# Comprehensive analyses of partially methylated domains and differentially methylated regions in esophageal cancer reveal both cell-type- and cancer-specific epigenetic regulation

Yueyuan Zheng[1,2], Benjamin Ziman[2,3], Allen S. Ho[4], Uttam K. Sinha[5], Li-Yan Xu[6], En-Min Li[6], H Phillip Koeffler[2], Benjamin P. Berman[7,*], De-Chen Lin[2,3,*]

[1]Clinical Big Data Research Center, Scientific Research Center, The Seventh Affiliated Hospital of Sun Yat-sen University, Shenzhen 518107, P.R. China.

[2]Department of Medicine, Samuel Ochin Comprehensive Cancer Institute, Cedars-Sinai Medical Center, Los Angeles, USA.

[3]Center for Craniofacial Molecular Biology, Herman Ostrow School of Dentistry, and Norris Comprehensive Cancer Center, University of Southern California, Los Angeles, CA 90033, USA

[4]Division of Otolaryngology-Head and Neck Surgery, Department of Surgery, Samuel Oschin Cancer Center, Cedars-Sinai Medical Center, Los Angeles, CA, USA

[5]Department of otolaryngology, Keck School of Medicine, University of Southern California

[6]The Key Laboratory of Molecular Biology for High Cancer Incidence Coastal Chaoshan Area, Shantou University Medical College, Guangdong, China

[7]Department of Developmental Biology and Cancer Research, Institute for Medical Research Israel-Canada, Faculty of Medicine, The Hebrew University of Jerusalem, Israel

*, To whom correspondence should be addressed:

Benjamin P. Berman, Ph.D: Department of Developmental Biology and Cancer Research, Institute for Medical Research Israel-Canada, Faculty of Medicine, The Hebrew University of Jerusalem, Israel. email: ben.berman@mail.huji.ac.il

De-Chen Lin, Ph.D: Center for Craniofacial Molecular Biology, University of Southern California, 2250 Alcazar Street – CSA 207D, Los Angeles, CA 90033. Phone: (323) 442-1220, Fax: (323) 442-2981, email: dechenli@usc.edu

## Abstract

28    As one of the most common malignancies, esophageal cancer has two subtypes, squamous cell

29    carcinoma (ESCC) and adenocarcinoma (EAC), arising from distinct cells-of-origin. However,

30    distinguishing cell-type-specific molecular features from cancer-specific characteristics has been

31    challenging. Here, we analyze whole-genome bisulfite sequencing (WGBS) data on 45

32    esophageal tumor and nonmalignant samples from both subtypes. We develop a novel sequence-

33    aware method to identify large partially methylated domains (PMDs), revealing profound

34    heterogeneity at both the methylation level (depth) and genomic distribution (breadth) of PMDs

35    across tumor samples. We identify subtype-specific PMDs, which are associated with repressive

36    transcription, chromatin B compartments and high somatic mutation rate. While the genomic

37    locations of these PMDs are pre-established in normal cells, the degree of loss is significantly

38    higher in tumors. We find that cell-type-specific deposition of H3K36me2 may underlie the

39    genomic distribution PMDs. At a smaller genomic scale, both cell-type- and cancer-specific

40    differentially methylated regions (DMRs) are identified for each subtype. Using binding motif

41    analysis within these DMRs, we show that a cell-type-specific transcription factor such as HNF4A

42    can maintain the binding sites that it establishes in normal cells, while being recruited to new

43    binding sites with novel partners such as FOSL1 in cancer. Finally, leveraging pan-tissue single-

44    cell and pan-cancer epigenomic datasets, we demonstrate that a substantial fraction of the cell-

45    type-specific PMDs and DMRs identified here in esophageal cancer, are actually markers that co-

46    occur in other cancers originating from related cell types. These findings advance our

47    understanding of the DNA methylation dynamics at various genomic scales in normal and

48    malignant states, providing novel mechanistic insights into cell-type- and cancer-specific

49    epigenetic regulations.

50

## Introduction

Ranking seventh in cancer incidence and sixth in mortality worldwide, esophageal carcinoma is highly aggressive and its patients have poor outcomes, with a 5-year survival rate lower than 20%[1,2]. Esophageal cancer comprises two major histologic subtypes: squamous cell carcinoma (ESCC) and adenocarcinoma (EAC). These two subtypes have distinct clinical characteristics. ESCC occurs predominantly in the upper and mid-esophagus; EAC is prevalent in the lower esophagus near the gastroesophageal junction (GEJ) and is associated with the precursor lesion known as Barrett's esophagus (BE). Biologically, ESCC arises from the squamous epithelial cells and has common features with other squamous cell carcinomas (SCC), such as head and neck SCC (HNSCC). In comparison, EAC has columnar cell features and shares many characteristics with tubular gastrointestinal adenocarcinomas. In particular, EAC is almost indistinguishable from GEJ adenocarcinoma in terms of genomic, biological and clinical features.

Epigenetically, multiple studies have reported molecular changes in esophageal cancer, especially at the DNA methylation level[3–9]. For example, methylation differences across thousands of loci between ESCC and EAC were noted by The Cancer Genome Atlas (TCGA)[3] consortium. However, these prior works focused largely on the analyses of DNA methylation in gene promoter regions, which only make up ~6% of all CpG sites across the human genome. DNA methylation is known to play important roles in other noncoding regions, such as enhancers[10], partially methylated domains (PMDs)[11], as well as repetitive elements[12]. Therefore, the DNA methylome of esophageal cancer awaits further and comprehensive characterization through genome-wide single-base resolution approaches such as whole-genome bisulfite sequencing (WGBS).

CpG island (CGI) promoter hypermethylation and global DNA hypomethylation are two epigenomic hallmarks in cancer[13]. In most healthy tissues, the vast majority of CpG sites (>80%) across the genome are fully methylated, except for the CpG-rich regions (e.g., CGIs) and other regulatory elements (predominantly enhancers)[14]. Indeed, focal demethylation is a reliable signature of gene promoters and enhancers, and their methylation levels are robustly maintained across healthy tissues. Additionally, methylation patterns of CpG sites across the genome are notably variable across various normal cell types, and can be grouped into cell-type-specific differentially methylated regions (DMRs), which are linked to cell-type-specific regulatory regions[14,15]. By contrast, abnormal CGI promoter hypermethylation is frequently observed in cancer, which is commonly associated with long-term and stable gene repression[14].

83  With respect to the global methylation loss, large hypomethylated blocks, also known as
84  PMDs, cover more than one-third of the genome and coincide with heterochromatin, chromatin
85  "B" compartment (determined by HiC) and nuclear lamina associated domains[16–18]. We and others
86  recently found that accumulation of PMD hypomethylation is linked to cumulative mitotic cell
87  divisions, late replication timing as well as the deposition of the histone mark H3K36me3[19,20].
88  Functionally, PMDs are associated with inactive gene transcription, heightened genomic
89  instability and may be accompanied by activation of transposable elements (TEs)[19,21]. While
90  incompletely understood, the majority of the PMD regions are possibly shared across
91  developmental lineages[19]. However, there are enough cell-type specific PMDs to differentiate
92  between different cancer cell types[17,22,23] and between different healthy cell types[24].

93  Several important questions on cell-type- and cancer-specific DMRs and PMDs await further
94  characterization, including: i) the degree of the regional specificity of these domains (i.e, the
95  proportions of DMR/PMD that are cell-type- and cancer-specific), ii) the functional significance of
96  DMRs and PMDs in cancer biology, and iii) underlying mechanisms of the alteration of DMRs and
97  PMDs during tumorigenesis. To address these questions, we performed analyses of WGBS data
98  generated from a cohort of 45 esophageal samples, including 21 ESCC and 5 nonmalignant
99  esophageal squamous (NESQ) tissues, as well as 12 EAC/GEJ tumors and 7 nonmalignant GEJ
100 (NGEJ) tissues (**Fig. 1A**). We chose esophageal cancer as the disease model considering that
101 the two subtypes are developed from distinct cell-of-origins, and we hypothesized that
102 characterization of their methylome profiles might reveal cell-type- and cancer-specific
103 methylation changes, together with underlying epigenetic mechanisms.

## Results

**Development of a novel sequence-aware calling method to identify PMDs**

106 To characterize the esophageal cancer methylome, we analyzed WGBS profiles of 45
107 esophageal samples from two different cancer subtypes and their corresponding nonmalignant
108 tissues[25] (**Fig. 1A, Supplementary Fig. 2A**). All of the nonmalignant esophageal squamous
109 (NESQ) tissues showed high inter-sample correlation despite that they were from two different
110 cohorts (**Supplementary Fig. 2B** and **Supplementary Table 1**). To analyze the overall
111 methylation pattern, we first investigated the methylation level at various genomic domains (**Fig.
112 1B**). As anticipated, both global hypomethylation (especially in common PMDs, defined as shared
113 PMDs identified from 40 different cancer types[19]) and CGI promoter hypermethylation were
114 observed in tumor samples. EAC tumors harbored notably higher methylation levels in CGI

115   promoters than ESCC tumors, in line with TCGA results showing that gastrointestinal
116   adenocarcinoma had higher frequency of CGI hypermethylation than cancers from most other
117   tissues[26]. Interestingly, most NGEJ tissues showed higher CGI promoter methylation levels than
118   NESQ tissues, and usually even higher than ESCC tumor samples. Similar to EAC, BE samples
119   (a recognized precursor lesion of EAC) were reported to have a hypermethylation pattern at CGI
120   promoters[7]. Since our NGEJ tissues were pathologically confirmed as inflammatory tissues but
121   devoid of apparent BE, this result suggests that CGI hypermethylation may occur in inflamed GEJ.
122   Interestingly, CGI hypermethylation has been observed in long-term-cultured colon organoids and
123   cells upon prolonged exposure to cigarette smoke extract[27,28]. These data suggest that prolonged
124   extrinsic pressure may result in DNA methylation changes at CGIs. Repetitive elements,
125   especially from the LINE and LTR classes, lost DNA methylation in tumors compared with
126   nonmalignant tissues (**Fig. 1B**), which might be accompanied with the activation of repetitive
127   elements in tumor samples[21,29].

128   Considering the importance of PMDs in cancer biology[17,19,22,23], we sought to characterize
129   this epigenomic domain in depth. Computational tools have been developed for the identification
130   of PMDs, including MethPipe[30] and MethylSeekR[31]. However, they sometimes fail or return
131   unsatisfactory results for WGBS samples, either from tissues which have very slight
132   hypomethylation (see Sample 1 in **Fig. 1C**) or tumors with near-complete methylation loss (see
133   Sample 2 in **Fig. 1C**).

134   We recently used a deep learning neural network approach to establish universal sequence
135   features that are almost entirely predictive of CpG methylation loss or retention in PMD regions
136   of the human genome[32]. We hypothesized that utilizing sequence features associated with DNA
137   methylation loss and exploiting the variation patterns among different CpGs within PMDs could
138   improve the predictive models used in these tools (**Supplementary Fig. 1A-D;** see **Methods)**.
139   To this end, we developed a sequence-aware PMD calling method based on the Hidden Markov
140   Model (HMM) used in MethylSeekR (**Fig. 1C**; see **Methods**), which was termed Multi-model PMD
141   SeekR (MMSeekR). Importantly, using tumor samples from the Blueprint consortium, we showed
142   that MMSeekR outperformed both MethylSeekR and MethPipe (**Supplementary Fig. 1E-F**).
143   Indeed, MMSeekR successfully identified PMD fractions consistently across all samples (using
144   common PMDs as benchmark, top bar, **Supplementary Fig. 1E** and **Supplementary Table 2**).
145   MethylSeekR performed well in general, but was noisier and failed on several samples
146   (**Supplementary Fig. 1E**, red arrows). MethPipe performs poorly on samples with a small degree
147   of PMD methylation loss; indeed, this tool failed to identify PMD in almost half of these samples
148   (**Supplementary Fig. 1E** and **Supplementary Table 2**). PMD has been shown to exhibit cancer

149  type specificity[22,23], which can also be used to evaluate the performance of these methods.

150  Notably, MMSeekR almost completely separated different cancer types, while both MethylSeekR

151  and MethPipe produced much less clean separation (**Supplementary Fig. 1F**).

152  Encouraged by these results, we next applied MMSeekR to our esophageal samples.

153  Importantly, Principal Component Analysis (PCA) using PMDs identified by three different

154  methods again confirmed that MMSeekR outperformed MethylSeekR and MethPipe, completely

155  separating EAC and ESCC samples (**Fig. 1D** and **Supplementary Fig. 1G**). Interestingly,

156  nonmalignant samples clustered together with the corresponding cancer subtype. We also

157  provided exemplary PMDs that failed to be identified by either MethPipe (**Fig. 1E**) or MethylSeekR

158  (**Fig. 1F**).

**Characterization of shared and subtype-specific PMDs in esophageal samples**

160  We performed a genome-wide annotation of PMDs on a sample-by-sample basis (**Fig. 2A**).

161  Consistent with our earlier report[19] and the genome-wide analysis (**Fig. 1B**), PMDs showed a

162  slight decrease of DNA methylation in nonmalignant samples and lost methylation further in

163  tumors. Notably, PMDs exhibited high inter-sample heterogeneity in both their depth (i.e., DNA

164  methylation beta value) and breadth (i.e., genomic location). Indeed, the genome fraction covered

165  by PMDs varied markedly across samples, ranging from 24.3% to 63.4% (**Supplementary Fig.

166  2C**). We categorized these methylation domains into 4 groups based on the frequencies of their

167  occurrence in our cohort: shared PMDs, EAC-specific PMDs, ESCC-specific PMDs and shared

168  HMDs (**Fig. 2B** and **Supplementary Fig. 2D-E;** also see **Methods**). Interestingly, EAC-specific

169  PMDs covered significantly more of the genome than ESCC-specific PMDs (121.9Mb *vs.*

170  12.4Mb). To verify our results, we used solo-WCGW CpGs, which lose methylation faster than

171  other CpGs[19], to measure the average methylation loss within the 4 domain groups. In EAC

172  samples, shared PMDs and EAC-specific PMDs had lower methylation levels than the other two

173  groups, as expected (**Fig. 2C, left panel**). Reciprocally in ESCC samples, shared PMDs and

174  ESCC-specific PMDs had lower methylation levels (**Fig. 2C, right panel**). Independent cohorts

175  from either the TCGA (**Fig. 2D)** or other individual studies (**Supplementary Fig. 2F-G**) further

176  validated these subtype-specific patterns of DNA methylation loss. Since PMDs are associated

177  with the HiC B compartment[17,23], we next mathematically modeled the A/B chromatin

178  compartments for each esophageal cancer subtype using a method based on the HM450k

179  array[33]. Indeed, subtype-specific PMDs were enriched in B compartments in a subtype-specific

180  manner (**Fig. 2E**). By contrast, shared PMDs showed, as anticipated, no such specificity

181  (**Supplementary Fig. 2H**). PMD regions were also reported to have higher somatic mutation rate

182    compared with non-PMD regions in cancer[34,35]. We analyzed the whole-genome sequencing
183    (WGS) dataset from the OCCAMS (which has the largest number of EAC samples), finding a
184    significantly higher somatic mutation rate in EAC-specific PMDs than in either ESCC-specific
185    PMDs or HMDs (**Fig. 2F, left panel**). A reciprocal pattern was observed in the largest ESCC WGS
186    cohort (**Fig. 2F, right panel**).

187        At the transcription level, PMDs are reported to be less transcriptionally active than HMDs.
188    We confirmed that subtype-specific PMDs were associated with low levels of gene expression
189    specifically in the corresponding subtypes **(Fig. 3A-B)**. To explore the biological implication of
190    subtype-specific PMDs, we performed Cistrome-GO analysis using genes which were under-
191    expressed in the subtype-specific PMD regions, finding that biological processes characteristic
192    for the other subtype were enriched and repressed **(Fig. 3C-D)**. Specifically, pathways of
193    cornification, keratinocyte differentiation and epidermis development, which are central to
194    squamous cell differentiation and function, were enriched and inactive in EAC-specific PMDs **(Fig.**
195    **3C)**. For example, many keratinocyte-specific genes were clustered within EAC-specific PMDs
196    (**Fig. 3E**, **left panel**) and downregulated in EAC tumors (**Fig. 3F**, **upper panel**). On the other
197    hand, pathways important for gastrointestinal cell function, such as digestive system process,
198    intestinal absorption, lipid metabolic process and O−glycan processing, were enriched and
199    suppressed in ESCC-specific PMDs (**Fig. 3D**). The right panel of **Fig. 3E** shows as an example
200    that SLC2A2, which contributes to digestive system process and absorption, was located in
201    ESCC-specific PMDs and downregulated in ESCC samples (**Fig. 3F**, **lower panel**). These results
202    suggest that subtype-specific PMDs contain inactive genes which are associated with cell-type-
203    specific functions.

**H3K36me2 is inversely associated with PMDs in a cell-type-specific manner**

205        Both H3K36me2 and H3K36me3 were observed to recruit DNA methyltransferases
206    (DNMT3A[36] and DNMT3B[37], respectively) to maintain DNA methylation levels in large chromatin
207    domains. H3K36me3 is enriched in gene bodies of active transcripts, while H3K36me2 covers
208    larger multi-gene domains. Indeed, we have previously shown that the deposition of H3K36me3
209    is inversely associated with PMD distribution[19]. Here, we further hypothesized that H3K36me2
210    also contributed to maintaining DNA methylation levels, and the histone modification by this mark
211    might affect the genomic distribution of PMDs and HMDs. To test this, we performed H3K36me2
212    ChIP-seq in both EAC and ESCC cell lines. Indeed, shared HMDs (black line) showed high
213    H3K36me2 intensity in both cell types, while shared PMDs (purple line) exhibited the lowest
214    signals (**Fig. 4A**). EAC-specific PMDs (blue line) had low H3K36me2 levels in EAC cells but high

215  H3K36me2 levels in ESCC cells. The reciprocal pattern was observed in ESCC-specific PMDs

216  (red line). For example, H3K36me2 signals were undetectable in an EAC-specific PMD covering

217  the loci of *XR_945002.2* and *XR_945004.2* in EAC cells, but were strong in ESCC (**Fig. 4B, right**

218  **panel**). On the other hand, shared HMDs such as the one covering the *VSP8* gene were

219  decorated highly with H3K36me2 in both cell types (**Fig. 4B, left panel**).

220       To further verify these results, we interrogated public H3K36me2 ChIP-seq data from

221  HNSCC cell lines (squamous cancer highly similar to ESCC in terms of cell-of-origin and

222  epigenome). Indeed, a similar pattern of H3K36me2 distribution to ESCC was observed in Cal27

223  and Det562 HNSCC cells. Specifically, both shared PMDs and ESCC-specific PMDs harbored

224  low signals in HNSCC cell lines, while high H3K36me2 levels were found in HMDs and EAC-

225  specific PMDs (**Fig. 4C**). However, FaDu appeared to be an outlier, showing invariably high levels

226  across different regions (**Fig. 4C**), which warrants further investigation. Together, these results

227  demonstrate a prominent depletion of H3K36me2 mark in PMDs in a cell-type-specific manner,

228  which is likely owing to the finding that H3K36me2 promotes the maintenance of DNA methylation

229  by recruiting DNMT3A.


230  **Subtype-specific differentially methylated regions (DMRs) in esophageal cancer**

231       We next sought to investigate differentially methylated regions (DMRs) at small genomic

232  scales, given their direct roles in transcriptional regulation. However, our above results suggest

233  an overwhelming, global effect of PMD hypomethylation in tumor samples, which can strongly

234  affect the calling of focal DMRs. Indeed, PCA analysis of the most variable CpGs genome-wide

235  revealed that PC1, the most significant component, was clearly driven by methylation loss at

236  PMDs (**Supplementary Fig. 3A**).

237       To factor out the effect of PMD hypomethylation, we masked any PMD found within two thirds

238  of either EAC or ESCC samples (**Supplementary Fig. 3B**). We re-performed the PCA analysis,

239  finding that the two cancer subtypes were completely separated by PC1, which was the most

240  significant component and accounted for 42.2% of the total methylation variance (**Supplementary**

241  **Fig. 3C, left panel**). In addition, nonmalignant and tumor samples were separated along PC2,

242  and all NESQ samples were clustered closely together despite being generated from two different

243  cohorts. Notaly, this approach removed most correlation with the global methylation level

244  (**Supplementary Fig. 3C, right panel**). Thus, it is critical to remove the effects of global

245  hypomethylation when investigating cancer-associated methylation features outside PMDs.

246       We next identified DMRs between EAC and ESCC samples within the PMD-subtracted

247  genome described above (~46.5% of the genome). Under the cutoff of q value < 0.05 and absolute

248    delta methylation change > 0.2, a total of 7,734 DMRs were hypomethylated in EAC and 5,470 in

249    ESCC (**Fig. 5A**). As expected, hypomethylated DMRs (hypoDMRs) had low average methylation

250    levels in corresponding subtypes (**Supplementary Fig. 3D-E**). The majority of DMRs were about

251    1-2 kb long and located mostly in intronic and intergenic regions (**Fig. 5B**), similar to that of the

252    random background (**Supplementary Fig. 3F**). To investigate the epigenomic characteristics of

253    hypoDMRs, we systematically evaluated the chromatin accessibility at these regions, using the

254    ATAC-seq data from the TCGA[38] and H3K27ac ChIP-seq data from previous studies[39–42]. Relative

255    to random background regions, EAC hypoDMRs were accessible exclusively in EAC samples,

256    and ESCC hypoDMRs exclusively in ESCC samples (**Fig. 5C-D**). Additionally, EAC hypoDMRs

257    had high H3K27ac signals in 70% (5/7) of EAC cell lines (**Supplementary Fig. 3G**). A similar

258    observation was made in ESCC cell lines (**Supplementary Fig. 3H**). These data demonstrate

259    that hypoDMR regions are associated with accessible chromatin and active histone marks.

260        To explore the relevance of DMRs in gene transcription, we assigned each hypoDMR to the

261    closest genes annotated by HOMER[43,44], and performed correlational analyses using TCGA

262    transcriptomic data of esophageal cancers. Consistent with prior findings[43], about 30%

263    (3,986/13,204) of the DMRs were associated with differentially expressed genes

264    **(Supplementary Fig. 3I)**. Expectedly, an inverse correlation between DNA methylation and gene

265    expression accounted for the majority (~59%) of these associations, and these DMRs had a larger

266    overlap with promoter and enhancer regions **(Supplementary Fig. 3J)**. Importantly, functional

267    annotation using the Cistrome-GO method revealed that subtype hypoDMRs were enriched in

268    cell-type-specific biological processes. For example, lipid metabolic process, digestive system

269    process and O−glycan processing, which are housekeeping functions for gastrointestinal

270    columnar cells, were specifically enriched in EAC hypoDMRs **(Fig. 5E)**. On the other hand,

271    epidermis development, cornification and epithelial cell differentiation, which are unique to

272    squamous cells, were enriched in ESCC hypoDMRs **(Fig. 5F)**. These results indicate that a large

273    number of hypoDMRs regulate the transcription of cell-type-specific genes.

274        We next performed sequence motif enrichment analysis of hypoDMRs, which have

275    previously been associated with transcription-factor-binding sites[17,22,45]. A number of known

276    esophageal cell-specific transcription factors were identified, including GATA4/6, HNF4A/G,

277    HNF1B, ELF3, EHF in EAC[39,46,47] and TP63, SOX2 and MAFB in ESCC[41,48] (**Fig. 5G-H**). To

278    validate these results, we focused on the top-ranking transcription factors (GATA4 for EAC, TP63

279    for ESCC). Specifically, we performed WGBS in an EAC cell line (ESO26) where we previously

280    generated ChIP-seq data for GATA4 and H3K27ac. Indeed, GATA4 ChIP-seq peaks were

281    associated with high H3K27ac signal, DNA hypomethylation and GATA4 binding motif sequence

282    (**Fig. 5I**). Moreover, ~20% of GATA4 peaks overlapped with EAC hypoDMRs. In sharp contrast,

283    almost no GATA4 peaks were found in ESCC hypoDMRs (**Fig. 5I, left bars**). We similarly

284    performed WGBS on an ESCC cell line (TE5), and analyzed TP63 ChIP-Seq data that we

285    generated in the same sample. We noted consistent patterns and significant overlap with ESCC

286    hypoDMRs in this ESCC-specific transcription factor, and almost no overlap with EAC hypoDMRs

287    (**Fig. 5J**). These results demonstrate that subtype-specific DMRs are occupied by cell-type-

288    specific transcription factors and contribute to regulation of cell-type-specific functions.

**Identification of tumor-specific hypoDMRs**

290    To identify tumor-specific hypoDMRs from the above subtype-specific DMRs and to

291    investigate their role in cancer biology, we next performed a methylation comparison between

292    tumors and their corresponding nonmalignant samples for each hypoDMR. We found that 25.5%

293    (1,972/7,734) of EAC hypoDMRs (**Fig. 6A**) and 12.0% (654/5,470) of ESCC hypoDMRs

294    (**Supplementary Fig. 4A**) had significantly lower (FDR<0.05) methylation levels in tumors than

295    corresponding nonmalignant samples, which were referred to as "tumor specific hypoDMRs (ts-

296    hypoDMRs)", while the rest were referred to as "cell-type-specific DMRs (cts-hypoDMRs)". Ts-

297    hypoDRMs were distributed in both intergenic and intronic domains, similar to hypoDMRs overall

298    and the random background (**Fig. 6B** and **Supplementary Fig. 4B**). Between 18.0-21.4% of ts-

299    hypoDMRs were correlated with the expression of nearest genes (**Supplementary Fig. 4C-D**).

300    Importantly, ts-hypoDMRs were strongly enriched in cancer-related pathways such as cell cycle

301    progression (in both EAC and ESCC), and extracellular structure organization in ESCC (**Fig. 6C-

302    D**). These data suggest that ts-hypoDMRs are associated with genes which contribute to tumor-

303    specific functions.

304    The identification of ts-hypoDMRs and cts-hypoDMRs allowed us to further investigate

305    properties of tumor-specific regulatory regions *vs.* cell-type-specific regulatory regions. This is

306    particularly helpful for the epigenetic understanding of ESCC and EAC, which contain both tumor-

307    and cell-type-specific features. In addition, lineage-specific developmental factors have been

308    shown to promote malignant cell states[49,50], and thus it is important to distinguish their functional

309    contribution to normal development *vs.* cancer biology. To this end, we performed motif

310    enrichment analysis to identify transcription-factor-binding sites that were unique to either ts- or

311    cts-hypoDMRs, and integrated expression patterns of the corresponding transcription factors. For

312    EAC, this approach revealed cancer-upregulated transcription factors which favored binding ts-

313    hypoDMRs, including HNF4A, HNF4G, and FOSL1 (upper right corner of **Fig. 6E**). In comparison,

314    the lower left corner of **Fig. 6E** contained cancer-downregulated transcription factors which

315  preferred occupying cts-hypoDMRs, including GATA4/6 and FOXA, which are well-recognized for

316  their key roles in the development of gastrointestinal cell lineage[51,52]. The top factor for ts-

317  hypoDMR, HNF4A, had its binding motif in 46.6% ts-hypoDMRs but only 32.6% cts-hypoDMRs

318  (**Fig. 6F**). Indeed, ChIP-seq data of HNF4A in EAC cell lines (ESO26 and OE19) validated this

319  bias: HNF4A binding peaks overlapped with 14.2% ts-hypoDMRs but only 7.6% cts-hypoDMRs

320  (**Fig. 6G**). To identify factors that may facilitate recruitment of HNF4A specifically to hypoDMRs,

321  we performed enrichment analyses restricted within HNF4A-motif-containing hypoDMRs.

322  Interestingly, AP-1 motifs (such as JUN, FOSL1, FOSL2 and FOSB) were enriched in these

323  HNF4A$^+$ ts-hypoDMRs, while FOXA1/2 in cts-hypoDMRs (**Fig. 6H**). A parallel analysis was

324  performed in ESCC, which identified a number of tumor-specific factors, including RUNX1/3,

325  SOX2/4 and CEBPA/B (**Supplementary Fig. 4E**). This distinct pattern of co-occurring motifs

326  between ts- and cts-hypoDMRs in EAC is noteworthy, considering that AP-1 family transcription

327  factors contribute to EAC tumor development[53] while FOXA1/2 are required for normal

328  gastrointestinal cell development[52]. It is also notable that our analysis identified FOSL1 as an AP-

329  1 factor due to its high tumor expression **(Fig. 6E)**.

330  **PMDs and hypoDMRs exhibit strong cell-type-specific epigenomic features**

331  The above data identified both cell-type- and cancer-specific methylation differences in tumor

332  hypoDMRs, and we next asked whether tumor PMDs likewise harbor both of these two types of

333  methylation differences. In subtype-specific PMDs that were defined based on tumor methylomes

334  alone, nonmalignant tissues notably exhibited the same pattern of methylation changes as their

335  malignant counterparts (**Fig. 7A**). For example, EAC-specific PMDs had low methylation levels in

336  NGEJ but high in NESQ **(Fig. 7A, left)**, and a reciprocal pattern was found in ESCC-specific

337  PMDs **(Fig. 7A, right)**. Statistically, a large subset of subtype-specific PMDs (33.0% for EAC and

338  26.5% for ESCC) were already hypomethylated in their respective nonmalignant samples (**Fig.

339  7B**). The same analyses for hypoDMRs confirmed that more than 80% of subtype hypoDMRs

340  significantly decreased DNA methylation in their corresponding nonmalignant samples (**Fig. 7C-

341  D**). These data demonstrate that a substantial fraction of both subtype-specific PMDs and

342  hypoDMRs identified from tumor samples reflect methylation differences present in normal

343  counterparts. Nonetheless, while the genomic locations of PMDs are established in normal

344  samples, the degree of methylation loss is significantly higher in tumors (**Fig. 2C** and

345  **Supplementary Fig. 3D**).

346  To understand further PMDs and hypoDMRs in normal samples, we analyzed public single-

347  cell ATAC-seq data from 146,305 normal epithelial cells across 24 tissues (including esophageal

348   samples)[54], by measuring the chromatin accessibility of our subtype-specific PMDs or hypoDMRs.

349   This is premised on the fact that focal ATAC-seq peaks are almost always DNA demethylated[38],

350   and reduced ATAC-seq signals measured in large genomic windows reflect the Hi-C B

351   compartment which results in PMD hypomethylation[17,23]. The published single-cell unsupervised

352   clustering contains a cluster of esophageal squamous epithelial cells (red dots in **Fig. 7E, left**

353   **panel**), the recognized cell-of-origin for ESCC. With respect to EAC, although its cell-of-origin is

354   still under intense investigation, the epigenome is likely close to gastrointestinal epithelial cells

355   (blue dots **Fig. 7E, left panel**). Importantly, normal esophageal squamous cells showed the

356   lowest chromatin accessibility in ESCC-specific PMDs; reciprocally, normal gastrointestinal

357   epithelial cells had the lowest ATAC-Seq signals in EAC-specific PMDs (**Fig. 7E, middle panel;**

358   **quantified in Fig. 7F**). In addition, keratinocytes, which belong to squamous cell type, also had

359   low ATAC-Seq signals in ESCC-specific PMDs. In sharp contrast to subtype-specific PMDs, no

360   difference was observed in either shared PMDs or HMDs in this single-cell analysis

361   (**Supplementary Fig. 5C**). We performed the same analysis for hypoDMRs, finding that ESCC

362   hypoDMRs had the highest accessibility in squamous cells while EAC hypoDMRs were more

363   open in gastrointestinal epithelial cells (**Fig. 7E, right panel; quantified in Fig. 7G**). These single-

364   cell results confirmed that both PMDs and hypoDMRs have strong normal cell-type-specificity.

365   **Pan-cancer analysis of subtype-specific PMDs and hypoDMRs**

366   The above results also suggest that PMDs and hypoDMRs that we identified in ESCC and

367   EAC may be shared with other squamous and gastrointestinal adenocarcinomas, respectively.

368   To test this, we analyzed TCGA pan-cancer samples, since the TCGA multi-omic clustering

369   scheme[55] has identified the pan-gastrointestinal cluster (adenocarcinomas from esophagus,

370   stomach and colon, blue samples in **Fig. 8A**) and the pan-squamous cluster (squamous cancers

371   from esophagus, head and neck, lung, cervix and bladder, red samples in **Fig. 8A**). We first

372   measured the methylation changes between subtype-specific PMDs and hypoDMRs across all

373   33 cancer types (**Fig. 8B-E)**. Importantly, most pan-gastrointestinal tumors lost DNA methylation

374   in EAC-specific PMDs, while most pan-squamous tumors had reduced methylation in ESCC-

375   specific PMDs (**Fig. 8B and 8D**). Highly consistent results were observed in subtype hypoDMRs

376   (**Fig. 8C and 8E**). In contrast, no specific pattern was found in shared PMDs and HMDs

377   (**Supplementary Fig. 5D**), as anticipated.

378   We next analyzed the ATAC-seq data, which is available from a small subset of TCGA bulk

379   tumors[38], shown based on multi-omic clustering from ref[55] in **Fig 8F**. Importantly,  consistent with

380   the single-cell ATAC-Seq results from healthy tissues, pan-squamous cancers showed the lowest

381    chromatin accessibility in ESCC-specific PMDs and highest accessibility in ESCC hypoDMRs,

382    and the reciprocal results were obtained in pan-gastrointestinal cancers (**Fig. 8G-J**). Again, as

383    negative controls, shared PMDs and HMDs failed to generate this distinguishing epigenetic

384    pattern (**Supplementary Fig. 5E**).

385    These results prompted us to further investigate premalignant lesions, with the hypothesis

386    that these methylation changes are pre-established in normal cells and preserved during the

387    onset of neoplastic transformation. To address this, we interrogated public methylation data on

388    BE, a recognized precursor to EAC, from two different studies[7,8]. Importantly, the methylation

389    patterns of BE samples were highly comparable with EAC tumors, showing reduced methylation

390    levels in both EAC-specific PMDs and hypoDMRs in two different cohorts (**Fig. 8K-L**). Overall,

391    these data strongly suggest that epigenomic changes of PMDs and hypoDMRs occur in normal

392    cells and are maintained in cancer, which further loses methylation in PMDs and gains additional

393    DMRs. Moreover, these region-specific epigenomic regulations are shared across related cell

394    types.

## 395    Discussion

396    We generated one of the largest WGBS datasets in esophageal cancer to date, and here

397    we focused on the analyses of PMDs (large scale) and DMRs (small scale) and revealed novel

398    epigenomic properties of these regions. PMDs are megabase-long genomic regions with

399    decreased DNA methylation, coinciding with heterochromatic late-replicating domains and Hi-C

400    B domains[17]. PMDs reflect long-range chromatin organization that help orchestrate gene

401    expression programs and can influence replication timing and 3D genome organization[24,33,56–

402    58]. In addition, PMDs are associated with increased genomic instability and possibly activation of

403    transposable elements (TEs)[19,21]. Nevertheless, apart from these correlational observations, we

404    have only limited mechanistic understanding of the origin and regulation of cancer PMD.

405    Moreover, direct mechanisms linking PMDs to gene transcription remain to be established. Thus,

406    a deeper characterization of PMD is warranted, which first requires an accurate and sensitive

407    identification of these large domains from WGBS data. However, current PMD callers, including

408    MethylSeekR and MethPipe, either are insensitive for the identification of shallow PMDs, or fail to

409    call PMDs in tumor samples with extreme hypomethylation.

410    We have previously demonstrated that a local sequence context (solo-WCGW) is a strong

411    determinant of DNA methylation loss at CpGs[19]. Extending this finding, we recently performed

412    deep learning using the neural network method, and established universal sequence context

413    features influencing the hypomethylation of CpGs across the genome[32]. Here, we integrated this

414    sequence code into the MethylSeekR program and developed a novel multi-model PMD caller,

415    MMSeekR. Using both the Blueprint tumor WGBS dataset and our esophageal samples, we

416    demonstrated a superior performance of MMSeekR over other current tools. In order to facilitate

417    methodological development in the field of methylome investigation, we have made MMSeekR

418    available at Github as a free software package (https://github.com/yuanzi2/MMSeekR).

419         The degree of variation of PMD methylation levels (depth) and genomic distribution (breadth)

420    between cancer types was hitherto unclear. Here we observed strong heterogeneity at the PMD

421    methylation level across cancer samples, while nonmalignant samples harbored expectedly

422    shallow PMDs. Moreover, the genome fraction covered by PMDs varied profoundly among

423    different samples, ranging from 24.3% to 63.4%. We identified and characterized subtype-specific

424    PMDs, finding that they were associated with repressive transcription, B compartments and high

425    somatic mutation rate. We previously identified replication timing as a key determinant for

426    methylation loss in PMDs[19]. However, this does not account for the variation in PMD genomic

427    distribution across cell types. By investigation of the genome-wide occupancy of H3K36me2 in

428    different cell types, we noted that H3K36me2 deposition correlated positively with HMD

429    localization, while negatively with PMD in a cell-type-specific manner. Considering that

430    H3K36me2 is able to recruit DNMT3A to maintain the level of DNA methylation[36], these results

431    suggest that cell-type-specific deposition of H3K36me2 mark facilitates the maintenance of DNA

432    methylation, thereby dictating the genomic distribution of HMDs and PMDs.

433         At a smaller genomic scale, we identified over ten thousand hypoDMRs between the two

434    subtypes of esophageal cancer. Utilizing their matched nonmalignant samples, we further defined

435    cell-type- *vs.* cancer-specific hypoDMRs. Using motif sequence analysis combined with ChIP-

436    seq, we identified and validated candidate upstream regulators associated with either cell-type-

437    or cancer-specific hypoDMRs. This approach is important for understanding of the transcriptional

438    regulation during tumor development, particularly because increasing evidence has shown that

439    tumor-driving transcription factors are often lineage-specific developmental regulators functionally

440    co-opted to promote malignant cellular states[49,50]. For example, our top candidate, HNF4A, is

441    essential for the epithelial differentiation of the gastrointestinal tract. Consistently, we found that

442    a substantial subset of cell-type-specific hypoDMRs contained HNF4A-binding sequence; these

443    HNF4A[+] cell-type-specific hypoDMRs were also co-enriched for transcript factors indispensable

444    for normal gut development, such as FOXA1 (**Fig. 6H**). Importantly, compared with cell-type-

445    specific hypoDMRs, HNF4A-binding sequence was significantly more enriched in tumor-specific

446    hypoDMRs (**Fig. 6H)**. Moreover, instead of FOXA1, these HNF4A[+] tumor-specific hypoDMRs

447     were co-enriched for AP-1 factors, which are well-recognized for their function in promoting EAC

448     malignancy[53], similar to HNF4A itself[46,47]. Consistently, one of the AP1 factors, FOSL1, has highly

449     enriched binding sites in tumor-specific hypoDMRs as well as upregulated mRNA expression in

450     EAC tumors relative to NGEJ. Together, careful dissection of cell-type- and cancer-specific

451     hypoDMRs suggest that lineage master regulators control both normal and tumor cell

452     transcriptomes, likely by occupying different genomic regions through cooperating with different

453     transcriptional factor partners.

454     We further characterized the cell-type-specificity of PMDs and DMRs in normal cells. Starting

455     from esophageal samples, we found that a large fraction of methylation changes in both PMDs

456     and DMRs were already evident in normal samples. Pan-tissue single-cell ATAC-seq with

457     145,594 normal epithelial cells further showed that both PMDs and DMRs identified in esophageal

458     cancer had strong specificity that was evident in related cell types. This was also observed in pan-

459     cancer analyses of both methylation and ATAC-seq data from primary tumors, wherein cancers

460     originating from related cell types exhibited similar profiles of both PMDs and DMRs. Moreover,

461     by measuring cancer precursor lesions, we demonstrated that epigenomic changes of PMDs and

462     DMRs were preserved during the onset of neoplastic transformation. Nonetheless, PMDs in

463     normal samples were much shallower than tumors (**Fig. 2A** and **Fig 2C *vs*. Fig.7A**). Overall,

464     these data highlight the presence of cell-type-specific PMDs and DMRs in normal cell types, which

465     are preserved in malignant cells. To our knowledge, this is the first demonstration of the prominent

466     cell-type-specificity of PMDs across normal, precursor and malignant states. While prior studies

467     have revealed that DMRs contain tissue-specific regulatory regions, here we present a paradigm

468     for distinguishing cell-type- *vs.* cancer-specific regions, and use those to identify tumor-specific

469     regulatory mechanisms.

470     # Methods

471     **Cell culture**

472     Esophageal cancer cell lines, TE5, KYSE70, OE19 and ESO26, were grown in RPMI-1640

473     medium (Gibco, USA), supplemented with 10% FBS (Omega Scientific, USA) and 1% penicillin-

474     streptomycin (Thermo Scientific, USA). All cultures were maintained in a 37 °C incubator

475     supplemented with 5% CO2.

476     **Whole genome bisulfite sequencing (WGBS)**

477      WGBS of ESO26 or TE5 cells was performed at Novogene, Inc. Briefly, after DNA extraction

478    and quality control (QC), 3 ug DNA of ESO26 or TE5 cells spiked with 26 ng lambda DNA were

479    fragmented by sonication. The sonicated DNA was ligated with different cytosine-methylated

480    molecular barcodes. Next, bisulfite conversion was performed using EZ DNA Methylation-GoldTM

481    Kit (Zymo Research). PCR amplification with KAPA HiFi HotStart Uracil+Ready Mix (Kapa

482    Biosystems) was then applied to the DNA fragments. The clustering of index-coded DNA samples

483    were sequenced using the Illumina Hiseq 2500 platform.

484    **H3K36me2 chromatin immunoprecipitation sequencing (ChIP-Seq)**

485      Ten million esophageal cancer cells were harvested and transferred into 15 ml tubes,

486    followed by fixing with 4 ml of 1% paraformaldehyde for 10 min under room temperature. The

487    reaction was stopped by 2 ml of 250 mM of glycine. Cell samples were rinsed twice by 1X PBS

488    and lysed by 1 ml of 1X lysis/wash buffer (150 mM NaCl, 0.5 M EDTA pH 7.5, 1M Tris pH 7.5,

489    0.5% NP-40). Cell pellets were next resuspended using shearing buffer (1% SDS, 10 mM EDTA

490    pH 8.0, 50 nM Tris pH 8.0) followed by sonication using a Covaris sonicator. Subsequently, debris

491    was removed by centrifuge and supernatants were diluted five times with the buffer containing

492    0.01% SDS, 1% Triton X-100, 1.2 mM EDTA pH 8.0, 150 nM NaCl. 1 ug of the H3K36me2

493    antibody (Cell Signaling Technology, USA, Cat# 2901S) was added and incubated by rotation at

494    4℃ overnight. Protein G Dynabeads (Life Technologies, USA) were added the next morning and

495    incubated by rotation for an additional 4 hours. Dynabeads were next washed with 1X wash buffer

496    followed by cold TE buffer. DNAs were reverse crosslinked, purified, followed by library

497    preparation and deep sequencing using the Illumina HiSeq platform.

498    **Data sources**

499      DNA methylome of esophageal samples were obtained from our recent work[25], including

500    WGBS on 21 ESCC, 3 NESQ, 5 EAC, 7 GEJ tumors and 7 NGEJ tissues. We obtained additional

501    two NESQ samples from the ENCODE consortium to ensure statistical power. Considering the

502    indistinguishable clinical and molecular characteristics between EAC and GEJ tumors, in the

503    present study they were combined as the same subtype (referred to as EAC), which is a common

504    strategy in the field[3]. TCGA Pan-cancer DNA methylome derived from HM450k methylation array

505    was downloaded from GDC v16.0 by TCGAbiolinks package (version 2.13.6)[59]. Other DNA

506    methylation data from individual studies, including EAC EPIC array data from the Oesophageal

507    Cancer Clinical and Molecular Stratification (OCCAMS) consortium (EGAD00010001822)[9], EAC

508    and BE methylome from GSE72874[7] and GSE81334[8], along with ESCC tumor WGBS data
509    (GSE149608)[6], were analyzed for validation purposes in this study.

510        Other public datasets which were analyzed included: bulk ATAC-seq data of pan-cancer
511    samples from TCGA[38], single-cell ATAC-seq data across different adult human tissues
512    (GSE184462)[54], H3K27ac ChIP-seq in EAC samples (GSE132680)[39], EAC cell lines (ESO26,
513    FLO1, JH-EsoAd1, OACp4C, OE19, OE33, SKGT4 from GSE132680)[39] and ESCC cell lines
514    (KYSE140, KYSE70, TE5 from GSE106563[40]; KYSE150, KYSE180, KYSE200 from
515    GSE131490[41]; TE7 from GSE106433[42]), HNF4A ChIP-seq in OE19 (E-MTAB-6858)[46] and ESO26
516    cell lines (GSE132813)[47], GATA4 ChIP-seq in ESO26 cell line (GSE132813)[47] and TP63 ChIP-
517    seq in TE5 cell line (GSE148920)[41]. H3K36me2 bigwig files of wildtype (NSD1-WT) HNSCC cell
518    lines were downloaded from GSE149670[60]. Somatic mutation datasets were downloaded from
519    individual studies[9,61]. We also retrieved the transcriptomic data of esophageal cancer from the
520    TCGA consortium and GSE149609[6].

521        CGI promoters are annotated as regions ranging from 250 bp upstream to 500 bp
522    downstream of any TSSs overlapping with Takai CGIs[62]. Repetitive elements, including long
523    interspersed nuclear elements (LINE), short interspersed nuclear elements (SINE) and long
524    terminal repeats (LTR), were extracted from UCSC website (http://hgdownload.soe.ucsc.edu).
525    We downloaded the annotation of common PMDs  (defined as shared PMDs identified from 40
526    different cancer types)[19] as well as solo-WCGW from https://zwdzwd.github.io/pmd[19] and
527    ENCODE blacklist regions from https://github.com/Boyle-Lab/Blacklist/tree/master/lists[63]. All of
528    the annotations were converted to the hg38 version using the UCSC LiftOver script
529    (https://genome.ucsc.edu/cgi-bin/hgLiftOver). The human core transcription-factor-binding
530    sequences in the HOMOCOMO database (version 11) were used for motif annotation[64].

531    **DNA methylation data analysis**

532        For WGBS data, raw reads were mapped to the human genome (GRCh38) by Biscuit align
533    command (version 0.1.4, https://www.githubcom/zwdzwd/biscuit) with default settings. Mapped
534    reads were sorted by genome position, and duplicates were marked using Picard MarkDuplicates
535    tool (version 1.136, http://broadinstitute.github.io/picard/). Biscuit pileup and vcf2bed command
536    were then used to extract DNA methylation information. All CpG sites with a coverage >=3
537    informative reads and outside of the ENCODE blacklist regions were retained for downstream
538    analyses. For EPIC and HM450K array data, methylation of each probe was extracted using the
539    SeSAME package with noob and dyeBiasCorrTypeINorm function for background subtraction and
540    dye bias correction[65]. According to the annotation of Infinium DNA methylation arrays[66],

541     recommended general masking probes were removed. HM450K methylation data were used to

542     estimate the chromatin B compartments using minfi compartments function with

543     "resolution=100*1000, what = "OpenSea"" options[33].

**Development of a sequence-aware PMD calling method: Multi-model PMD SeekR (MMSeekR)**

546     We recently performed neural network-based machine learning to establish local DNA

547     sequence features of CpGs that were associated with global DNA methylation loss, and derived

548     a neural network (NN) score for each CpG across the human genome[32]. In order to exclude the

549     potential impact of high CpG density (such as CpG island), we reserved CpGs having 2 or fewer

550     neighboring CpGs within the 151 bp window centered on the reference CpG. We investigated the

551     correlation between NN scores and methylation in individual samples in non-overlapping 201-

552     CpG windows across the genome. As expected, due to the greater degree of methylation loss

553     within PMDs, there was a strong negative correlation between DNA methylation levels and NN

554     scores within windows in PMDs, in contrast to much more modest correlations within highly

555     methylated domains (HMD) windows (**Supplementary Fig. 1A**).

556     We next applied Pearson correlation coefficient (PCC) between our NN score and DNA

557     methylation, as well as the "alpha score" used in the MethylSeekR model, to 201-CpG windows

558     genome-wide. Compared with the NN score, the MethylSeekR alpha score is a very different

559     measurement, returning a high score if the distribution of methylation values is closer to a

560     unimodal beta distribution centered on 0.5 (typical of PMDs) than it is to a bimodal methylation

561     value distribution close to 0 and 1 (typical of HMDs). Specifically, we applied a Hidden Markov

562     Model (HMM) segmentation (as in MethylSeekR) to each model independently, and found that

563     both the PCC and MethylSeekR alpha score showed bimodal distributions for the testing sample

564     (**Supplementary Fig. 1B-C**). We hypothesized that since the PCC and the alpha score were very

565     different models, combining them might improve the performance of PMD calling

566     (**Supplementary Fig. 1D**). Thus we developed a "2-dimensional (2D)" model accordingly (**Fig.**

567     **1C**). This 2D model performed comparably well or better than either MethylSeekR or MethPipe in

568     most cases, returning results consistently and highly overlapping with common PMDs

569     (**Supplementary Table 2**).

570     While the 2D model generally performed well, we did note that it failed in a few samples with

571     extreme methylation loss. Interestingly, these failed cases universally showed PMD methylation

572     values very close to 0, which would be expected to violate the assumptions of both the PCC model

573     and alpha model due to lack of variance within PMDs (**Fig. 1C right part**). We thus postulated

574 the raw methylation values (transformed to an M-value to disperse scores close to 0 and 1) might

575 provide additional predictive power in certain samples with extreme methylation loss, and we

576 developed a 3D model accordingly by adding the M-value model to the 2D model. In order to

577 decide whether the 2D or 3D model should be applied for any given sample, we first measured

578 the methylation values of all CpGs with 2 or fewer neighboring CpGs within a 151bp window,

579 which excludes most CpG islands, and contains a set of CpGs that are strongly associated with

580 PMD hypomethylation[19]. If the bottom 10th percentile of these CpGs had a methylation value

581 below 0.025, the 3D model was selected, otherwise, the 2D model was selected. This was based

582 on the observation that the majority of samples with extreme methylation loss failed under both

583 the MethylSeekR and MMSeekR 2D model (**Fig. 1C**).

**Application of MMSeekR to WGBS data**

585 MMSeekR was applied to call PMDs in each WGBS sample. Before PMD calling, CpG sites

586 with coverage of fewer than 5 informative reads were excluded. Then ENCODE blacklist regions

587 were subtracted from the resulting PMDs. Within each esophageal cancer subtype, PMDs

588 generated from each sample were integrated using bedtools multiinter function (version 2.27.1,

589 https://bedtools.readthedocs.io/en/latest/). The common PMD set for each subtype contained

590 those occurring in at least two-thirds of samples from that subtype. We further defined subtype-

591 specific PMDs as those common PMDs from one subtype that were detected in fewer than one-

592 third of samples in the other subtype. Meanwhile, PMDs that were in both the common EAC set

593 and the common ESCC set were denoted as shared PMDs. Regions that were PMDs in <1/3

594 samples of both subtypes were denoted as shared HMDs.

**Identification and characterization of DMRs**

596 Regions belonging to either the common ESCC or common EAC PMD sets were masked

597 out from the DMR analysis. The Dmrseq package (version 1.10.0)[67] was used to identify DMRs

598 between ESCC and EAC tumors with the following parameters: cutoff =0.1, bpSpan=1000,

599 minInSpan=30, maxPerms=500. Since the coverage information of each CpG site is required by

600 dmrseq for statistical inference, here we included all CpG sites with >= 3 informative reads.

601 Regions with q value < 0.05 and absolute delta methylation change > 0.2 were identified as DMRs.

602 For hypomethylated DMRs (hypoDMRs) from each cancer subtype, we further performed one-

603 tailed t-tests comparing the mean methylation within the DMR in nonmalignant *vs.* tumor samples,

604 and those with FDR<0.1 were considered as tumor-specific (ts)-hypoDMRs. Both hypoDMRs and

605 ts-hypoDMRs were annotated using HOMER annotatePeaks.pl script (version 4.9.1)[44].

**Calculation of mean DNA methylation levels**

CpG sites with a coverage of at least 5 informative reads were used for this calculation. Average methylation levels of CpG sites across the genome (global level), within CGI promoters, commonPMDs, SINE, LINE and LTR in each sample were calculated independently. Besides, we obtained the mean methylation of CpG sites in non-PMD regions. For genome/domain-wide visualization, the average methylation of 10-kb consecutive non-overlapping tiles was shown. To calculate the mean methylation levels within shared PMDs/HMDs, EAC-specific PMDs and ESCC-specific PMDs, solo-WCGW CpG sites/probes were used.

**Principal component analysis of WGBS data**

PMDs were identified by either MethPipe, MethylSeekR or MMseekR (**Fig. 1D**). The whole genome was split into 30-kb consecutive but non-overlapping tiles. For each tile, the ratio overlapping with any PMD was calculated for each caller. The top 5,000 most variable 30-kb tiles from each PMD caller were used in Principal component analysis (PCA). In **Supplementary Fig. 3A** and **3B**, CpG sites with at least 7 reads across all esophageal samples were used. Then the top 8,000 most variable CpG sites were selected for PCA using the R prcomp function. PCA was performed before and after masking the combined common PMDs from EAC and ESCC, and generated the point plots by ggplot2 package (version 3.1.0).

**RNA-seq data analysis**

According to the raw read counts obtained from the TCGA, we identified significant upregulated genes by DESeq2 package (version 1.22.2) with adjusted p-value < 0.05, fold change > 1.5 and mean FPKM >1 in the corresponding sample groups[68]. For expression datasets of nonmalignant squamous and ESCC tissues, raw reads were aligned to GRCh38 using HISAT2 (version 2.0.4)[69] and quantified by htseq-count program (version 0.11.2) at default setting. Significant upregulated genes were identified using the same method as for the TCGA datasets.

**ChIP-seq data analysis**

Raw reads were mapped to GRCh38 (ENSEMBL release 84) using BWA mem program (version 0.7.15) with the default options[70]. Then the mapped reads were sorted using SAMtools program (version 1.3.1)[71], followed by removing PCR duplicates and blacklist regions by Picard MarkDuplicates tool and bedtools (version 2.27.1). MACS2 (Model-Based Analysis of ChIP-Seq, version 2.1.2) were applied to call peaks with the default setting for transcription factors, "-q 0.01– extsize = 146 –nomodel" options for H3K27ac and "--broad -p 0.01 --extsize=146 --nomodel" for

637 H3K36me2[72]. Bigwig files were generated by deepTools bamCompare function (version 3.1.3)

638 with "--operation subtract --normalizeUsing CPM --extendReads 146 --binSize 20" parameters[73].

639 Average signals of shared PMDs/HMDs, EAC-only PMDs and ESCC-only PMDs in each

640 H3K27ac or H3K36me2 ChIP-seq sample were extracted from bigwig files using deepTools

641 computeMatrix function with "scale-regions" option.

**ATAC-seq data analysis**

642

643 For bulk pan-cancer ATAC-seq data obtained from the TCGA project, the average

644 accessibility of regions/domains was extracted from the available bigwig files using deepTools

645 computeMatrix function[38]. To avoid the influence of scaling factors across different samples and

646 batches, the mean accessibility across the whole genome in each sample was calculated and

647 used for normalization. For single cell ATAC-seq data, based on the clustering and annotation

648 results from the publication[54], only epithelial cell types were used for further analysis. Similarly,

649 the average accessibility of regions/domains was derived for each cell in each sample and

650 normalized by the mean signal across the whole genome.

**DMR motif enrichment analysis**

651

652 For each hypoDMR or ts-hypoDMR, we randomly sampled 10 regions with the same size

653 and number of CpGs to define the background set. Then motif searching of both DMRs and

654 background regions was performed using HOMER annotatePeaks.pl with "-noann -m

655 HOCOMOCOv11_core_HUMAN_mono_homer_format_0.0001.motif" parameters[44]. The

656 ELMER method was next applied to identify potential transcription-factor-binding sequences and

657 the top 15 transcription factors with q-value < 0.05 and FPKM > 5 in the corresponding cancer

658 subtype were reserved for further analysis[74].

**Pathway enrichment analysis**

659

660 We performed the pathway (Biological Process) enrichment analysis by Cistrome-GO[75] using

661 candidate regions with methylation changes and differential expression analysis results. For

662 hypoDMR analysis, subtype-specific DMRs and upregulated genes in the corresponding tumors

663 were used as input data. For subtype-specific PMDs, the input data contained PMD regions and

664 downregulated genes in the corresponding tumors. The top 15 enriched pathways with q value <

665 0.05 were shown.

**Code Availability**

666

667    Source code for MMSeekR is available at https://github.com/yuanzi2/MMSeekR. Source

668    code for WGBS data analysis and figure reproduction is in

669    https://github.com/yuanzi2/ESCA_WGBS_analysis.

670    **Data Availability**

671    WGBS data and ChIP-seq data for H3K36me2 in EAC and ESCC cell lines were available

672    at GSE210220.

673

# Reference

675    1.    Sung, H. *et al.* Global Cancer Statistics 2020: GLOBOCAN Estimates of Incidence and

676          Mortality Worldwide for 36 Cancers in 185 Countries. *CA Cancer J. Clin.* **71**, 209–249

677          (2021).

678    2.    Siegel, R. L., Miller, K. D., Fuchs, H. E. & Jemal, A. Cancer Statistics, 2021. *CA Cancer J.*

679          *Clin.* **71**, 7–33 (2021).

680    3.    Cancer Genome Atlas Research Network *et al.* Integrated genomic characterization of

681          oesophageal carcinoma. *Nature* **541**, 169–175 (2017).

682    4.    Talukdar, F. R. *et al.* Genome-Wide DNA Methylation Profiling of Esophageal Squamous

683          Cell Carcinoma from Global High-Incidence Regions Identifies Crucial Genes and Potential

684          Cancer Markers. *Cancer Res.* **81**, 2612–2624 (2021).

685    5.    Teng, H. *et al.* Inter- and intratumor DNA methylation heterogeneity associated with lymph

686          node metastasis and prognosis of esophageal squamous cell carcinoma. *Theranostics* **10**,

687          3035–3048 (2020).

688    6.    Cao, W. *et al.* Multi-faceted epigenetic dysregulation of gene expression promotes

689          esophageal squamous cell carcinoma. *Nat. Commun.* **11**, 3675 (2020).

690    7.    Krause, L. *et al.* Identification of the CIMP-like subtype and aberrant methylation of

691          members of the chromosomal segregation and spindle assembly pathways in esophageal

692      adenocarcinoma. *Carcinogenesis* **37**, 356–365 (2016).

693    8.   Yu, M. *et al.* Subtypes of Barrett's oesophagus and oesophageal adenocarcinoma based

694      on genome-wide methylation analysis. *Gut* **68**, 389–399 (2019).

695    9.   Jammula, S. *et al.* Identification of Subtypes of Barrett's Esophagus and Esophageal

696      Adenocarcinoma Based on DNA Methylation Profiles and Integration of Transcriptome and

697      Genome Data. *Gastroenterology* **158**, 1682–1697.e1 (2020).

698   10.   Angeloni, A. & Bogdanovic, O. Enhancer DNA methylation: implications for gene regulation.

699      *Essays Biochem.* **63**, 707–715 (2019).

700   11.   Lister, R. *et al.* Human DNA methylomes at base resolution show widespread epigenomic

701      differences. *Nature* **462**, 315–322 (2009).

702   12.   Slotkin, R. K., Keith Slotkin, R. & Martienssen, R. Transposable elements and the

703      epigenetic regulation of the genome. *Nature Reviews Genetics* vol. 8 272–285 (2007).

704   13.   Baylin, S. B. & Jones, P. A. Epigenetic Determinants of Cancer. *Cold Spring Harb.*

705      *Perspect. Biol.* **8**, (2016).

706   14.   Luo, C., Hajkova, P. & Ecker, J. R. Dynamic DNA methylation: In the right place at the right

707      time. *Science* **361**, 1336–1340 (2018).

708   15.   Karlow, J. A., Miao, B., Xing, X., Wang, T. & Zhang, B. Common DNA methylation

709      dynamics in endometriod adenocarcinoma and glioblastoma suggest universal epigenomic

710      alterations in tumorigenesis. *Commun Biol* **4**, 607 (2021).

711   16.   Hansen, K. D. *et al.* Increased methylation variation in epigenetic domains across cancer

712      types. *Nat. Genet.* **43**, 768–775 (2011).

713   17.   Berman, B. P. *et al.* Regions of focal DNA hypermethylation and long-range

714      hypomethylation in colorectal cancer coincide with nuclear lamina-associated domains. *Nat.*

715      *Genet.* **44**, 40–46 (2011).

716   18.   Hon, G. C. *et al.* Global DNA hypomethylation coupled to repressive chromatin domain

717      formation and gene silencing in breast cancer. *Genome Res.* **22**, 246–258 (2012).

718    19. Zhou, W. *et al.* DNA methylation loss in late-replicating domains is linked to mitotic cell

719        division. *Nat. Genet.* **50**, 591–602 (2018).

720    20. Duran-Ferrer, M. *et al.* The proliferative history shapes the DNA methylome of B-cell tumors

721        and predicts clinical outcome. *Nat Cancer* **1**, 1066–1081 (2020).

722    21. Hur, K. *et al.* Hypomethylation of long interspersed nuclear element-1 (LINE-1) leads to

723        activation of proto-oncogenes in human colorectal cancer metastasis. *Gut* **63**, 635–646

724        (2014).

725    22. Hovestadt, V. *et al.* Decoding the regulatory landscape of medulloblastoma using DNA

726        methylation sequencing. *Nature* **510**, 537–541 (2014).

727    23. Brinkman, A. B. *et al.* Partially methylated domains are hypervariable in breast cancer and

728        fuel widespread CpG island hypermethylation. *Nat. Commun.* **10**, 1749 (2019).

729    24. Salhab, A. *et al.* A comprehensive analysis of 195 DNA methylomes reveals shared and

730        cell-specific features of partially methylated domains. *Genome Biol.* **19**, 150 (2018).

731    25. Pan, F. *et al.* Characterization of epigenetic alterations in esophageal cancer by whole-

732        genome bisulfite sequencing. *bioRxiv* 2021.12.05.471340 (2021)

733        doi:10.1101/2021.12.05.471340.

734    26. Liu, Y. *et al.* Comparative Molecular Analysis of Gastrointestinal Adenocarcinomas. *Cancer*

735        *Cell* **33**, 721–735.e8 (2018).

736    27. Tao, Y. *et al.* Aging-like Spontaneous Epigenetic Silencing Facilitates Wnt Activation,

737        Stemness, and Braf-Induced Tumorigenesis. *Cancer Cell* **35**, 315–328.e6 (2019).

738    28. Vaz, M. *et al.* Chronic Cigarette Smoke-Induced Epigenomic Changes Precede

739        Sensitization of Bronchial Epithelial Cells to Single-Step Transformation by KRAS

740        Mutations. *Cancer Cell* **32**, 360–376.e6 (2017).

741    29. Ehrlich, M. & Lacey, M. DNA hypomethylation and hemimethylation in cancer. *Adv. Exp.*

742        *Med. Biol.* **754**, 31–56 (2013).

743    30. Decato, B. E. *et al.* Characterization of universal features of partially methylated domains

744      across tissues and species. *Epigenetics Chromatin* **13**, 39 (2020).

745   31. Burger, L., Gaidatzis, D., Schübeler, D. & Stadler, M. B. Identification of active regulatory

746      regions from DNA methylation data. *Nucleic Acids Res.* **41**, e155 (2013).

747   32. Bar, D. *et al.* A local sequence signature defines a subset of heterochromatin-associated

748      CpGs with minimal loss of methylation in healthy tissues but extensive loss in cancer.

749      *bioRxiv* 2022.08.16.504069 (2022) doi:10.1101/2022.08.16.504069.

750   33. Fortin, J.-P. & Hansen, K. D. Reconstructing A/B compartments as revealed by Hi-C using

751      long-range correlations in epigenetic data. *Genome Biol.* **16**, 180 (2015).

752   34. Schuster-Böckler, B. & Lehner, B. Chromatin organization is a major influence on regional

753      mutation rates in human cancer cells. *Nature* **488**, 504–507 (2012).

754   35. Lawrence, M. S. *et al.* Mutational heterogeneity in cancer and the search for new cancer-

755      associated genes. *Nature* **499**, 214–218 (2013).

756   36. Weinberg, D. N. *et al.* The histone mark H3K36me2 recruits DNMT3A and shapes the

757      intergenic DNA methylation landscape. *Nature* **573**, 281–286 (2019).

758   37. Neri, F. *et al.* Intragenic DNA methylation prevents spurious transcription initiation. *Nature*

759      **543**, 72–77 (2017).

760   38. Corces, M. R. *et al.* The chromatin accessibility landscape of primary human cancers.

761      *Science* **362**, (2018).

762   39. Chen, L. *et al.* Master transcription factors form interconnected circuitry and orchestrate

763      transcriptional networks in oesophageal adenocarcinoma. *Gut* **69**, 630–640 (2020).

764   40. Jiang, Y. *et al.* Co-activation of super-enhancer-driven CCAT1 by TP63 and SOX2

765      promotes squamous cancer progression. *Nat. Commun.* **9**, 3619 (2018).

766   41. Jiang, Y.-Y. *et al.* TP63, SOX2, and KLF5 Establish a Core Regulatory Circuitry That

767      Controls Epigenetic and Transcription Patterns in Esophageal Squamous Cell Carcinoma

768      Cell Lines. *Gastroenterology* **159**, 1311–1327.e19 (2020).

769   42. Xie, J.-J. *et al.* Super-Enhancer-Driven Long Non-Coding RNA LINC01503, Regulated by

770        TP63, Is Over-Expressed and Oncogenic in Squamous Cell Carcinoma. *Gastroenterology*

771        **154**, 2137–2151.e1 (2018).

772   43. Espinet, E. *et al.* Aggressive PDACs Show Hypomethylation of Repetitive Elements and the

773        Execution of an Intrinsic IFN Program Linked to a Ductal Cell of Origin. *Cancer Discov.* **11**,

774        638–659 (2021).

775   44. Heinz, S. *et al.* Simple combinations of lineage-determining transcription factors prime cis-

776        regulatory elements required for macrophage and B cell identities. *Mol. Cell* **38**, 576–589

777        (2010).

778   45. Aran, D., Sabato, S. & Hellman, A. DNA methylation of distal regulatory sites characterizes

779        dysregulation of cancer genes. *Genome Biol.* **14**, R21 (2013).

780   46. Rogerson, C. *et al.* Identification of a primitive intestinal transcription factor network shared

781        between esophageal adenocarcinoma and its precancerous precursor state. *Genome Res.*

782        **29**, 723–736 (2019).

783   47. Pan, J. *et al.* Lineage-Specific Epigenomic and Genomic Activation of Oncogene HNF4A

784        Promotes Gastrointestinal Adenocarcinomas. *Cancer Res.* **80**, 2722–2736 (2020).

785   48. Lopez-Pajares, V. *et al.* A LncRNA-MAF:MAFB transcription factor network regulates

786        epidermal differentiation. *Dev. Cell* **32**, 693–706 (2015).

787   49. Reddy, J. *et al.* Predicting master transcription factors from pan-cancer expression data.

788        *Sci Adv* **7**, eabf6123 (2021).

789   50. Sanda, T. *et al.* Core transcriptional regulatory circuit controlled by the TAL1 complex in

790        human T cell acute lymphoblastic leukemia. *Cancer Cell* **22**, 209–221 (2012).

791   51. Walker, E. M., Thompson, C. A. & Battle, M. A. GATA4 and GATA6 regulate intestinal

792        epithelial cytodifferentiation during development. *Dev. Biol.* **392**, 283–294 (2014).

793   52. Ye, D. Z. & Kaestner, K. H. Foxa1 and Foxa2 control the differentiation of goblet and

794        enteroendocrine L- and D-cells in mice. *Gastroenterology* **137**, 2052–2062 (2009).

795   53. Britton, E. *et al.* Open chromatin profiling identifies AP1 as a transcriptional regulator in

796      oesophageal adenocarcinoma. *PLoS Genet.* **13**, e1006879 (2017).

797   54. Zhang, K. *et al.* A single-cell atlas of chromatin accessibility in the human genome. *Cell*

798      **184**, 5985–6001.e19 (2021).

799   55. Hoadley, K. A. *et al.* Cell-of-Origin Patterns Dominate the Molecular Classification of 10,000

800      Tumors from 33 Types of Cancer. *Cell* **173**, 291–304.e6 (2018).

801   56. Nothjunge, S. *et al.* DNA methylation signatures follow preformed chromatin compartments

802      in cardiac myocytes. *Nat. Commun.* **8**, 1667 (2017).

803   57. Du, Q. *et al.* DNA methylation is required to maintain both DNA replication timing precision

804      and 3D genome organization integrity. *Cell Rep.* **36**, 109722 (2021).

805   58. Johnstone, S. E. *et al.* Large-Scale Topological Changes Restrain Malignant Progression in

806      Colorectal Cancer. *Cell* **182**, 1474–1489.e23 (2020).

807   59. Mounir, M. *et al.* New functionalities in the TCGAbiolinks package for the study and

808      integration of cancer data from GDC and GTEx. *PLoS Comput. Biol.* **15**, e1006701 (2019).

809   60. Farhangdoost, N. *et al.* Chromatin dysregulation associated with NSD1 mutation in head

810      and neck squamous cell carcinoma. *Cell Rep.* **34**, 108769 (2021).

811   61. Cui, Y. *et al.* Whole-genome sequencing of 508 patients identifies key molecular features

812      associated with poor prognosis in esophageal squamous cell carcinoma. *Cell Res.* **30**,

813      902–913 (2020).

814   62. Takai, D. & Jones, P. A. Comprehensive analysis of CpG islands in human chromosomes

815      21 and 22. *Proc. Natl. Acad. Sci. U. S. A.* **99**, 3740–3745 (2002).

816   63. Amemiya, H. M., Kundaje, A. & Boyle, A. P. The ENCODE Blacklist: Identification of

817      Problematic Regions of the Genome. *Sci. Rep.* **9**, 9354 (2019).

818   64. Kulakovskiy, I. V. *et al.* HOCOMOCO: towards a complete collection of transcription factor

819      binding models for human and mouse via large-scale ChIP-Seq analysis. *Nucleic Acids*

820      *Res.* **46**, D252–D259 (2018).

821   65. Zhou, W., Triche, T. J., Jr, Laird, P. W. & Shen, H. SeSAMe: reducing artifactual detection

822    of DNA methylation by Infinium BeadChips in genomic deletions. *Nucleic Acids Res.* **46**,

823    e123 (2018).

824    66. Zhou, W., Laird, P. W. & Shen, H. Comprehensive characterization, annotation and

825    innovative use of Infinium DNA methylation BeadChip probes. *Nucleic Acids Res.* **45**, e22

826    (2017).

827    67. Korthauer, K., Chakraborty, S., Benjamini, Y. & Irizarry, R. A. Detection and accurate false

828    discovery rate control of differentially methylated regions from whole genome bisulfite

829    sequencing. *Biostatistics* **20**, 367–383 (2019).

830    68. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for

831    RNA-seq data with DESeq2. *Genome Biol.* **15**, 550 (2014).

832    69. Kim, D., Langmead, B. & Salzberg, S. L. HISAT: a fast spliced aligner with low memory

833    requirements. *Nat. Methods* **12**, 357–360 (2015).

834    70. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.

835    (2013).

836    71. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–

837    2079 (2009).

838    72. Zhang, Y. *et al.* Model-based analysis of ChIP-Seq (MACS). *Genome Biol.* **9**, R137 (2008).

839    73. Ramírez, F. *et al.* deepTools2: a next generation web server for deep-sequencing data

840    analysis. *Nucleic Acids Res.* **44**, W160–5 (2016).

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

842    networks from DNA methylation and transcriptome profiles. *Bioinformatics* **35**, 1974–1977

843    (2019).

844    75. Li, S. *et al.* Cistrome-GO: a web server for functional enrichment analysis of transcription

845    factor ChIP-seq peaks. *Nucleic Acids Res.* **47**, W206–W211 (2019).

846    76. Irizarry, R. A. *et al.* Genome-wide methylation analysis of human colon cancer reveals

847    similar hypo- and hypermethylation at conserved tissue-specific CpG island shores. *Nat.*

848      *Genet.* **41**, 178 (2009).

849   77. Silva, T. C. *et al.* ELMER v.2: An R/Bioconductor package to reconstruct gene regulatory

850      networks from DNA methylation and transcriptome profiles. doi:10.1101/148726.

## Acknowledgement

## Author contribution

860 D.-C.L. and B.P.B. conceived and devised the study. D.-C.L., B.P.B., Y.Y.Z., and B.Z. designed

861 experiments and analyses. Y.Y.Z and B.P.B. performed bioinformatics and statistical analysis.

862 B.Z performed the experiments. Y.Y.Z., B.P.B., and D.-C.L. analyzed the data. B.P.B., D.-C.L.

863 supervised the research. A.S.H, U.K.S, L.Y.X, E.M.L and H.P.K. contributed the data and

864 materials. Y.Y.Z., and D.-C.L. wrote the manuscript with input from B.P.B. The last two authors

865 (D.-C.L. and B.P.B.) are co-senior authors who jointly supervised the work, and they have the

866 right to list their names last in their CV.

## Supplementary information

868 **Supplementary Figures.docx**

869 **Supplementary Table 1.** WGBS data sets used in the current study.

870 **Supplementary Table 2.** The percent of PMDs identified by three different callers overlapping

871 with common PMDs or HMDs in each tumor sample from the Blueprint consortium or esophageal

872 tissue.

873

874

875 **Figure 1. Identification of PMDs in esophageal samples by a sequence-aware multi-model**

876 **PMD caller (MMSeekR). (A)** A graphic model of the present study design. **(B)** Dot plots showing

877 average methylation levels for all CpGs across the whole genome, CpGs within CGI promoters,

878 common PMDs, SINE, LINE and LTR in different samples. The annotations from Takai et al[62].

879 were used for CGI methylation quantification. **(C)** Development of a new PMD caller. The

880 MethylSeekR α score measures the distribution of methylation levels in sliding windows with 201

881 consecutive CpGs across the genome. α score < 1 corresponds to a polarized distribution towards

882 a high or low methylation level (that is, HMDs), while α score >=1 corresponds to the distribution

883 towards intermediate methylation levels (that is, PMDs). PCC shows the correlation between the

884 predicted hypomethylation score based on a NN model, and the actual methylation level. A strong

885 negative correlation indicates regions favoring PMDs, while weak/null correlation favors HMDs.

886 **(D)** PCA analysis of 45 esophageal samples using the top 5,000 most variable 30-kb tiles for the

887 three PMD callers. **(E-F)** Representative windows showing PMDs successfully identified by

888 MMSeekR but failed to be detected by either MethPipe **(E)** or MethylSeekR **(F)**.
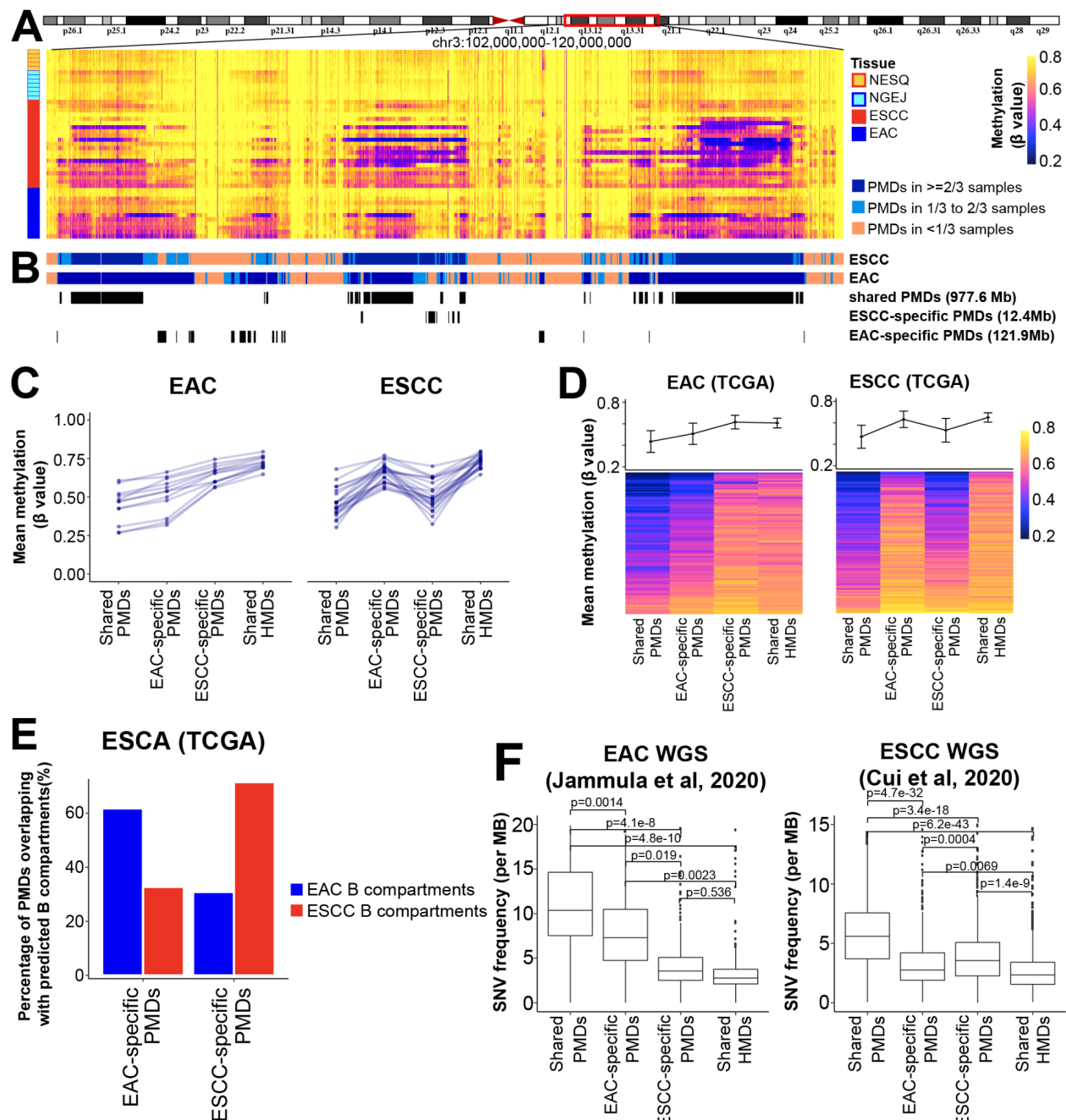
889

890

**Figure 2. Characterization of shared and subtype-specific PMDs. (A)** A representative window of DNA methylation profiles from 45 esophageal samples. Average methylation values are shown in consecutive and non-overlapping 10-kb tiles. CGI regions were masked using the annotation from Irizarry et al[76]. **(B)** Different PMD categories were identified based on the frequency and overlap between the two esophageal cancer types. **(C)** Line plots showing average methylation levels for different PMD categories in esophageal tumors, where each line represents one sample. **(D)** Similar line plot patterns were observed using TCGA methylation datasets, showing the mean and standard deviation across samples. Each row in the heatmap below shows

900     an individual sample. **(E)** Bar plots showing the percentage of WGBS PMDs overlapping with

901     chromatin B compartments, which were predicted using TCGA methylation datasets and analyzed

902     by minfi package. **(F)** Somatic mutation rates based on WGS in the indicated studies, calculated

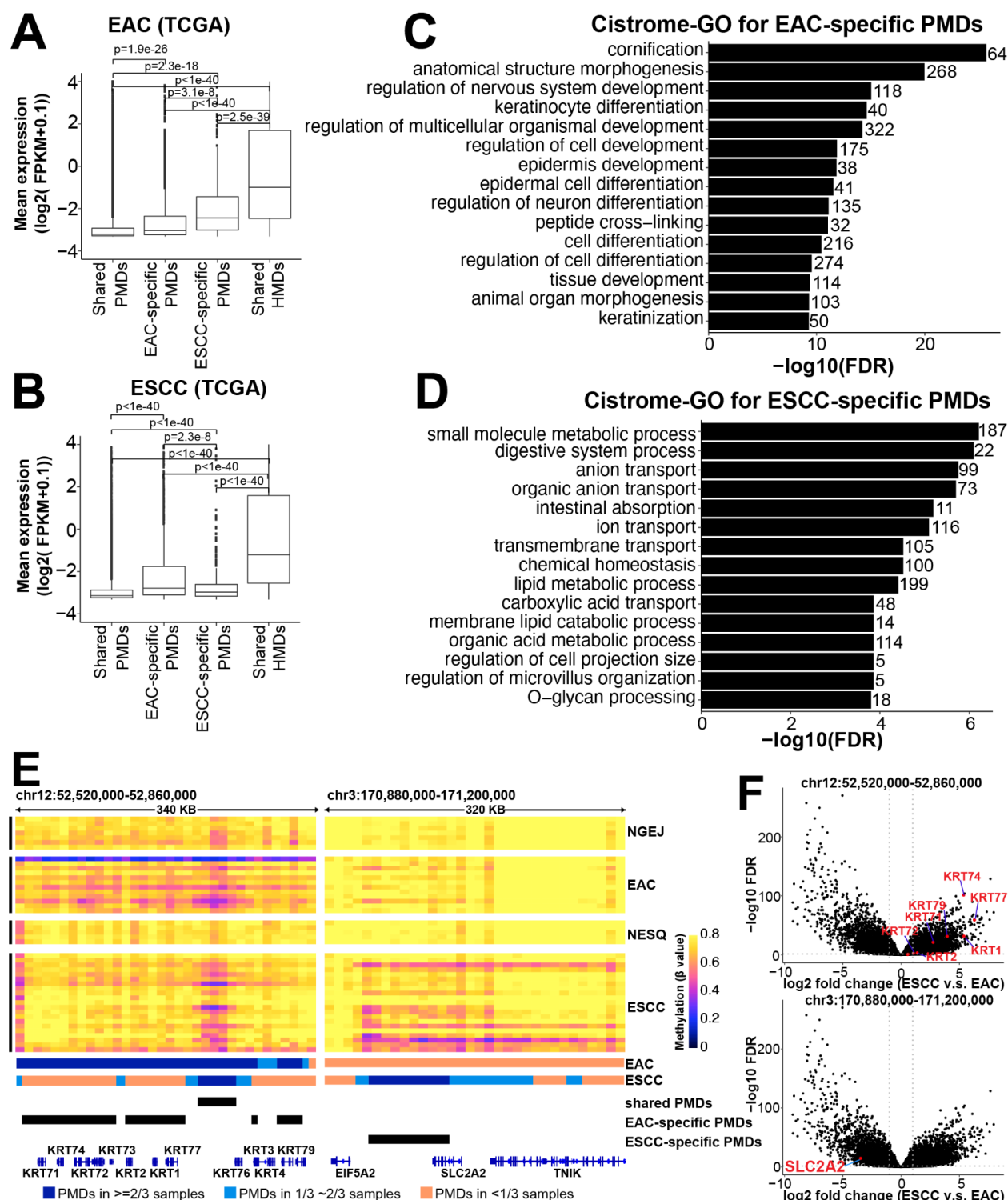903     separately for each of the WGBS PMD categories.

904

**Figure 3. Subtype-specific PMDs control cell-type-specific functions. (A-B)** In both EAC **(A)** and ESCC **(B)**, genes covered by PMDs are expressed at lower levels than those in non-PMDs in a cancer-specific manner. **(C-D)** Cistrome-GO enrichment analyses using either EAC-specific **(C)** or ESCC-specific **(D)** PMDs and the downregulated genes within them. The top 15 most

910     significant pathways are shown, and the number of genes enriched in each pathway is shown on
911     the right. **(E)** Two representative genome windows showing the methylation profiles of EAC-
912     specific (left panel) and ESCC-specific PMDs (right panel). CGI regions were masked using the
913     annotation from Irizarry et al[76]. **(F)** Volcano plots showing that genes residing within genome
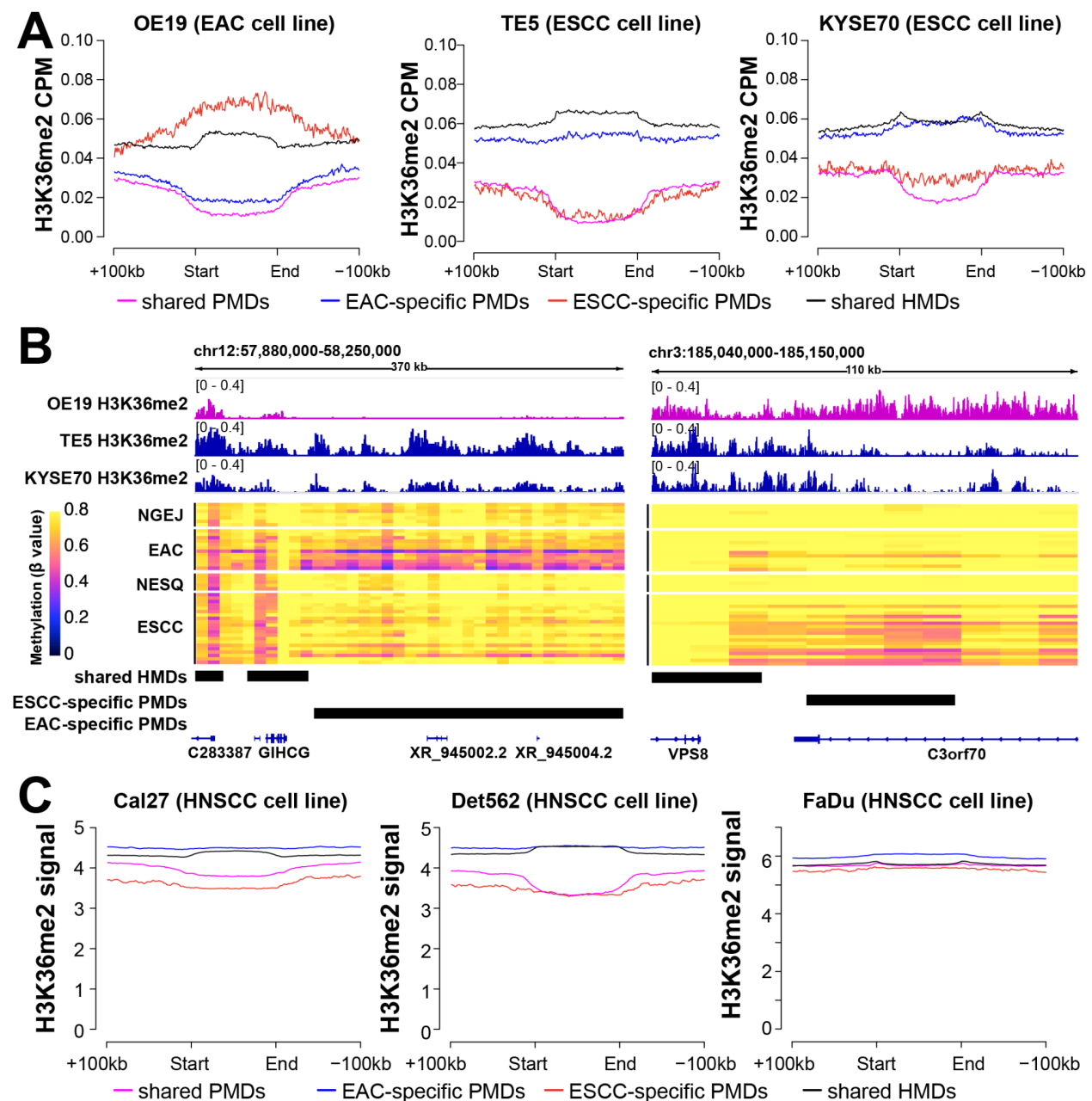914     domains in **(E)** are downregulated in corresponding cancer subtypes.

915
916



**Figure 4. The H3K36me2 mark is inversely associated with PMDs in a cell-type-specific manner. (A)** Aggregation plots of H3K36me2 ChIP-seq levels in esophageal cancer cell lines

920    across four different PMD categories: shared PMDs, EAC-specific PMDs, ESCC-specific PMDs,

921    shared HMDs. **(B)** Representative genomic loci showing H3K36me2 signal from ChIP-seq, and

922    subtype-specific PMDs from WGBS data. CGI regions were masked using the annotation from

923    Irizarry et al[76]. **(C)** Aggregation plots of H3K36me2 ChIP-seq levels in HNSCC cell lines across

924    four different PMD categories. Bigwig files of the H3K36me2 ChIP-seq signal were obtained from

925    GSE149670.

926

**Figure 5. Subtype-specific DMRs in esophageal cancer. (A)** Cancer hypoDMRs were identified from the comparison between EAC and ESCC tumors. Regions with FDR < 0.05 and absolute delta methylation levels > 0.2 were identified as DMRs. **(B)** Density plots showing the

931    size distribution of hypoDMRs; stacked bar plots displaying fractions of hypoDMRs that overlap

932    with different genomic features. **(C-D)** Aggregation plots of ATAC-seq signals from esophageal

933    cancer samples within EAC **(C)** or ESCC **(D)** hypoDMRs or random genomic regions

934    (background), which contained 10-times randomly selected regions with the same CpG density.

935    ATAC-seq signals were obtained from the TCGA and normalized with the CPM method. **(E-F)**

936    Cistrome-GO enrichment analyses using EAC **(E)** or ESCC **(F)** hypoDMRs and upregulated

937    genes in the corresponding subtype. Top 15 most significant pathways are shown. The number

938    of genes enriched in each pathway is shown on the right. **(G-H)** Transcription-factor-binding motif

939    sequences were identified by the ELMER[77] method using EAC **(G)** or ESCC **(H)** hypoDMRs as

940    the foreground and random regions as the background. **(I-J)** The most strongly enriched TFs in

941    EAC (GATA4) **(I)** and ESCC (TP63) **(J)** were chosen for the experimental validation, using TF

942    ChIP-seq, H3K27ac ChIP-seq and WGBS in matched cell lines. Peaks overlapping with subtype

943    hypoDMRs are shown on the left; the percentages of overlapped peaks are expressed in the
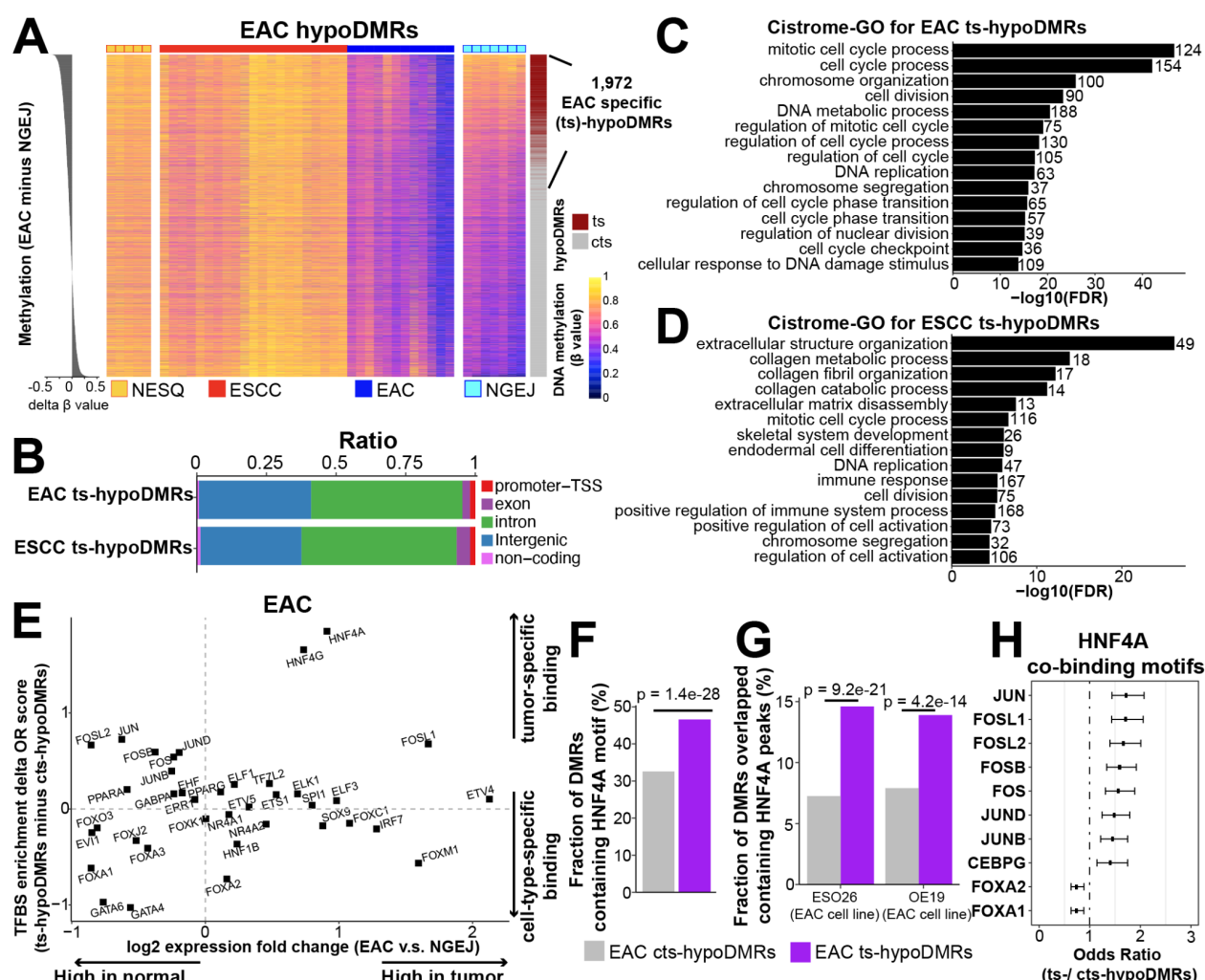
944    column plots.

945

946

**Figure 6. Identification of tumor-specific hypoDMRs. (A)** Heatmaps showing DNA methylation levels for each EAC hypoDMR. Each column denotes one sample and the row was ordered by the delta mean methylation between EAC and NGEJ (left). EAC ts-hypoDMRs were identified using a one-tailed t test between EAC tumor and NGEJ samples (right) with the FDR cutoff < 0.05. **(B)** Stacked bar plots showing fractions of ts-hypoDMRs that overlap with different genomic features. **(C-D)** Cistrome-GO enrichment analyses using either EAC **(C)** or ESCC **(D)** ts-hypoDMRs and the upregulated genes in each subtype compared with corresponding nonmalignant samples. Top 15 most significant pathways are shown. **(E)** Scatter plots showing transcription-factor-binding sites that were enriched in EAC ts-hypoDMRs compared with cts-hypoDMRs. The X axis represents the expression fold change between EAC and matched nonmalignant GEJ samples. The Y axis shows the delta enrichment score of transcription-factor-binding sites

959     between EAC ts- and cts-hypoDMRs. Expression data were from the TCGA and motif

960     enrichment analyses were performed by the ELMER method. **(F)** EAC ts-hypoDMRs

961     contained significantly more HNF4A-recognition motifs compared with cts-hypoDMRs.

962     **(G)** More HNF4A peaks overlapped with ts-hypoDMRs than cts-hypoDMRs. Peaks were

963     called from HNF4A ChIP-seq in ESO26 and OE19 cell lines. **(H)** HNF4A was predicted

964     to co-occupy with the AP-1 family in ts-hypoDMRs, while with FOXA1/2 in cts-hypoDMRs.

965     Sequence motif analysis was performed using ts- *vs.* cts-hypoDMRs containing HNF4A

966     motifs. Significant transcription factors with FDR < 0.05 are shown. OR value over 1

967     represents higher enrichment in ts-hypoDMRs, while below 1 represents higher

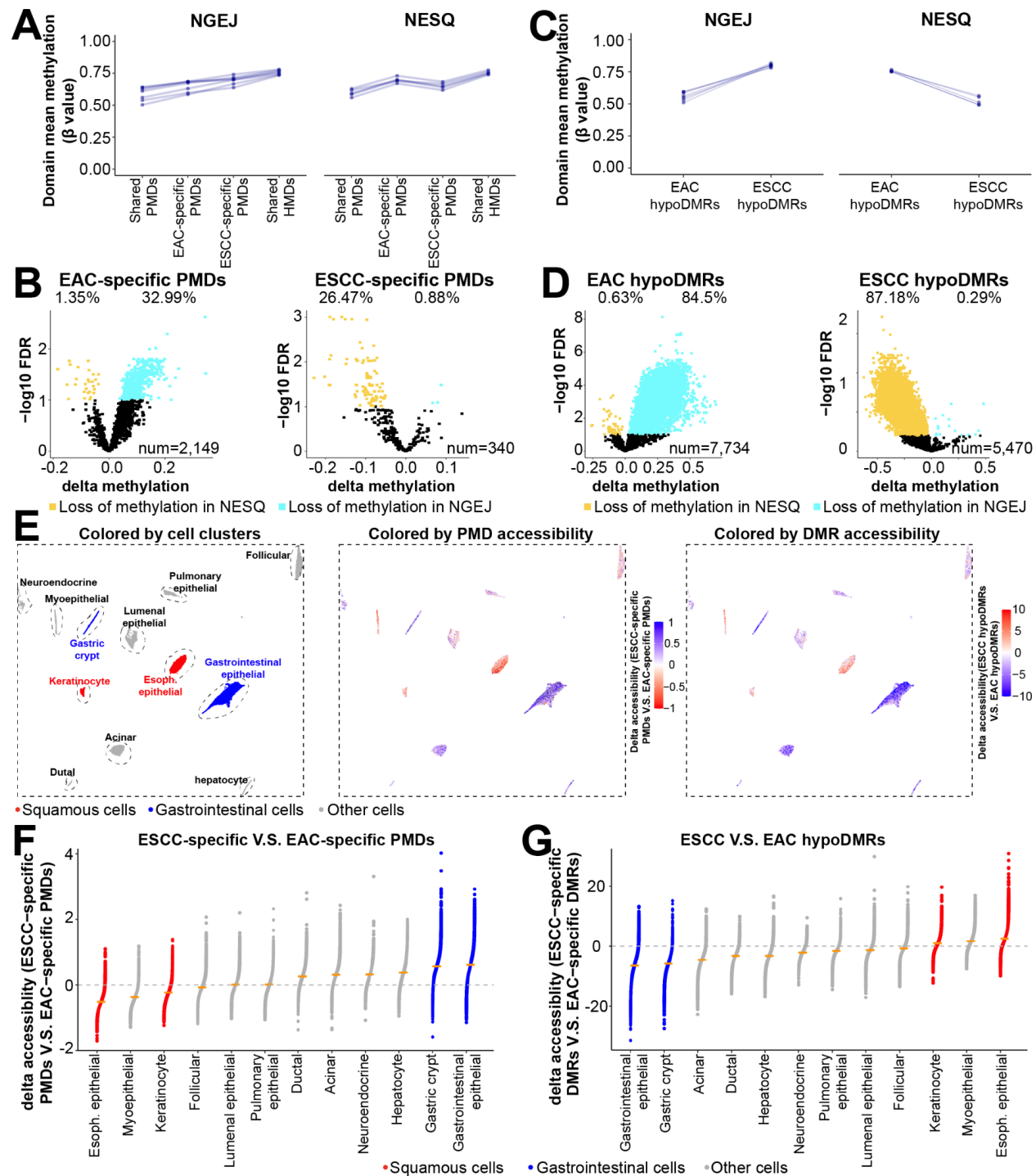968     enrichment in cts-hypoDMRs.

969

**Figure 7. PMDs and hypoDMRs exhibit strong cell-type-specific epigenomic features. (A)** Line plots showing average methylation levels for different PMD or **(C)** hypoDMR categories comparing two types of nonmalignant esophageal samples; these changes in nonmalignant samples are similar to those seen in tumors **(Fig. 2C, Supplementary Fig. 3D-E). (B)** Volcano plots showing average methylation levels for

976    different PMD or **(D)** hypoDMR categories in nonmalignant esophageal samples. Regions

977    with significant differences were determined by two-tailed t test with the FDR cutoff < 0.1.

978    **(E)** UMAP plots showing cell clusters (left), ATAC-seq levels in ESCC- *vs.* EAC-specific

979    PMDs (middle) or in ESCC- *vs.* EAC-specific hypoDMRs (right). Single-cell ATAC-seq

980    values and the cluster scheme were from Zhang et al. Total cell number is 146,305. **(F-**

981    **G)** Dot plots showing, at the sample level, delta ATAC-seq values in ESCC- *vs.* EAC-

982    specific PMDs **(F)** or in ESCC- *vs.* EAC-specific hypoDMRs **(G)**.
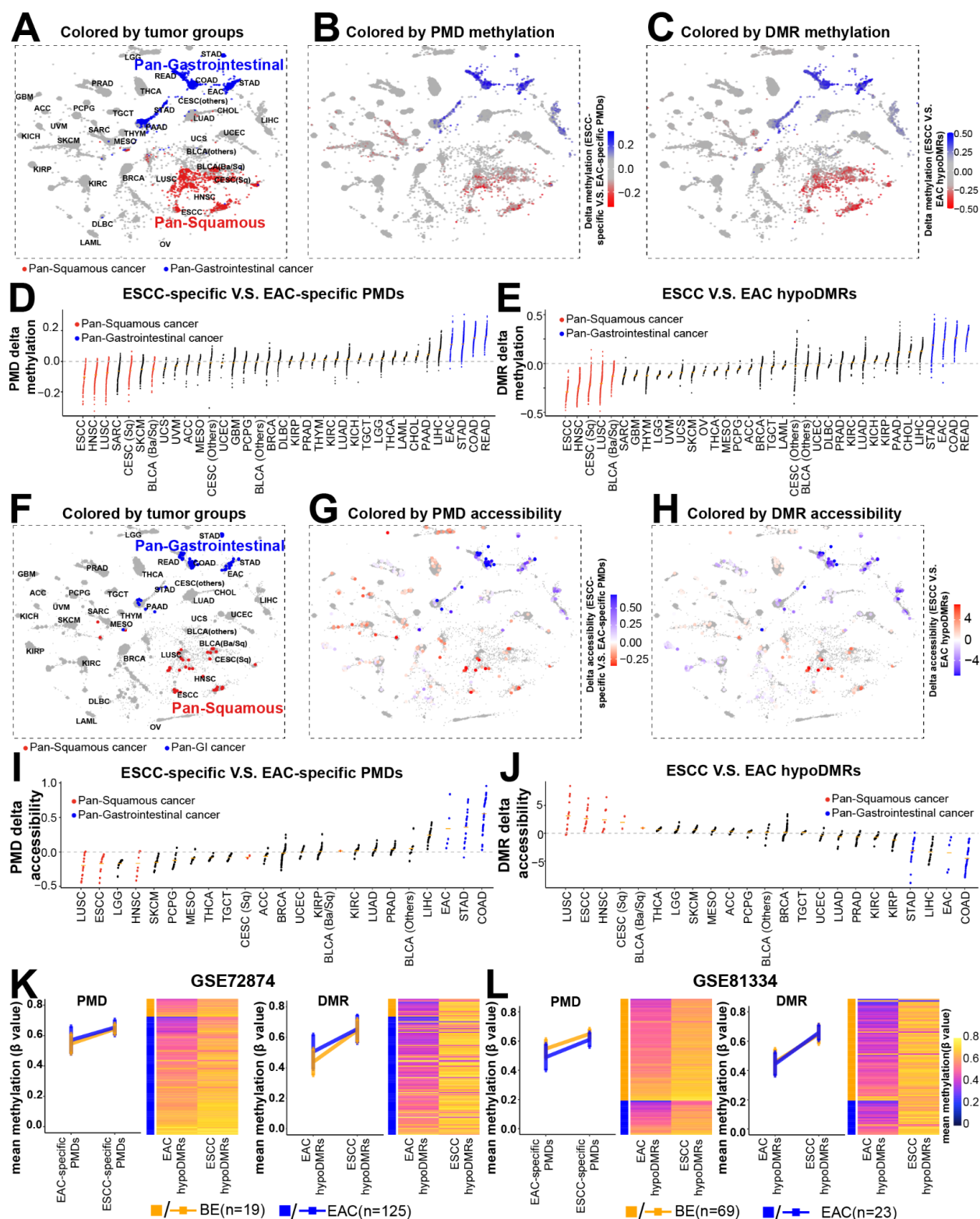
983

984

985

**Figure 8. Analyses of PMDs and hypoDMRs in pan-cancer datasets. (A-C)** TCGA

tumormap showing cancer type clusters **(A)**, DNA methylation levels in ESCC- *vs.* EAC-

988    specific PMDs **(B)**, or in ESCC- *vs.* EAC-specific hypoDMRs **(C)**. DNA methylation data

989    were obtained from the TCGA project. The TCGA-based clustering scheme denotes Pan-

990    Gastrointestinal cancers (COAD, READ, STAD and EAC) and Pan-squamous cancers

991    (ESCC, HNSC, LUSC and a subset of CESC and BLCA) are shown **(A)**. The number of

992    samples is 8,915. The detailed study name of TCGA Study Abbreviations are listed in

993    https://gdc.cancer.gov/resources-tcga-users/tcga-code-tables/tcga-study-abbreviations

994    **(D)** and **(E)** Dot plots quantification of the methylation differences in **(B)** and **(C)**,

995    respectively. **(F)** t-SNE plots showing cancer type clusters, **(G)** ATAC-seq levels in ESCC-

996    *vs.* EAC-specific PMDs or in **(H)** ESCC- *vs.* EAC-specific hypoDMRs across tumor

997    samples. ATAC-seq data were downloaded from the TCGA project. The number of

998    samples is 362. **(I)** and **(J)** Dot plots quantification of the ATAC-seq values in **(G)** and **(H)**,

999    respectively. **(K-L)** Line plots and heatmaps respectively showing average and individual

1000   methylation levels in BE and EAC samples from two different public datasets.