**, Fisher's exact test, $P = 0.003$). Using publicly available

374    ChIP-seq data of lung, we further determined that the differentially methylated regions

375    were enriched with histone modifications (i.e., H3K27ac, H3K4me1 and H3K4me3)

376    marking for actively transcribed promoters and enhancers (**Supplementary Fig.1d**).

377    The prevalence of hypomethylation ranged from 1.1% to 66.7%, with a mean

378    prevalence of 18.0% (**Fig. 5b**). Similarly, the prevalence of hypermethylation ranged

379    from 2.6% to 83.7%, with a mean prevalence of 24.9% (**Fig. 5c**). These results

18

380    indicated that the majority of differential DNAme associated with miRNA genes

381    occurred in less than 25% of the patient population.



382    **Fig. 5 Identifications of differentially methylated miRNA-coding genes in human lung cancers. a,** Proportions
383    of the hypo-, hyper- and dual methylated CpGs of miRNAs in lung cancer. **b-c,** Density plots showing the
384    prevalence distribution of the differentially methylated miRNAs in lung cancers (n = 457), **(b)** prevalence of
385    hypomethylation and **(c)** prevalence of hypermethylation. **d,** Mixture model of the *MIR30A* gene (left panel) and
386    Kaplan-Meier survival curves of patients in different mixtures (right panel). Red indicates normal methylation and
387    blue indicates hypermethylation. Gene expression levels were compared by Wilcoxon rank-sum test. **e,** Mixture
388    model of the *MIR1292* gene (left panel) and Kaplan-Meier survival curves of patients in different mixtures (right
389    panel). Red indicates hypomethylation and blue indicates normal methylation. **f-g-h,** Network visualization of (**f**)
390    the gene targets of miR-34a, (**g**) differentially methylated miRNAs related to the cell cycle pathway, and (**h**) focal
391    adhesion pathway. Blue squares: miRNAs, green circles: protein-coding genes targeted by miRNAs.

392

393    MicroRNAs play an important role in regulating cell proliferation, invasion and cancer

394    metastasis[46,47]. We next investigated whether the DNAme of miRNAs were associated

395    with patient survival. Of the 92 methylation-driven miRNAs, we identified 22 miRNAs

396    whose methylation states were significantly correlated with patient survival

19

397 (**Supplementary table 6,** log-rank test, *P* < 0.05). Half (11/22, 50%) of the survival-

398 associated miRNAs were hypomethylated and the others (11/22, 50%) were

399 hypermethylated. Some of the miRNAs, such as *MIR29C*[48], *MIR30A*[49], *MIR34A*[50] and

400 *MIR148A*[51], were known to be associated with lung cancer survival. For instance,

401 MIR30A, a tumor suppressor miRNA[49], was hypermethylated in 8.6% of the patients,

402 and the hypermethylated patients showed a significantly worse survival than the

403 normally methylated patients (**Fig.5d,** Hazard Ratio = 1.50, *P* = 0.001). In addition,

404 EpiMix identified many new survival-associated miRNAs. For instance, *MIR1292* was

405 hypomethylated in 8.6% of the patients, and the hypomethylated patients showed

406 significantly worse survival (**Fig.5e,** Hazard Ratio = 1.39, *P* = 0.0008). These results

407 demonstrated that EpiMix was able to identify survival-associated miRNAs that were

408 differentially methylated in only small subsets of the patients, and this feature can be

409 used to discover novel epigenetic biomarkers for prognosis.

410

411 To gain systematic insight into the biological functions of the methylation-driven

412 miRNAs, we queried miRTarBase[52] to obtain experimental validated target genes of

413 the miRNAs. We then performed pathway analyses of the target gene list. The

414 differentially methylated miRNAs were related to Wnt signaling pathway, cell cycle,

415 p53 signaling, focal adhesion and apoptosis (**Fig.5f-h** and **Supplementary Table 7**).

416 These results provided mechanistic insights into how abnormal DNAme of miRNAs

417 was involved in the development and progression of lung cancer. The data also

418 suggested that targeting miRNA expression can be a therapeutic strategy to inhibit

419 tumor progression and to improve patient survival.

420

421 **Identification of methylation-driven lncRNAs in human lung cancer**

422

423    To demonstrate the lncRNA mode of EpiMix, we used the same lung adenocarcinoma

424    dataset[41], and we aimed to identify differentially methylated lncRNA genes in tumors

425    compared to normal tissues. Compared to protein-coding genes, lncRNAs are shorter,

426    lower-expressed, less evolutionarily conserved, and expressed in a more tissue-

427    specific manner[53]. To precisely quantify lncRNA expression from RNA-Seq, we used

428    our previously developed pipeline[54]. With this pipeline, we combined the transcriptome

429    annotations from GENCODE and NONCODE[55]. Raw sequencing reads were aligned

430    to the combined transcriptome reference and quantified using the Kallisto-Sleuth

431    algorithm[56,57]. Using this pipeline, we were able to detect the expression of 2,475

432    lncRNAs in both tumors and normal tissues. This number was three times higher

433    compared to the lncRNAs detected by the traditional STAR-HTSeq pipeline. We then

434    computationally deconvoluted bulk DNAme data and lncRNA expression data to cell-

435    type-specific data (**Supplementary Fig. 3**). Since over 95% of the functional

436    differential DNAme was found in epithelial cells (**Supplementary Fig. 4c-d**), we next

437    focused our analysis on epithelial cells.

438

439    EpiMix identified 397 CpGs functionally associated with the expression of 132

440    lncRNAs in epithelial cells (**Fig.6a** and **Supplementary Table 8**). Of these CpGs, 146

441    (36.8%) CpGs associated with 69 genes were hypomethylated and 187 (47.1%) CpGs

442    associated with 73 genes were hypermethylated. Seventy-two percent (72.0%) of the

443    functional CpGs were located at the promoters, and this proportion was significantly

444    higher than randomly selected CpGs (**Supplementary Fig.1e,** Fisher's exact test, *P*

445    < 0.0001). The differentially methylated regions were enriched with histone

446    modifications marking for actively transcribed promoters and enhancers, including

447    H3K27ac, H3K4me1 and H3K4me3 (**Supplementary Fig.1f**).

448

449    The majority of differential methylation was identified in less 50% of the patients. The

450    prevalence for hypomethylation ranged from 1.8% to 53.0%, with a mean value of 19.8%

451    (**Fig.6b**). Similarly, the prevalence for hypermethylation ranged from 0.6% to 68.2%,

452    with a mean value of 18.9% (**Fig.6c**). For instance, one of the hypermethylated

453    lncRNAs was *LINC00881*. *LINC00881* was hypermethylated at CG11931463 in 15.7%

454    of the patients and CG00673344 in 7.9% of the patients (**Fig.6d**). Both CpGs were

455    located within the promoter (**Fig.6e**). Integrative analysis with clinical data showed that

456    *LINC00881* hypermethylation was associated with significantly worse patient survival

457    (**Figs.6f,** log-rank test, $P < 0.001$). These data demonstrated that many lncRNAs were

458    differentially methylated in only a subset of the lung cancer patients. In addition, EpiMix

459    was able to identify survival-associated lncRNAs that were differentially methylated in

460    small patient subsets.

461

462    One of the major outputs from EpiMix is a differential methylation or "DM" value matrix,

463    which reflects the homogeneous subpopulations of samples with a particular

464    methylation state (**Fig.6g**). An application of the DM value matrix is to identify DNAme-

465    associated subtypes, where patients are clustered into robust and homogenous

466    groups based on their differential DNAme profiles. Using unsupervised consensus

467    clustering, we discovered five DNAme subtypes (S1–S5) (**Fig.6h**). S5 contained a

468    significantly higher proportion of females (89/133 = 66.9%) compared to S1 (54/120 =

469    45.0%), S2 (36/74 = 48.6%) and S4 (16/50 = 32.0%) (**Fig.6i,** Fisher's exact test, $P <$

470    0.01). In addition, patients from S5 had significantly better survival than patients of S2

22

471   (**Fig.6j,** log-rank test, $P$ = 0.007). We benchmarked the clustering results from using

472   the DM value matrix versus using the raw DNAme data (beta values) of the

473   differentially methylated CpGs. The patient subsets identified using raw DNAme data

474   had low cluster consensus (**Supplementary Fig.5**), and no significant association was

475   found between patient subsets and survival outcome. These results demonstrated that

476   the DNAme subtypes discovered by EpiMix had prognostic values.

477

478   To investigate the biological functions of the differentially methylated lncRNAs, we

479   utilized ncFANs, a functional annotation tool for lncRNAs[58]. We identified 4,552

480   protein-coding genes functionally associated with 76 lncRNAs. GO analysis showed

481   that the protein-coding genes were primarily associated with DNA replication, cell

482   cycle and regulation of cell activation (**Fig.6k** and **Supplementary Table 9**). These

483   results indicated how differential methylation of lncRNAs were involved in the

484   regulation of lung cancer development and progression.

485

486

487
488

**Fig. 6 Identifications of differentially methylated lncRNA-coding genes in human lung cancers. a,** Proportions of the hypo-, hyper- and dual methylated CpGs of lncRNA genes in epithelial cells from lung cancers compared to normal tissues. **b-c,** Density plot showing the prevalence distribution of the **(b)** hypo- and **(c)** hyper-methylated lncRNAs in the lung cancer cohort (n = 457). **d,** Mixture models of the *LINC00881* gene at two different CpG sites. Red indicates normal methylation and blue indicates hypermethylation. **e,** Integrative visualization of the transcript structure, DM values and chromatin state associated with the *LINC00881* gene. DM = 0: normal methylation; DM > 0: hypermethylation. **f,** Kaplan-Meier survival curves of patients in the normally methylated and the hypermethylated mixtures. Red indicates normal methylation and blue indicates hypermethylation. **g,** Schematic representation of the DM value matrix. The rows correspond to CpG sites, and the columns correspond to patients. DM values represent the mean differences in DNAme levels between patients in each mixture component identified in the experimental group compared to the control group. At each CpG site, patients in the same mixture component have the same DM values. DM < 0: hypomethylation, DM = 0: normal methylation, DM > 0: hypermethylation. **h,** Consensus matrix showing patient clusters based on the DM values of lncRNAs. **i,** Proportions of male and female patients in different patient clusters (n1 = 120, n2 = 74, n3 = 72, n4 = 50, n5 = 133). **j,** Kaplan-Meier survival curves of patients in different patient clusters. **k,** Top 20 enriched GO terms of the methylation-driven lncRNAs in lung cancer. DM: differential methylation.

# Discussion

509    In this study, we present EpiMix, a comprehensive analytic framework for population-

510    level analysis of DNAme and gene expression. We packaged the EpiMix algorithms

511    in R, enabling end-to-end DNAme analysis. To enhance the user experience, we also

512    implemented a web-based application (https://epimix.stanford.edu) for interactive

513    exploration and visualization of EpiMix's results (**Fig.7**). EpiMix contains diverse

514    functionalities, including automated data downloading, preprocessing, methylation

515    modeling and functional analysis. The seamless connection of EpiMix to data from the

516    TCGA program and the GEO database enables DNAme analysis on a broad range of

517    diseases. Here, we showed that EpiMix identified novel methylation-driven pathways

518    in food allergy and lung cancer. However, EpiMix is not limited to these disease areas

519    and can be easily applied to any other diseases.



520  Fig. 7 Screenshots of the EpiMix web application. **a,** Interactive data filters and visualization of functional CpG-
521  gene pair matrix. **b,** Visualization of the mixture model of the SLC16A4 gene in lung cancer. **c,** Genome-browser
522  style visualization of the lncRNA gene LINC00881 in lung cancer. **d,** Kaplan-Meier survival curves of patients with
523  different methylation states of the miRNA gene miR-34a in lung cancer.

25

524  EpiMix uses a beta mixture model to decompose the DNAme profiles in a patient

525  population. Using EpiMix, we can resolve the epigenetic subtypes within the patient

526  population and pinpoint the individuals carrying differential DNAme profiles. In this

527  study, we identified five DNAme subtypes in lung cancers using the DM values of

528  lncRNAs. Patients of subtype 2 had worse survival than patients of subtype 5,

529  indicating that the DNAme subtypes discovered by EpiMix had prognostic values. The

530  biological interpretation of DNAme subtypes requires the integration of data from other

531  modalities, such as genetic mutations, lifestyle history, and other etiological features.

532

533  In addition, EpiMix was able to detect abnormal DNAme that was present in only small

534  subsets of a patient cohort. In our simulation study, EpiMix detected more differentially

535  methylated CpGs compared to existing methods, when the differential methylation

536  occurred in only a small patient subset. Using the real lung cancer dataset (n = 457),

537  we identified miRNAs that were differentially methylated in only 1.1% of the patient

538  population and lncRNAs differentially methylated in 0.6% of the patient population. We

539  showed that over half of the miRNAs and lncRNAs were differentially methylated in

540  only less than 20% of the patients. This unique feature of EpiMix to detect differential

541  DNAme in small patient subsets enables us to identify novel epigenetic mechanisms

542  underlying disease phenotypes. It can also be used to discover new epigenetic

543  biomarkers and drug targets for improving personalized treatment.

544

545  Another feature of EpiMix is its ability to model DNAme at functionally diverse genomic

546  elements. This includes *cis*-regulatory elements within or surrounding protein-coding

547  genes, distal enhancers, and genes encoding miRNAs and lncRNAs. To model

548  DNAme at distal enhancers, we selected the enhancers from the ENCODE and

26

549 ROADMAP consortiums, in which enhancers of over a hundred human tissues and

550 cell lines were identified using the chromatin-state discovery (ChromHMM)[59]. Since

551 enhancers are cell-type specific, EpiMix allows the users to select enhancers of

552 specific cell types or tissues. In this study, we selected the enhancers of human blood

553 and T cells, leading to the discovery of 40,311 CpG of enhancers. In addition to

554 enhancers, many other regulatory elements were identified from the ROADMAP

555 studies[59]. These include active transcription start site proximal promoters, zinc finger

556 protein genes, bivalent regulatory elements, polycomb-repressed regions and many

557 others. By customizing the "chromatin state" parameter of EpiMix, users can target the

558 DNAme analysis to any of these regulatory modules.

559

560 Despite the critical biological functions of non-coding RNAs, there are no existing tools

561 that specifically analyze DNAme regulating their transcription. To analyze DNAme of

562 miRNA genes, we utilized the miRNA annotation from miRBase, the largest and

563 consistently updated knowledge base of miRNAs[60]. In addition, we selected CpGs at

564 miRNA promoters by using a recent database that integrates the information of miRNA

565 TSSs from 14 genome-wide studies across different human cell types and tissues[61].

566 This led to the discovery of 17,192 CpGs associated with 1,484 miRNAs in the HM450

567 array and 23,379 CpGs associated with 1,759 miRNAs in the EPIC array. With miRNA-

568 Seq data provided, EpiMix can select differential DNAme that was associated with

569 miRNA expression. Different from profiling protein-coding gene expression, measuring

570 miRNA expression requires special library preparation strategies that capture small

571 RNAs from total RNAs[62]. Users are preferentially needed to supply miRNA expression

572 data obtained from proper library preparation strategies.

573

574    Similarly, custom methods are needed to accurately quantify lncRNA expression from

575    RNA-Seq. We adopted the data processing pipeline developed from our previous

576    study[54]. With this pipeline, we combined the transcriptome annotations from

577    GENCODE and NONCODE. Raw sequencing reads were aligned to the combined

578    transcriptome reference and quantified using the Kallisto-Sleuth algorithm[56,57]. Using

579    this pipeline, we detected the expression of over 2,400 lncRNA genes. In this study,

580    we have used our pipeline to generate lncRNA expression profiles for all the cancers

581    in the TCGA database, and users can retrieve these data with EpiMix. Note, if users

582    plan to use EpiMix on non-TCGA datasets, they are encouraged to use this pipeline

583    to profile lncRNA expression.

584

585    Future work will aim to extend the use of EpiMix to whole-genome bisulfite sequencing

586    and to further improve the scalability. Furthermore, the rapid development of single-

587    cell technologies enables co-assay of DNAme and gene expression in thousands of

588    cells. EpiMix can be used to identify differential DNAme that was present in only small

589    subsets of a cell population. Therefore, a joint analysis of single cell methylome and

590    transcriptome holds great promise for substantiating our goals, and the analytical

591    framework presented here will be a valuable component for future research and

592    applications.

593

594    # Methods

595

596    **Data downloading**

597

598    The downloading module enables automated data downloading from the GEO

599    database and TCGA project. Alternatively, users can supply custom datasets

600    generated from their own studies. To retrieve data from GEO, we utilized the *getGEO*

601    function from the GEOquery R package (version 2.62)[63]. In this study, we downloaded

602    DNAme data and gene expression data using GEO accession number GSE114135.

603    The DNAme data were beta values ranging from 0 to 1, representing the proportion of

604    the methylated signal to the total signal. The gene expression data were TMM values.

605    Other formats of gene expression data are also acceptable (e.g., RPKM, TPM, FPKM

606    etc.). To retrieve data from TCGA, we used the Broad Institute Firehose tool

607    (Firehose)[64]. We downloaded level three DNAme data and gene expression data. The

608    downloaded data have been preprocessed for several steps, including removing

609    problematic rows, removing redundant columns, reordering the columns and sorting

610    the data by gene name. With the Regular mode, we used log-transformed RSEM

611    values. With the miRNA mode, we used the pri-miRNA expression data with log-

612    transformed RPKM values.

613

614    **Preprocessing**

615

616    The majority of datasets obtained from the TCGA and GEO databases have already

617    been preprocessed for a few steps. EpiMix's contribution to preprocessing includes

618    missing value imputation, removal of single-nucleotide polymorphism (SNP) probe

619    and batch effect correction. Users can also select to remove CpGs on sex

620    chromosomes. We then removed CpGs and samples with more than 20% missing

621    values, and imputed missing values on the remaining dataset using the k-nearest

622    neighbor (KNN) algorithm with K = 15.

623

624    Data from large patient cohorts were typically collected in technical batches.

625    Systematic variances between technical batches may affect downstream data

626    analysis and interpretation. To correct batch effects, we implemented two alternative

627    approaches: (1) an anchor-based data integration approach adapted from the Seurat

628    package (version 4.0.1)[65] and (2) an empirical Bayes regression approach, Combat[66].

629    The anchor-based approach uses canonical correlation analysis and mutual nearest

630    neighbors to identify shared subpopulations (termed "anchors") across different

631    datasets and then uses a non-linear transformation to integrate the data. To identify

632    the anchors, we used the "vst" method to select the top 10% variable features.

633    Effective batch effect removal was confirmed using the PCA-based ANOVA analysis.

634    Alternatively, the batch effect can be corrected with the Combat algorithm[58]. We found

635    that the anchor-based approach was more time efficient compared to the Combat.

636    When tested on the lung cancer dataset, the former approach completed the batch

637    correction within 2 hours, whereas the Combat consumed more than 48 hours.

638

639    **CpG annotation and filtering**

640

641    *Regular mode*

642

643    The Regular mode aims to model DNAme at cis-regulatory elements within or

644    immediately surrounding protein-coding genes. We paired each CpG site to the

645    nearest genes based on the hg38 manifest generated from Zhou et al.[67]. Unique CpG-

646    gene pairs were identified, where a CpG was either within the gene body or at the

647    immediately surrounding area. Users can restrict the analysis to the promoters,

30

648     defined as 2 kb upstream and 500 bp downstream (-2000bp ~ +500bp) of the

649     transcription start sites (TSSs). TSS information was retrieved from Ensembl using the

650     *biomaRt* R package (version 2.50.1)[68].

651

652     *Enhancer mode*

653

654     The Enhancer mode aims to model DNAme specifically at distal enhancers. Therefore,

655     we selected the distal CpGs that were at least 2 kb away from any known TSSs. Users

656     can customize this distance based on their needs. To select the CpGs within

657     enhancers, we used the enhancer database established from the ENCODE and

658     ROADMAP consortiums, in which enhancers of over a hundred human tissues and

659     cell lines were identified using the chromatin-state discovery (ChromHMM)[59]. We

660     looked for the DNA elements associated with the chromatin states of active enhancers

661     ("EnhA1" and "EnhA2") and genic enhancers ("EnhG1" and "EnhG2"). Since

662     enhancers are cell-type specific, EpiMix allows users to select enhancers of specific

663     cell types or tissue groups. In this study, we selected the enhancers of human blood

664     and T cells, leading to the discovery of 40,311 CpGs of enhancers. For each CpG, we

665     retrieved 20 nearby genes as candidate genes targets. This gene number was

666     determined by the previous studies showing that many of the enhancers can regulate

667     a gene within a 10-gene distance[29,69,70]. Genes that are positively regulated by the

668     enhancers should have a negative relationship between DNAme and gene

669     expression[36,71,72]. Therefore, we performed a one-tailed Wilcoxon rank-sum test on

670     each enhancer-gene pair to select the enhancers whose methylation states were

671     inversely associated with the gene expression. The raw $P$ value from the Wilcoxon

672     rank-sum test was adjusted using a permutation approach[73], where an empirical $P$

31

673   value was determined by ranking the raw *P* value in a set of permutation *P* values from

674   testing the expression of a set of randomly selected 1,000 genes (**Supplementary**

675   **Fig.2**).

676

677   *miRNA mode*

678

679   MicroRNAs are commonly classified into "intergenic" or "intronic" based on their

680   genomic locations. Intergenic miRNAs are found at previously unannotated human

681   genome and are transcribed from their own unique promoters as independent entities.

682   In contrast, intronic miRNAs are believed to share promoters with their host genes and

683   co-transcribed from respective hosts. Recent evidence shows that some intronic

684   miRNAs can also be transcribed independently from their host genes, suggesting they

685   have their own independent promoters[74]. To select CpGs associated with miRNAs,

686   we used a combined strategy. First, we obtained the most recent annotation of

687   miRNAs from miRBase (version 22.1)[60]. For each miRNA gene, we selected CpGs

688   that were located within 5 kb upstream and 5 kb downstream. Second, we selected

689   CpGs at miRNA promoters by using a recent database that integrates miRNA TSS

690   information from 14 genome-wide studies across different human cell types and

691   tissues[61]. We included CpGs located with miRNA promoters defined as 2000 bp

692   upstream and 1000 bp downstream of the TSSs. This combined feature selection

693   strategy resulted in the discovery of 17,192 CpGs associated with 1,484 miRNAs in

694   the HM450 array and 23,379 CpGs associated with 1,759 miRNAs in the EPIC array.

695

696   *lncRNA mode*

697

698  The mechanisms for transcriptional regulation of lncRNAs are similar to protein-coding

699  genes. We first selected lncRNA-coding genes using the GENCODE annotation

700  (Version 36). We then selected CpGs associated with each lncRNA based on the

701  hg38 manifest generated from Zhou et al.[67]. Unique CpG-gene pairs were identified,

702  where a CpG was either located within the gene body or at the immediately

703  surrounding area. This resulted in the discovery of 98,320 CpGs associated with

704  11,280 lncRNAs in the HM450 array and 184,816 CpGs associated with 15,392

705  lncRNAs in the EPIC array. Alternatively, users can select to focus the analysis at

706  lncRNA promoters, defined as 2 kb upstream and 500 bp downstream (-2000bp ~

707  +500bp) of the TSSs. The TSS information was retrieved from Ensembl using the

708  *biomaRt* R package (version 2.50.1)[68].

709

710  **CpG site clustering and smoothing (optional features)**

711

712  *Clustering*

713

714  Modeling the DNAme at all individual CpG sites can be computationally expensive. In

715  addition, it can also lead to overfitting of DNAme data in identifications of patient

716  subsets. Since the DNAme at adjacent CpGs are strongly correlated, we implemented

717  an optional feature that allows users to group the correlated CpGs into CpG clusters.

718  First, we used the average linkage hierarchical clustering algorithm to cluster CpGs of

719  a single gene into clusters. Then we cut off the hierarchical tree at a Pearson

720  correlation threshold of 0.4 to define CpG clusters and single CpG sites when they do

721  not correlate with other sites. For each CpG site cluster, we used the mean levels of

722  DNAme of the CpGs to represent the cluster DNAme, resulting in potentially multiple

33

723    CpG site clusters representing a single gene. The DNAme modeling can then be

724    performed at each separate CpG site or CpG site cluster.

725

726    *Smoothing*

727

728    Smoothing is another technique frequently used in removing noise and increasing

729    statistical power in analyzing whole-genome bisulfite sequencing data[6]. This

730    technique estimates localized DNAme levels using data of adjacent CpGs at a user-

731    specified genomic window. EpiMix allows users to smooth the DNAme data using local

732    likelihood smoothing[75]. Since the number of CpGs is lower in array-based data than

733    in bisulfite sequencing data, using smoothing on array-based data should be taken

734    with cautions.

735

736    **Methylation modeling**

737

738    After preprocessing, the methylation data are beta values bounded between 0 and 1,

739    representing the proportion of the methylated signal to the total signal. When the study

740    population is large, the beta values can be assumed to come from multiple underlying

741    probability distributions, in our case, beta distributions. To model the DNAme, we fit a

742    beta mixture model to the methylation values at each CpG site (or CpG site cluster).

743    Let $y_i$ denote the beta value from subject $i$ at a CpG site, where $i \in \{1, \dots, n\}$, and $n$

744    represents the total number of subjects. Let $k$ denote the class membership of subject

745    $i$, where $k \in \{1, \dots, K\}$, and $K$ represents the total number of components in the mixture.

746    Assume subject $i$ belongs to component $k$ with probability $\eta_k$, we will have $\sum_{k=1}^{K} \eta_k = $

747    1. Subsequently, the likelihood contribution from subject $i$ is:

34

748
$$f(Y_i = y_i) = \sum_{k=1}^{K} \eta_k \frac{y_i^{\alpha_k-1}(1-y_i)^{\beta_k-1}}{B(\alpha_k, \beta_k)}$$

749

750 where $B(\alpha_k, \beta_k) = \int_0^1 t^{\alpha_k-1}(1-t)^{\beta_k-1} \, dt$ is the beta function. Since the population

751 contains $n$ subjects, the log-likelihood for the complete dataset is

752
$$l(\alpha, \beta, \eta) = \sum_{i=1}^{n} \log\{f(Y_i = y_i)\}$$

753 The goal of our modeling is to estimate the $\alpha, \beta, \eta$ parameters of each component that

754 best fit the methylation values. Let $\theta = \{\alpha_1, \beta_1, \eta_1 \dots, \alpha_k, \beta_k, \eta_k\}$ be a vector of

755 parameters that define the shape of each component in the mixture. We used the

756 expectation–maximization (EM) algorithm[76] to iteratively maximize the log-likelihood

757 and update the conditional probability that $y_i$ comes from the $k\,th$ component.

758

759 To determine the best number of components $K$, we used The Bayesian Information

760 Criterion (BIC) for model selection and to avoid overfitting:

761
$$BIC = \log(n)\,(3K) - 2 \times \sum_{i=1}^{n} \log\{f(Y_i = y_i)\}$$

762 This process involves iteratively adding a new mixture component if the BIC improves.

763 Each mixture component represents a subset of samples for whom a particular

764 DNAme state is observed.

765

766 **Identifications of differentially methylated CpGs**

767

768 If data of a control group are provided, we can determine whether a CpG site (or CpG

769 site cluster) was hypo- or hyper-methylated by comparing its methylation levels in the

35

770  experimental group to its counterpart in the control group. We first performed beta

771  mixture modeling on each CpG site (or CpG site cluster) to identify the mixture

772  components using data from the experimental group, and the methylation levels of

773  each of the mixture components were compared to the mean methylation levels of the

774  control group. This methodology is based on the assumption that the DNAme profile

775  is heterogenous across different subjects in the experimental (i.e., disease) group but

776  is homogenous in the control group. For instance, the DNAme profile is expected to

777  be different across cancer patients due to the difference in subtypes or driver

778  mutations, but in normal tissues the DNAme should be relatively homogenous. In

779  addition, the number of subjects in the experimental group is typically higher than the

780  control group (e.g., TCGA projects). To determine the significant difference between

781  the experimental and the control group, we used a Wilcoxon rank-sum to calculate the

782  *P*-value, and multiple comparison was corrected with the false discovery rate (FDR).

783  The Q-value threshold was set to 0.05. In addition, we required a minimum difference

784  of 0.10 based on the platform sensitivity reported previously[77].

785

786  **Identifications of differential DNAme that was associated with transcription**

787

788  If sample-matched gene expression data are provided, we can select the CpGs whose

789  methylation states were significantly associated with gene expression. In this study,

790  we focused on the identification of DNAme that represses gene expression. However,

791  users have the option to identify DNAme that is positively correlated with gene

792  expression. For each CpG-gene pair, we used a one-tailed Wilcoxon rank-sum test to

793  compare the mean levels of gene expression in patients showing an abnormal

794  methylation state (hypo- or hyper-methylation state) to those with a normal methylation

795    state. If a CpG was hypomethylated, we examined that the hypomethylated patients

796    have higher gene expression levels compared to the normally methylated patients.

797    Vice versa, if a CpG was hypermethylated, we tested that the hypermethylated

798    patients have lower gene expression levels compared to the normally methylated

799    patients. If a CpG was dual methylated (i.e., some samples were hypomethylated,

800    while some others were hypermethylated), we tested that the hypomethylated patients

801    have higher gene expression levels compared to the hypermethylated patients. Since

802    a gene is typically paired with multiple CpGs, we adjusted the *P*-value using FDR to

803    correct multiple comparisons. To select functionally significant CpG-gene pairs, we set

804    the maximum threshold of the adjusted *P*-value to 0.01.

805

806    **Simulation study**

807

808    The goal of the simulation studies was to assess the sensitivity of EpiMix to detect

809    differential DNAme present in only specific subsets of a population. The studies were

810    performed by creating synthetic CpG sites and synthetic populations. First, we filtered

811    CpGs showing statistically similar DNAme levels that fit a unimodal beta distribution

812    from the activation group and from the quiescent group (n = 103 samples per group).

813    We then randomly sampled a subset of CpGs (n = 300) from the quiescent group as

814    the baselines. The average DNAme levels (beta values) of the CpGs in the baseline

815    group ranged from 0.1 to 0.9, with a mean DNAme level of 0.6. Second, since the

816    magnitude of changes in DNAme levels can be a critical factor affecting sensitivity, we

817    created synthetic CpGs. For each CpG of the baseline group, we paired it with a

818    subset of CpGs from the activation group, such that the differences in the mean beta

819    values ($\Delta beta$) between the the activation group and the baseline group ranged from

37

820    0.1 to 0.7, where $\Delta beta \in \{0.10, 0.15, 0.20, 0.25, 0.30, 0.40, 0.50, 0.60, 0.70\}$. This

821    resulted in a total of 2,700 synthetic CpGs. Third, since our goal was to detect

822    differential DNAme that was present in only a subset of the population, we created

823    synthetic populations. For each synthetic CpG,  we controlled the number of samples

824    from the activation group to be combined with the baseline group, such that the final

825    proportion ($P$) of samples from the activation group in the combined datasets ranged

826    from            0.01            to            0.50,            where            $P \in$

827    $\{0.01, 0.02, 0.05, 0.08, 0.10, 0.15, 0.20, 0.25, 0.30, 0.35, 0.40, 0.45, 0.50\}$. Finally,  we  ran

828    the EpiMix algorithm on each synthetic CpG and assessed whether it could pick up

829    the differentially methylated signals in the synthetic populations.

830

831    **Benchmark with existing methods**

832

833    We benchmarked the performance of EpiMix with other existing methods, including

834    Minfi[10], iEVORA[27] and RnBeads[12,28].

835

836    Minfi includes a differential methylation step based on an F-test. We first transformed

837    beta values to M values, and the differential methylation analysis was performed with

838    the *dmpFinder* function. We set the significant *P*-value and *Q*-value thresholds to 0.05.

839

840    iEVORA is a two-step algorithm that selects differentially variable and differentially

841    methylated CpGs. The first step is to identify differentially variable CpGs using a

842    Bartlett's test. The Bartlett's test assesses the equity of variances between the

843    experimental and the control group. If in the experimental group, there are samples

844    showing large differences (outliers) in DNAme versus other samples, the Bartlett's test

38

845   can detect such abnormality. The second step is to select the differentially variable

846   CpGs that were also differentially methylated. The differential methylation analysis is

847   performed by comparing the mean levels of DNAme of all the samples in the

848   experimental group to the control group.  We used the default parameters of the

849   functions, with a *Q*-value (FDR) threshold of 0.001 for testing differential variability and

850   *P*-value threshold of 0.05 for testing differential methylation means. In our stimulation

851   studies, we found that iEVORA was able to identify differentially variable CpGs even

852   when the abnormal methylation was present in only a small subset of the experimental

853   group. However, since the algorithm does not identify which subjects were abnormally

854   methylated, and in the differential methylation step, it still compares the mean levels

855   of DNAme of the entire experimental group to the control group, the differential

856   methylation test could not generate statistically significant results.

857

858   RnBeads uses hierarchical linear models as implemented in the limma package to

859   identify differential methylated CpGs. We set the differential methylation *P*-value

860   threshold to 0.05.

861

862   **Imputation of cell-type-specific DNAme and gene expression data**

863

864   DNAme and gene expression are known to be cell-type specific. When the DNAme

865   were measured at the tissue ("bulk") level, the differential DNAme profiles between

866   patient subjects may result from the differences in tissue compositions. From a clinical

867   perspective, tissue composition is meaningful in classifications of tumor subtypes and

868   prediction of treatment response. However, from a biological perspective, users may

869   be interested in identifying the differential DNAme present in specific cell types. EpiMix

870   focuses on the identification of differential DNAme across patient individuals. To

871   resolve the confounding effect from tissue heterogeneity, we used previously validated

872   algorithms to infer cell-type proportions and cell-type specific methylomes and

873   transcriptomes (**Supplementary Fig.3**). First, we used CIBERSORTx[45], a reference-

874   based computational algorithm, to estimate cell-type proportions from bulk gene

875   expression data in each tumor and normal tissue, and deconvolute bulk gene

876   expression data into cell-type specific signals. This method leveraged the established

877   signature gene expression matrices for experimentally purified cells from normal

878   tissues and lung cancers[45]. Second, we used Tensor Composition Analysis (TCA)[44]

879   to deconvolute bulk DNAme data into cell-type-specific data based on the estimated

880   cell-type proportions in each tissue. The output from TCA was the methylome of each

881   cell type in each individual. In addition to these methods, users can leverage other

882   existing tools to adjust the effects from tissue compositions before inputting the data

883   to EpiMix[78–83].

884

885   **Genomic distribution of the differentially methylated CpGs**

886

887   Genomic coordinates of the TSSs of the methylation-driven genes were retrieved from

888   Ensembl using the *biomaRt* R package (version 2.50)[68]. Exons and Introns of the

889   protein-coding      genes      were      retrieved      from      the      TxDb      object

890   (*TxDb.Hsapiens.UCSC.hg38.knownGene*) (version 3.14)[84]. The GenomicRanges R

891   package (version 1.46)[85] was used to identify the differentially methylated CpGs

892   located within promoters, exons and introns.

893

894   **Motif enrichment analysis**

895

896    TF binding motifs were retrieved from HOCOMOCO, a comprehensive database for

897    TF binding sites[86]. HOMER (Hypergeometric Optimization of Motif EnRichment) was

898    used to find motif occurrences in a ±250bp region around each differentially

899    methylated regions (DMRs). We then combined all the DMRs to identify enriched

900    motifs. Enrichments were quantified using Fisher's exact test and multiple

901    comparisons were adjusted with the Benjamini-Hochberg procedure. To calculate the

902    enrichment Odds Ratio, we used all the distal CpGs as the background probes and

903    the functional CpGs of enhancers as the target probes. We set the significant $P$ value

904    cutoff to 0.05 and the smallest lower boundary of 95% confidence interval for Odds

905    Ratio to 1.1. The enrichment analysis was performed using the *get.enriched.motif*

906    function from the *ELMER* library (version 3.14) in R[11].

907

908    **Enrichment analysis of chromatin modifications**

909

910    Enrichment analysis of histone modifications at the DMRs was performed using the

911    Genomic Hyperbrowser GSUITE of tools[87]. A suite of tracks representing different

912    chromatin features for human naïve T cells (Epigenome ID: E038) and lung

913    (Epigenome ID: E096) were retrieved from the ENCODE and ROADMAP

914    consortiums[59]. To determine which tracks in the suite exhibit the strongest similarity

915    by co-occurrence to the DMRs, the Forbes coefficient was used to obtain rankings of

916    tracks, and Monte Carlo simulations were used to define a statistical assessment of

917    the robustness of the rankings using randomization of genomic regions covered by

918    the entire HM450 or EPIC array, and compute test statistics.

919

920 **Functional enrichment analysis**

921

922 Protein-coding genes

923

924 EpiMix provides an user interface to the *enrichGO* and *enrichKEGG* functions of the

925 *clusterProfiler* R package (version 4.2.1)[88]. This enables gene set analysis of the

926 methylation-driven genes using the gene ontology (GO) and KEGG datasets. Over-

927 represented biological pathways in the methylation-driven genes were identified using

928 the hypergeometric testing[88]. Enrichment results can be retrieved in a tabular format

929 or visualized in several different ways. To perform the GO analysis, we set the

930 significant $P$ value to 0.05 and Q value to 0.20. Highly similar GO terms were removed

931 with a cutoff $P$ value of 0.60 to retain the most representative terms.

932

933 miRNAs

934

935 To obtain the target genes of the differentially methylated miRNAs, we queried

936 miRTarBase with the *miRnetR* package[89]. Of the 144 differentially methylated miRNAs

937 in lung cancer, we identified 7,088 target protein-coding genes of 26 miRNAs. We

938 simplified this network by selecting the genes that were targeted by at least five

939 miRNAs. KEGG pathway analysis was then performed on the miRNA target genes

940 with hypergeometric testing.

941

942 lncRNAs

943

42

944    To carry out functional annotation and pathway analysis of the differentially methylated

945    lncRNAs, we used the ncFANs V2.0 server (http://ncfans.gene.ac/)[58]. The genes in

946    the significant CpG-gene pair matrix generated from EpiMix can be directly used as

947    an input to ncFANs. NcFANs assigns the functions of protein-coding genes to lncRNAs

948    based on pre-built co-expression networks in various normal tissues and cancers. We

949    used the co-expression network built in the lung adenocarcinoma dataset from TCGA,

950    and we set the correlation coefficient between lncRNAs and proteins-genes to 0.4 and

951    the cutoff of the topological overlap measure similarity to 0.01.

952

953    **Biomarker identification and survival analysis**

954

955    Patient clinical data were retrieved from TCGA using the Firehose tool[64]. Alternatively,

956    users can provide EpiMix with survival data if using their own datasets. We selected

957    the CpGs with at least two methylation states. For each CpG, we fit a Cox proportional

958    hazards regression model to assess the effect of methylation states on patient survival

959    time. The log-rank test was used to compare the survival curve and to calculate the

960    significant *P*-value. *P* < 0.05 was considered as significant. The Kaplan-Meier survival

961    plots were generated with the *survminer* R package (version 0.4.9).

962

963    **Genome browser-style visualization**

964

965    EpiMix enables genome browser-style visualization of the genomic coordinates and

966    chromatin states of the differentially methylated genes and regions. We implemented

967    two different forms of visualization. The gene-centric form shows the DM values of all

968    the CpGs associated with a specific gene (e.g., **Fig.3f).** The CpG-centric form shows

969    a differentially methylated CpG and its upstream and downstream genes (e.g., **Fig.4e)**.

970    Users can specify the number of nearby genes to display. Genes whose expression

971    levels were significantly associated with the DNAme levels of the CpG are shown in

972    red.

973

974    DNase I sensitivity and histone modification levels were retrieved from the ENCODE

975    and ROADMAP consortiums[59]. By providing the Epigenome ID, users can retrieve

976    data corresponding to the investigated tissue or cell type. In this study, we extracted

977    the chromatin features for human naïve T cells (Epigenome ID: E038) and fetal lung

978    (Epigenome ID: E088). The genomic coordinates (X-axis) were established on the

979    hg19 genome built, and the enrichment signal (Y-axis) represents negative log10 of

980    the Poisson *P*-values. Human transcript annotation was retrieved from the TxDb object

981    (*TxDb.Hsapiens.UCSC.hg19.knownGene*) (version 3.2.2)[90]. The genomic coordinates

982    of the adjacent genes of the differentially methylated CpGs were retrieved from

983    Ensembl using the *biomaRt* R package (version 2.50.1)[68]. The  visualization was

984    implemented with the *karyoploteR* package (version 1.20.0)[91].

985

986    **Identifications of DNAme subtypes**

987

988    DNAme subtypes can be discovered by applying consensus clustering to the DM-

989    value matrix, where patients were clustered into robust and homogenous groups

990    (putative subtypes) based on their abnormal methylation profiles. Consensus

991    clustering was performed with the ConsensusClusterPlus R package (version

992    1.58.0)[92]. We used 1,000 rounds of k-means clustering and a maximum of K = 10

44

993    clusters. Selection of the best number of clusters was based on the visual inspection

994    of ConsensusClusterPlus output plots.

995

## Code availability

997

998    EpiMix is available as an R package on Bioconductor

999    (https://bioconductor.org/packages/devel/bioc/html/EpiMix.html). In addition, we also

1000   developed a web application (https://epimix.stanford.edu) for users to interactively

1001   visualize and explore the results from EpiMix.

1002

## References

1004   1.  Li, J. *et al.* Insights Into the Role of DNA Methylation in Immune Cell

1005       Development and Autoimmune Disease. *Front Cell Dev Biol* **9**, 757318 (2021).

1006   2.  Si, J. *et al.* Epigenome-wide analysis of DNA methylation and coronary heart

1007       disease: a nested case-control study. *eLife* vol. 10 Preprint at

1008       https://doi.org/10.7554/elife.68671 (2021).

1009   3.  Zheng, Y., Luo, L., Lambertz, I. U., Conti, C. J. & Fuchs-Young, R. Early dietary

1010       exposures epigenetically program mammary cancer susceptibility through Igf1-

1011       mediated expansion of the mammary stem cell compartment. *Cells* **11**, 2558

1012       (2022).

1013   4.  Akalin, A. *et al.* methylKit: a comprehensive R package for the analysis of

1014       genome-wide DNA methylation profiles. *Genome Biol.* **13**, R87 (2012).

1015   5.  Park, Y. & Wu, H. Differential methylation analysis for BS-seq data under

1016       general experimental design. *Bioinformatics* **32**, 1446–1453 (2016).

1017    6.    Korthauer, K., Chakraborty, S., Benjamini, Y. & Irizarry, R. A. Detection and

1018          accurate false discovery rate control of differentially methylated regions from

1019          whole genome bisulfite sequencing. *Biostatistics* **20**, 367–383 (2019).

1020    7.    Wang, X., Hao, D. & Kadarmideen, H. N. GeneDMRs: An R Package for Gene-

1021          Based Differentially Methylated Regions Analysis. *J. Comput. Biol.* **28**, 304–316

1022          (2021).

1023    8.    Wang, D. *et al.* IMA: an R package for high-throughput analysis of Illumina's

1024          450K Infinium methylation data. *Bioinformatics* **28**, 729–730 (2012).

1025    9.    Warden, C. D. *et al.* COHCAP: an integrative genomic pipeline for single-

1026          nucleotide resolution DNA methylation analysis. *Nucleic Acids Res.* **41**, e117

1027          (2013).

1028    10.   Aryee, M. J. *et al.* Minfi: a flexible and comprehensive Bioconductor package for

1029          the analysis of Infinium DNA methylation microarrays. *Bioinformatics* **30**, 1363–

1030          1369 (2014).

1031    11.   Silva, T. C. *et al.* ELMER v.2: an R/Bioconductor package to reconstruct gene

1032          regulatory networks from DNA methylation and transcriptome profiles.

1033          *Bioinformatics* **35**, 1974–1977 (2019).

1034    12.   Müller, F. *et al.* RnBeads 2.0: comprehensive analysis of DNA methylation data.

1035          *Genome Biology* vol. 20 Preprint at https://doi.org/10.1186/s13059-019-1664-9

1036          (2019).

1037    13.   Shaknovich, R. *et al.* DNA methylation signatures define molecular subtypes of

1038          diffuse large B-cell lymphoma. *Blood* **116**, e81-9 (2010).

1039    14.   Chen, X., Zhang, J. & Dai, X. DNA methylation profiles capturing breast cancer

1040          heterogeneity. *BMC Genomics* **20**, 823 (2019).

1041  15. Schenkel, L. C. *et al.* DNA methylation epi-signature is associated with two

1042      molecularly and phenotypically distinct clinical subtypes of Phelan-McDermid

1043      syndrome. *Clin. Epigenetics* **13**, 2 (2021).

1044  16. Ghandi, M. *et al.* Next-generation characterization of the Cancer Cell Line

1045      Encyclopedia. *Nature* **569**, 503–508 (2019).

1046  17. Partridge, E. C. *et al.* Occupancy maps of 208 chromatin-associated proteins in

1047      one human cell type. *Nature* **583**, 720–728 (2020).

1048  18. Schoenfelder, S. & Fraser, P. Long-range enhancer–promoter contacts in gene

1049      expression control. *Nat. Rev. Genet.* **20**, 437–455 (2019).

1050  19. Yao, L., Shen, H., Laird, P. W., Farnham, P. J. & Berman, B. P. Inferring

1051      regulatory element landscapes and transcription factor networks from cancer

1052      methylomes. *Genome Biol.* **16**, 105 (2015).

1053  20. Cribbs, A. P. *et al.* Methotrexate restores regulatory T cell function through

1054      demethylation of the FoxP3 upstream enhancer in patients with rheumatoid

1055      arthritis. *Arthritis rheumatol.* **67**, 1182–1192 (2015).

1056  21. Wang, L. *et al.* A MicroRNA Linking Human Positive Selection and Metabolic

1057      Disorders. *Cell* **183**, 684-701.e14 (2020).

1058  22. Statello, L., Guo, C.-J., Chen, L.-L. & Huarte, M. Gene regulation by long non-

1059      coding RNAs and its biological functions. *Nat. Rev. Mol. Cell Biol.* **22**, 96–118

1060      (2021).

1061  23. Watanabe, K. *et al.* Genome structure-based screening identified epigenetically

1062      silenced microRNA associated with invasiveness in non-small-cell lung cancer.

1063      *Int. J. Cancer* **130**, 2580–2590 (2012).

1064  24. Zhang, M., Wu, J., Zhong, W., Zhao, Z. & He, W. DNA-methylation-induced

1065      silencing of DIO3OS drives non-small cell lung cancer progression via activating

1066      hnRNPK-MYC-CDC25A axis. *Mol Ther Oncolytics* **23**, 205–219 (2021).

1067  25. Gentleman, R. C. *et al.* Bioconductor: open software development for

1068      computational biology and bioinformatics. *Genome Biol.* **5**, R80 (2004).

1069  26. Martino, D. *et al.* Epigenetic dysregulation of naive CD4+ T-cell activation genes

1070      in childhood food allergy. *Nat. Commun.* **9**, 3308 (2018).

1071  27. Teschendorff, A. E. *et al.* DNA methylation outliers in normal breast tissue

1072      identify field defects that are enriched in cancer. *Nat. Commun.* **7**, 10478 (2016).

1073  28. Assenov, Y. *et al.* Comprehensive analysis of DNA methylation data with

1074      RnBeads. *Nat. Methods* **11**, 1138–1140 (2014).

1075  29. Mifsud, B. *et al.* Mapping long-range promoter contacts in human cells with high-

1076      resolution capture Hi-C. *Nat. Genet.* **47**, 598–606 (2015).

1077  30. Javierre, B. M. *et al.* Lineage-Specific Genome Architecture Links Enhancers

1078      and Non-coding Disease Variants to Target Gene Promoters. *Cell* **167**, 1369-

1079      1384.e19 (2016).

1080  31. Jin, F. *et al.* A high-resolution map of the three-dimensional chromatin

1081      interactome in human cells. *Nature* **503**, 290–294 (2013).

1082  32. Rao, S. S. P. *et al.* A 3D map of the human genome at kilobase resolution

1083      reveals principles of chromatin looping. *Cell* **159**, 1665–1680 (2014).

1084  33. Sanyal, A., Lajoie, B. R., Jain, G. & Dekker, J. The long-range interaction

1085      landscape of gene promoters. *Nature* **489**, 109–113 (2012).

1086  34. Bietz, A., Zhu, H., Xue, M. & Xu, C. Cholesterol Metabolism in T Cells. *Front.*

1087      *Immunol.* **8**, 1664 (2017).

1088  35. Rasmussen, K. D. *et al.* TET2 binding to enhancers facilitates transcription

1089       factor recruitment in hematopoietic cells. *Genome Res.* **29**, 564–575 (2019).

1090  36. Wang, L. *et al.* TET2 coactivates gene expression through demethylation of

1091       enhancers. *Sci Adv* **4**, eaau6986 (2018).

1092  37. Li, P. *et al.* BATF–JUN is critical for IRF4-mediated transcription in T cells.

1093       *Nature* **490**, 543–546 (2012).

1094  38. Glasmacher, E. *et al.* A Genomic Regulatory Element That Directs Assembly

1095       and Function of Immune-Specific AP-1–IRF Complexes. *Science* vol. 338 975–

1096       980 Preprint at https://doi.org/10.1126/science.1228309 (2012).

1097  39. Furuta, M. *et al.* miR-124 and miR-203 are epigenetically silenced tumor-

1098       suppressive microRNAs in hepatocellular carcinoma. *Carcinogenesis* **31**, 766–

1099       776 (2010).

1100  40. Baer, C., Claus, R. & Plass, C. Genome-wide epigenetic regulation of miRNAs

1101       in cancer. *Cancer Res.* **73**, 473–477 (2013).

1102  41. Cancer Genome Atlas Research Network. Comprehensive molecular profiling of

1103       lung adenocarcinoma. *Nature* **511**, 543–550 (2014).

1104  42. Bloushtain-Qimron, N. *et al.* Cell type-specific DNA methylation patterns in the

1105       human breast. *Proc. Natl. Acad. Sci. U. S. A.* **105**, 14076–14081 (2008).

1106  43. Schmidl, C. *et al.* Lineage-specific DNA methylation in T cells correlates with

1107       histone methylation and enhancer activity. *Genome Res.* **19**, 1165–1174 (2009).

1108  44. Rahmani, E. *et al.* Cell-type-specific resolution epigenetics without the need for

1109       cell sorting or single-cell biology. *Nat. Commun.* **10**, 3417 (2019).

1110  45. Newman, A. M. *et al.* Determining cell type abundance and expression from bulk

1111       tissues with digital cytometry. *Nat. Biotechnol.* **37**, 773–782 (2019).

46. Liang, G. *et al.* miR-196b-5p–mediated downregulation of TSPAN12 and GATA6 promotes tumor progression in non-small cell lung cancer. *Proc. Natl. Acad. Sci. U. S. A.* **117**, 4347–4357 (2020).

47. Long, X. *et al.* MicroRNA-99a Suppresses Breast Cancer Progression by Targeting FGFR3. *Front. Oncol.* **9**, 1473 (2019).

48. Liu, L. *et al.* MicroRNA-29c functions as a tumor suppressor by targeting VEGFA in lung adenocarcinoma. *Mol. Cancer* **16**, 50 (2017).

49. Tang, R. *et al.* Downregulation of MiR-30a is Associated with Poor Prognosis in Lung Cancer. *Med. Sci. Monit.* **21**, 2514–2520 (2015).

50. Zhao, K. *et al.* Circulating microRNA-34 family low expression correlates with poor prognosis in patients with non-small cell lung cancer. *J. Thorac. Dis.* **9**, 3735–3746 (2017).

51. Chen, Y. *et al.* miRNA-148a serves as a prognostic factor and suppresses migration and invasion through Wnt1 in non-small cell lung cancer. *PLoS One* **12**, e0171751 (2017).

52. Huang, H.-Y. *et al.* miRTarBase update 2022: an informative resource for experimentally validated miRNA-target interactions. *Nucleic Acids Res.* **50**, D222–D230 (2022).

53. Derrien, T. *et al.* The GENCODE v7 catalog of human long noncoding RNAs: analysis of their gene structure, evolution, and expression. *Genome Res.* **22**, 1775–1789 (2012).

54. Zheng, H., Brennan, K., Hernaez, M. & Gevaert, O. Benchmark of long non-coding RNA quantification for RNA sequencing of cancer samples. *Gigascience* **8**, (2019).

1136    55. Fang, S. *et al.* NONCODEV5: a comprehensive annotation database for long

1137        non-coding RNAs. *Nucleic Acids Res.* **46**, D308–D314 (2018).

1138    56. Bray, N. L., Pimentel, H., Melsted, P. & Pachter, L. Near-optimal probabilistic

1139        RNA-seq quantification. *Nat. Biotechnol.* **34**, 525–527 (2016).

1140    57. Pimentel, H., Bray, N. L., Puente, S., Melsted, P. & Pachter, L. Differential

1141        analysis of RNA-seq incorporating quantification uncertainty. *Nat. Methods* **14**,

1142        687–690 (2017).

1143    58. Zhang, Y. *et al.* ncFANs v2.0: an integrative platform for functional annotation of

1144        non-coding RNAs. *Nucleic Acids Res.* **49**, W459–W468 (2021).

1145    59. Roadmap Epigenomics Consortium *et al.* Integrative analysis of 111 reference

1146        human epigenomes. *Nature* **518**, 317–330 (2015).

1147    60. Kozomara, A., Birgaoanu, M. & Griffiths-Jones, S. miRBase: from microRNA

1148        sequences to function. *Nucleic Acids Res.* **47**, D155–D162 (2019).

1149    61. Wang, S., Talukder, A., Cha, M., Li, X. & Hu, H. Computational annotation of

1150        miRNA transcription start sites. *Brief. Bioinform.* **22**, 380–392 (2021).

1151    62. Motameny, S., Wolters, S., Nürnberg, P. & Schumacher, B. Next generation

1152        sequencing of miRNAs--strategies, resources and methods. *Genes* **1**, 70–84

1153        (2010).

1154    63. Davis, S. & Meltzer, P. S. GEOquery: a bridge between the Gene Expression

1155        Omnibus (GEO) and BioConductor. *Bioinformatics* **23**, 1846–1847 (2007).

1156    64. Center, B. I. T. G. D. A. Analysis-ready standardized TCGA data from Broad

1157        GDAC Firehose 2016_01_28 run. *Broad Institute of MIT and Harvard* (2016).

1158    65. Stuart, T. *et al.* Comprehensive Integration of Single-Cell Data. *Cell* **177**, 1888-

1159        1902.e21 (2019).

1160   66.  Johnson, W. E., Li, C. & Rabinovic, A. Adjusting batch effects in microarray

1161       expression data using empirical Bayes methods. *Biostatistics* **8**, 118–127

1162       (2007).

1163   67.  Zhou, W., Laird, P. W. & Shen, H. Comprehensive characterization, annotation

1164       and innovative use of Infinium DNA methylation BeadChip probes. *Nucleic Acids*

1165       *Res.* **45**, e22 (2017).

1166   68.  Durinck, S., Spellman, P. T., Birney, E. & Huber, W. Mapping identifiers for the

1167       integration of genomic datasets with the R/Bioconductor package biomaRt. *Nat.*

1168       *Protoc.* **4**, 1184–1191 (2009).

1169   69.  Ernst, J. *et al.* Mapping and analysis of chromatin state dynamics in nine human

1170       cell types. *Nature* **473**, 43–49 (2011).

1171   70.  Grubert, F. *et al.* Landscape of cohesin-mediated chromatin loops in the human

1172       genome. *Nature* **583**, 737–743 (2020).

1173   71.  Aran, D., Sabato, S. & Hellman, A. DNA methylation of distal regulatory sites

1174       characterizes dysregulation of cancer genes. *Genome Biol.* **14**, R21 (2013).

1175   72.  Cho, J.-W. *et al.* The importance of enhancer methylation for epigenetic

1176       regulation of tumorigenesis in squamous lung cancer. *Exp. Mol. Med.* **54**, 12–22

1177       (2022).

1178   73.  Sham, P. C. & Purcell, S. M. Statistical power and significance testing in large-

1179       scale genetic studies. *Nat. Rev. Genet.* **15**, 335–346 (2014).

1180   74.  Ramalingam, P. *et al.* Biogenesis of intronic miRNAs located in clusters by

1181       independent transcription and alternative splicing. *RNA* **20**, 76–87 (2014).

1182   75.  Loader, C. *Local regression and likelihood*. (Springer, 1999).

76. Dempster, A.P., Laird, N.M. and Rubin, D.B. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society, Series B (Methodological).*

77. Bibikova, M. *et al.* High-throughput DNA methylation profiling using universal bead arrays. *Genome Res.* **16**, 383–393 (2006).

78. Zou, J., Lippert, C., Heckerman, D., Aryee, M. & Listgarten, J. Epigenome-wide association studies without the need for cell-type composition. *Nat. Methods* **11**, 309–311 (2014).

79. Rahmani, E. *et al.* Sparse PCA corrects for cell type heterogeneity in epigenome-wide association studies. *Nat. Methods* **13**, 443–445 (2016).

80. Houseman, E. A. *et al.* Reference-free deconvolution of DNA methylation data and mediation by cell composition effects. *BMC Bioinformatics* **17**, 259 (2016).

81. Lutsik, P. *et al.* MeDeCom: discovery and quantification of latent components of heterogeneous methylomes. *Genome Biol.* **18**, (2017).

82. Teschendorff, A. E., Breeze, C. E., Zheng, S. C. & Beck, S. A comparison of reference-based algorithms for correcting cell-type heterogeneity in Epigenome-Wide Association Studies. *BMC Bioinformatics* **18**, 105 (2017).

83. Rahmani, E. *et al.* BayesCCE: a Bayesian framework for estimating cell-type composition from DNA methylation without the need for methylation reference. *Genome Biol.* **19**, 141 (2018).

84. Team, B. C. & Maintainer, B. P. TxDb. Hsapiens. UCSC. hg38. knownGene: Annotation package for TxDb object (s). *R package* (2019).

85. Lawrence, M. *et al.* Software for computing and annotating genomic ranges. *PLoS Comput. Biol.* **9**, e1003118 (2013).

1207    86. Kulakovskiy, I. V. *et al.* HOCOMOCO: towards a complete collection of

1208         transcription factor binding models for human and mouse via large-scale ChIP-

1209         Seq analysis. *Nucleic Acids Res.* **46**, D252–D259 (2018).

1210    87. Simovski, B. *et al.* GSuite HyperBrowser: integrative analysis of dataset

1211         collections across the genome and epigenome. *Gigascience* **6**, 1–12 (2017).

1212    88. Wu, T. *et al.* clusterProfiler 4.0: A universal enrichment tool for interpreting

1213         omics data. *Innovation (N Y)* **2**, 100141 (2021).

1214    89. Chang, L., Zhou, G., Soufan, O. & Xia, J. miRNet 2.0: network-based visual

1215         analytics for miRNA functional analysis and systems biology. *Nucleic Acids Res.*

1216         **48**, W244–W251 (2020).

1217    90. Carlson, M. & Maintainer, B. Txdb. hsapiens. ucsc. hg19. knowngene:

1218         Annotation package for txdb object (s). *R package version* **3**, (2015).

1219    91. Gel, B. & Serra, E. karyoploteR: an R/Bioconductor package to plot

1220         customizable genomes displaying arbitrary data. *Bioinformatics* **33**, 3088–3090

1221         (2017).

1222    92. Wilkerson, M. D. & Hayes, D. N. ConsensusClusterPlus: a class discovery tool

1223         with confidence assessments and item tracking. *Bioinformatics* **26**, 1572–1573

1224         (2010).

1225

# Acknowledgements

1230

# Funding

# Author contributions

Y.Z: study design, implementation, data analysis and interpretation of results, draft manuscript preparation

J.J: implementation and data analysis

K.B: study conception and design

O.G: study conception and design, resources, supervision

# Competing interests

The authors declare no competing interests.