



## Machine learning identifies activation of RUNX/AP-1 as drivers of mesenchymal and fibrotic regulatory programs in gastric cancer

Milad Razavi-Mohseni, Weitai Huang, Yu Amanda Guo, et al.

*Genome Res.* published online May 22, 2024

Access the most recent version at doi:[10.1101/gr.278565.123](https://doi.org/10.1101/gr.278565.123)

---

<b>P&lt;P</b>	Published online May 22, 2024 in advance of the print journal.
<b>Accepted Manuscript</b>	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
<b>Creative Commons License</b>	This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see <a href="https://genome.cshlp.org/site/misc/terms.xhtml">https://genome.cshlp.org/site/misc/terms.xhtml</a> ). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at <a href="http://creativecommons.org/licenses/by-nc/4.0/">http://creativecommons.org/licenses/by-nc/4.0/</a> .
<b>Email Alerting Service</b>	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or <a href="#">click here</a> .

---

---

Advance online articles have been peer reviewed and accepted for publication but have not yet appeared in the paper journal (edited, typeset versions may be posted when available prior to final publication). Advance online articles are citable and establish publication priority; they are indexed by PubMed from initial publication. Citations to Advance online articles must include the digital object identifier (DOIs) and date of initial publication.

---

To subscribe to *Genome Research* go to:  
<https://genome.cshlp.org/subscriptions>

---

Published by Cold Spring Harbor Laboratory Press

1 **Machine learning identifies activation of RUNX/AP-1 as drivers of mesenchymal and**  
2 **fibrotic regulatory programs in gastric cancer**

3  
4 Milad Razavi-Mohseni<sup>1</sup>, Weitai Huang<sup>2</sup>, Yu A Guo<sup>2</sup>, Dustin Shigaki<sup>1</sup>, Shamaine Wei Ting Ho<sup>3</sup>,  
5 Patrick Tan<sup>3,4,5,6</sup>, Anders J Skanderup<sup>\*2</sup>, Michael A Beer<sup>\*1</sup>

6  
7 <sup>1</sup> Department of Biomedical Engineering and McKusick-Nathans Department of Genetic  
8 Medicine, Johns Hopkins University, Baltimore MD, USA

9 <sup>2</sup> Laboratory of Computational Cancer Genomics, Genome Institute of Singapore (GIS), Agency  
10 for Science, Technology and Research (A\*STAR), Singapore

11 <sup>3</sup> Laboratory of Cancer Epigenetic Regulation, Genome Institute of Singapore (GIS), Agency for  
12 Science, Technology and Research (A\*STAR), Singapore

13 <sup>4</sup> Cancer and Stem Cell Biology Program, Duke-NUS Medical School, Singapore

14 <sup>5</sup> Cancer Science Institute of Singapore, National University of Singapore, Singapore

15 <sup>6</sup> Department of Physiology, Yong Loo Lin School of Medicine, National University of Singapore,  
16 Singapore

17  
18 \* AJS and MAB are co-corresponding authors

19

20 **Abstract**

21  
22 Gastric cancer (GC) is the fifth most common cancer worldwide and is a heterogeneous  
23 disease. Among GC subtypes, the mesenchymal phenotype (Mes-like) is more invasive than  
24 the epithelial phenotype (Epi-like). While gene expression of the epithelial-to-mesenchymal  
25 transition (EMT) has been studied, the regulatory landscape shaping this process is not fully

26 understood. Here we use ATAC-seq and RNA-seq from a compendium of gastric cancer cell  
27 lines and primary tumors to detect drivers of regulatory state changes and their transcriptional  
28 responses. Using the ATAC-seq, we developed a machine learning approach to determine the  
29 transcription factors (TFs) regulating the subtypes of GC. We identified TFs driving the  
30 mesenchymal (RUNX2, ZEB1, SNAI2, AP-1 dimer) as well as the epithelial states (GATA4,  
31 GATA6, KLF5, HNF4A, FOXA2, GRHL2) in gastric cancer. We identified DNA copy number  
32 alterations associated with dysregulation of these TFs, specifically deletion of *GATA4* and  
33 amplification of *MAPK9*. Comparisons with bulk and single-cell RNA-seq datasets identified  
34 activation toward fibroblast-like epigenomic and expression signatures in Mes-like GC. The  
35 activation of this mesenchymal fibrotic program is associated with differentially accessible DNA  
36 *cis*-regulatory elements flanking upregulated mesenchymal genes. These findings establish a  
37 map of TF activity in GC and highlight the role of copy number driven alterations in shaping  
38 epigenomic regulatory programs as potential drivers of gastric cancer heterogeneity and  
39 progression.

## 40 Introduction

41 Gastric cancer (GC) is the fifth most common cancer worldwide and is more prevalent in Asia as  
42 well as among Asian Americans in the USA (Ajani et al. 2017; Taylor et al. 2015; Lui et al.  
43 2014). Despite advances in early diagnosis and treatments, the overall 5-year survival rate is  
44 still about 30% (Badgwell et al. 2017). Gastric cancer is a heterogeneous disease. A landmark  
45 TCGA study proposed four molecular subtypes of GC, including Epstein–Barr Virus (EBV)-  
46 associated, Microsatellite Instable (MSI), Genomically Stable (GS), and Chromosomal Instability  
47 (CIN) (Sohn et al. 2017). Histological subtyping using the Lauren Classification is also  
48 associated with therapeutic recommendations (Berlth et al. 2014).

49

50 To determine personalized therapeutic responses to different treatments, previous studies have  
51 demonstrated how treatment responses correlate with the gene expression profiles of GC (Tan  
52 et al. 2011) as well as the activity and mutation status of frequently-mutated genes such as  
53 *CDH1* (E-cadherin), *TP53*, and *ARID1A* (Park et al. 2016; Luo et al. 2018; Xu et al. 2023).  
54 Among such classifications are the mesenchymal (Mes) and epithelial (Epi) phenotypes, where  
55 the former confers more invasive disease (Oh et al. 2018). The process by which epithelial cells  
56 gain a mesenchymal phenotype is called the epithelial to mesenchymal transition (EMT)  
57 (Brabletz et al. 2018). Expression profiles of Mes GC and their correlation with response to  
58 therapeutic agents such as 5-Fluorouracil and PI3K inhibitors have been investigated (Tan et al.  
59 2011; Oh et al. 2018; Lei et al. 2013). In addition, the activity of transcription factors (TFs)  
60 regulating normal stomach and GC (GATA4, GATA6, KLF5, HNF4A, and EBF1) have been  
61 implicated in GC invasiveness using TF and histone ChIP-seq (Chia et al. 2015; Sheng et al.  
62 2021; Xing et al. 2020). In mesenchymal (Mes) GC, TEAD1 was identified as a major TF  
63 regulator (Ho et al. 2023) using genome-wide chromatin accessibility assays (ATAC-seq).  
64 However, there is currently a lack of systematic and unbiased approaches to uncover the  
65 complete range of TFs and epigenomic alterations underlying GC heterogeneity.

66

67 Here, we expand upon this earlier work by generating the first systematic map of gastric cancer  
68 TF activity derived from an unbiased machine learning approach applied to ATAC-seq datasets  
69 from a broad compendium of GC cell lines and primary TCGA tumors (Corces et al. 2018; Bass  
70 et al. 2014), coupled with RNA-seq analysis for interpretation of the altered regulatory pathways.  
71 Our analysis aims to infer the complete set of TFs whose activity explains the heterogeneity of  
72 chromatin landscapes across GC and connect TF activity through enhancers to transcriptional  
73 activation of specific gene responses. We further investigate mutational and copy number  
74 alterations as genomic aberrations underlying the dysregulation of these key TFs. Finally, we

75 relate GC cell line subtypes to other tissues using models trained on accessibility and single-cell  
76 transcriptional profiles.

77

## 78 Results

### 79 Sequence-based machine learning model of ATAC-seq profiles 80 identify distinct regulatory programs driving GC heterogeneity

81 To understand differences in epigenomic state profiles in gastric-cancer derived cell lines, we  
82 generated ATAC-seq data from a collection of 25 cell lines (Ho et al. 2023). Principal  
83 Component Analysis (PCA) of variation in the activity of distal peaks (Methods) showed  
84 significant heterogeneity in chromatin accessibility across the cell lines, forming three clusters  
85 (Fig. 1A) with similar genomic accessibility profiles. We used a machine learning approach,  
86 gkm-SVM (Ghandi et al. 2014a; Beer et al. 2020; Lee et al. 2015), to identify the sequence  
87 features (transcription factor binding sites, TFBS) driving these variations in chromatin  
88 accessibility. We chose gkm-SVM for its combination of predictive accuracy and interpretability  
89 (Ghandi et al. 2014a; Beer et al. 2020; Lee et al. 2015; Shigaki et al. 2019; Ghandi et al. 2016;  
90 Beer 2017; Yan et al. 2021; Kreimer et al. 2017; Lee 2016; Gate et al. 2018; Mo et al. 2016),  
91 and because of its advantages relative to DNNs when training on smaller numbers of  
92 differentially active peaks (~2000). For each cell line, we trained a gkm-SVM model on the top  
93 10k distal ATAC-seq peaks against random genomic sequences, following our standard pipeline  
94 (Beer et al. 2020). On each training, following the default training method in the gkmSVM-R  
95 package (Ghandi et al. 2016), the AUROC was evaluated on 5-fold cross validation (CV) and  
96 the reported AUROC is the average of 5 CV test sets. The median cross-validation AUROC is

97 0.925. All cross-fold validation sets produce similar sets of features and have very similar weight  
98 vectors. Because this test set AUROC is high, there is little evidence of overfitting. Also,  
99 following (Ghandi et al. 2016), for feature interpretation, but not for quantitative assessment, we  
100 use the gapped *k*-mer weight vector when training on all the data to reduce variability in the  
101 inferred features. Each gkm-SVM model summarizes the TFBS required to classify the ATAC-  
102 seq peaks in a vector of gapped *k*-mer weights which can predict the accessibility of the  
103 sequences. We clustered the gkm-SVM weight vector from each cell line (Fig. 1B) and found  
104 nearly identical clusters to Fig. 1A, thus the similar genomic accessibility profiles are explained  
105 by activation of shared TFs.

106

107 Based on these differentially active TFs and their resultant differentially expressed target genes,  
108 which largely align with EMT and mesenchymal-related pathways (see RNA-seq analysis  
109 below), we label these three sets of GC cell lines as Epithelial-like (Epi-like), Intermediate, and  
110 Mesenchymal-like (Mes-like). To isolate the TFs driving most of the variability between these  
111 sets of cell lines, we next trained gkm-SVM on the most differentially accessible peaks when  
112 averaged over the Epi-like and Mes-like cell lines, with Mes-like peaks as the positive set, and  
113 Epi-like peaks as the negative set (Supplemental Fig. S1, S2). Again, this generated a highly  
114 predictive sequence model and the gkm-SVM weights with the largest difference between the  
115 Epi-like and Mes-like cell lines were associated with differential activity of TFs binding the  
116 specific motifs: RUNX, AP-1 (FOS/JUN), TEAD, ZEB, HNF4A, FOXA, GRHL, GATA, and KLF  
117 (Fig. 1C, Supplemental Fig. S3). Since ZEB1 and ZEB2 are known transcriptional repressors  
118 (Aghdassi et al. 2012; Byles et al. 2012), ZEB activity is indicated by a negative motif score. To  
119 assess if these TFs also drive variation in tumors, we next trained gkm-SVM on differential  
120 peaks in normal stomach vs. primary stomach tumor samples (TCGA-STAD, The Cancer  
121 Genome Atlas for Stomach Adenocarcinoma) (Fig. 1C, Supplemental Fig. S4, S5). The TFs  
122 detected largely overlap with those in the cell lines, indicating that these TFs also drive variation

123 in the regulatory state between normal stomach tissue and GC tumors. While there is significant  
124 heterogeneity among the tumor samples, and only a fraction of the cells in a tumor are likely in a  
125 mesenchymal state, we find significant evidence that the Mes-like regulatory program identified  
126 in the cell lines is activated in the STAD samples relative to normal stomach. 70% of the high  
127 quality TCGA-STAD ATAC-seq samples detect RUNX, and 70% detect AP-1 (Supplemental  
128 Fig. S5, Methods). Finally, we trained gkm-SVM on differentially accessible peaks in each pair  
129 of cell lines (300 pairs, see Methods) to detect the full set of TFs explaining the ATAC-seq  
130 variation across the cell line collection. RUNX and AP-1 (and lack of ZEB) were strong TFBS  
131 signals driving accessibility in Mes-like cell lines, while GRHL, GATA, and KLF TFBS were  
132 strongly associated with most Epi-like samples (Fig. 1D, 1E). In summary, this ATAC-seq  
133 analysis of open chromatin identifies RUNX, AP-1, ZEB, GRHL, GATA, and KLF as novel  
134 regulatory drivers of the Epi-like and Mes-like transcriptional states in gastric cancer. These  
135 predictions from gkm-SVM are validated by experiments showing strong RUNX2 and JUN  
136 (which binds AP-1) binding at Mes-like enhancers from ChIP-seq in LPS141 (Bevill et al. 2023)  
137 and strong KLF5, GATA4, and GATA6 binding at Epi-like enhancers from ChIP-seq generated  
138 in AGS (Chia et al. 2015; Liu et al. 2020) (Fig. 1F, Supplemental Fig. S6).

## 139 **Regulatory and expression signatures of EMT and advanced-vs-** 140 **early GC explain the differences between GC groups**

141 We used RNA-seq data from the 25 GC cell lines to understand how the differential TF activity  
142 inferred from the ATAC-seq data affects gene expression. The three groups of GC cell lines,  
143 previously clustered based on their chromatin accessibility, were recapitulated using the  
144 expression of tissue-specific genes (Fig. 2A, Supplemental Fig. S7, Methods), demonstrating  
145 that the variation in RNA-seq is concordant with that of ATAC-seq. The only exception is  
146 SNU484, where the ATAC-seq and RNA-seq profiles are less consistent than they are for every

147 other cell line. Because we can use gkm-SVM AUROC to assess the quality of the ATAC-seq  
148 data, we are more confident of SNU484 ATAC-seq since the AUROC is 0.963. Because we  
149 don't have an independent metric to assess the quality of the RNA-seq, we weighted the ATAC-  
150 seq more heavily for SNU484 and labelled it Epi-like. For all other samples the ATAC-seq and  
151 RNA-seq based cluster labels are consistent.

152

153 We next compare gene expression in GC cell lines with that of primary stomach tissue from  
154 TCGA-STAD and TCGA normal adjacent tissue (Methods). The combined PCA with cell lines  
155 and TCGA shows that there are inherent differences between the stomach-derived cell lines  
156 and stomach primary tissue on the first principal component (PC1, Fig. 2B), as noted in previous  
157 studies (Aran et al. 2015; Yu et al. 2019). PC2 shows TCGA-Normal samples are more aligned  
158 with Epi-like GC cell lines, while TCGA-STAD cancer samples are closer to Mes-like cell lines  
159 (Supplemental Fig. S8).

160

161 To associate these expression differences with biological processes, we performed gene set  
162 enrichment analysis using PCA gene weights and GSEA (Mootha et al. 2003; Subramanian et  
163 al. 2005; Liberzon et al. 2011, 2015). The first principal component was strongly associated with  
164 mesenchymal and EMT phenotypes, advanced-vs-early gastric cancer and cancer invasiveness  
165 (Fig. 2C, hypergeometric FDR <  $10^{-40}$ ) (Liberzon et al. 2015; Charafe-Jauffret et al. 2006;  
166 Schuetz et al. 2006; Vecchi et al. 2007). EMT is also the top hit when using expression changes  
167 between Mes-like and Epi-like clusters (Supplemental Fig. S9, S10) instead of PC1 weights. We  
168 discovered a similar pattern of enrichment comparing TCGA-STAD and normal tissues with the  
169 cell lines (Fig. 2D, which is gene set enrichment analysis on PC2 of Fig. 2B). Taken together,  
170 this demonstrates that the biological differences between the Mes-like and Epi-like clusters of  
171 GC cell lines overlap with the EMT and invasiveness phenotypes. Because of these  
172 associations, we label the red group of cell lines in Fig. 1A and Fig. 2A "*Mes-like*" for their

173 mesenchymal expression signature (EMT); the green group “*Epi-like*” for their more normal  
174 epithelial expression profile; and the gray group of cell lines is labeled “*Intermediate*”. The  
175 Intermediate set lies along the continuum between Epi-like and Mes-like, but is distinguishable  
176 from the Mes-like set through the stronger activity of ZEB and GRHL motifs, and distinguishable  
177 from the Epi-set by the lack of GATA activity, as shown in (Fig. 1E). Our ATAC-seq analysis  
178 allows us to identify TF regulators with these EMT and invasiveness gene expression  
179 signatures.

180

181 TF regulators detected in ATAC-seq are markers of EMT and invasiveness  
182 and are differentially expressed in Mes-vs-Epi GC cell lines

183 To further explore the connection between regulatory states and cancer invasiveness, we  
184 generated a set of 1770 differentially expressed (DE) genes between Mes-like and Epi-like cell  
185 lines (Methods). These DE genes include the regulator TFs implicated from ATAC-seq, and  
186 many genes previously implicated in epithelial GC or EMT (Supplemental Table S1 and  
187 Supplemental Table S2). We observed that TFs such as *GATA4*, *GATA6*, *KLF5*, *HNF4A*,  
188 *GRHL2*, *FOXA1*, *FOXA2*, and *FOXA3* in the Epi-like cell lines, as well as *SNAI2*, *ZEB1*, *FOSL1*  
189 (AP-1), and *RUNX2* in Mes-like, not only have a significantly higher expression (Fig. 2E, 2F,  
190 Supplemental Table S1) but also have differentially accessible DNA-binding sites in Mes-like vs.  
191 Epi-like GC cell line ATAC-seq (Fig. 1C). Further, we assessed the expression of these DE TFs  
192 in an independent patient cohort from the Asian Cancer Research Group (ACRG) (Cristescu et  
193 al. 2015) (Fig. 2G, Methods), as well as in TCGA-STAD (Supplemental Fig. S11), and found  
194 that differential TF expression is associated with overall survival, suggesting that activation of  
195 these regulator TFs has both biological and disease prognostic utility.

196

197 Many of the differentially expressed genes have been previously implicated in cancer. In the  
198 Epi-like upregulated group, *KLF5*, and *HNF4A* along with *GATA4* and *GATA6* have been shown  
199 to promote GC (Chia et al. 2015). *GRHL2* plays an anti-EMT role in GC (Sheng et al. 2021; Luo  
200 et al. 2019) through TGFB signaling, whereas, *GATA6* suppresses EMT by promoting the  
201 expression of *FOXA2* and *CDH1* while decreasing the activity of *SNAI2*, *TWIST1*, *ZEB1* and  
202 *VIM* in pancreatic cancer (Martinelli et al. 2017). TFs in the Mes-like upregulated group, *ZEB1*,  
203 *RUNX2*, and *SNAI2* are known regulators of EMT (Brabletz et al. 2018; Niu et al. 2012; Roche  
204 2018). Additionally, *FOSL1* dimerizes with *JUN* family members to bind the AP-1 motif and  
205 drives EMT by upregulating *ZEB1/2* and *SNAI2* (Bakiri et al. 2015; Guneri-Sozeri et al. 2023).  
206 The combined differential expression of these TF genes and the identification of their TF binding  
207 sites as the largest signals in predictive sequence modeling of the chromatin accessibility  
208 profiles strongly suggest that the variability in activation of this small set of TFs is responsible for  
209 driving the EMT program. Although this finding is in the context of GC, the fact that LPS141  
210 arrived at a similar transcriptional state and has *RUNX2* and *JUN* binding at the Mes-like  
211 enhancers suggests that this may be a more general mechanism of activating EMT across other  
212 tumor types (Supplemental Fig. S6).

213

## 214 Mutations and copy number alterations are associated with 215 distinct regulatory states in GC

216 We next sought to identify possible mechanisms driving the variation in regulatory state across  
217 the GC cell lines. To understand how the activity of TFs in Mes-like and Epi-like cell lines  
218 correlated with known genetic variation in gastric cancer, we examined TF gene expression in  
219 TCGA-STAD patients. *CDH1*, a recurrently mutated GC driver gene (Supplemental Table S3)  
220 (Luo et al. 2018; Xu et al. 2023; Lei et al. 2013; Fodde 2002), was downregulated in Mes-like

221 GC cell lines compared to Epi-like. We examined the expression of DE TFs in *CDH1*-mutated  
222 vs. *CDH1*-WT TCGA-STAD patients (Fig. 3A, Supplemental Fig. S12, Methods). *CDH1*-mutated  
223 TCGA-STAD tumors showed upregulation of the Mes-like TFs. Conversely, the *CDH1*-WT  
224 tumors displayed upregulation of Epi-like TFs (Fig. 3A). This result indicates that *CDH1*  
225 mutations in GC correlate with upregulation of key EMT-promoting TFs and downregulation of  
226 EMT-inhibiting TFs.

227

228 In addition, the TFs *GATA4*, *GATA6*, *SMAD4*, *KLF5*, and *HNF4A*, are among the TFs with the  
229 highest rates of DNA Copy Number (CN) alterations in TCGA-STAD patients (Fig. 3B,  
230 Supplemental Fig. S13). *GATA4*, *GATA6*, and *SMAD4* are known to be important in endoderm  
231 and stomach development (Li et al. 2019), and *SMAD4* has been identified as a stomach cancer  
232 driver gene (Martínez-Jiménez et al. 2020). *GATA4* and *GATA6* have lower expression in Mes-  
233 like GC cell lines (Fig. 2F), consistent with gkm-SVM analysis (Fig. 1E). ATAC-seq profiles for  
234 Mes-like, Intermediate, and Epi-like GC cell lines also identify two distal peaks (*GATA6\_E1* and  
235 *GATA6\_E2*) in the *GATA6* locus (Fig. 3C) whose accessibility correlates with *GATA6*  
236 expression (Supplemental Fig. S14, S15). *GATA6\_E1* is upstream of *GATA6* and bound by  
237 *GATA6* in AGS, and *GATA6\_E2* is intronic and is also bound by *GATA6* in AGS, and by *SMAD4*  
238 and *ZEB1* in HepG2 cells. The low ATAC signal in Mes-like cell lines at *GATA6\_E2* is consistent  
239 with the upregulation of the repressor *ZEB1* and downregulation of *GATA6*. In further support of  
240 their regulatory role, both *GATA6\_E1* and *GATA6\_E2* are contained in loop extrusion (LE)  
241 model predicted CTCF loops containing the *GATA6* promoter (Xi and Beer 2021; Luo et al.  
242 2023). Overall, the up-regulation of the epithelial regulator *GATA6* is consistent with the Mes-  
243 like and Epi-like ATAC-seq profiles through *GATA6\_E1* and *GATA6\_E2*.

244 Reduced GATA and increased AP-1 regulatory activity in GC cell lines are  
245 accompanied by DNA copy number alterations

246 The gkm-SVM inferred differential GATA activity can be attributed to either GATA6 or GATA4  
247 dysregulation, and both *GATA4* and *GATA6* have lower expression in Mes-like GC (Fig. 2F). To  
248 identify potential driving mechanisms, we compared CN changes in Epi-like cell lines vs. non-  
249 Epi cell lines (Mes and Intermediate), using DNA sequencing (Methods). On average, *GATA4*  
250 has a higher copy number in Epi-like cell lines (Fig. 3D), which correlates with lower expression  
251 of *GATA4* in Mes-like cell lines (Fig. 3E,  $r = 0.56$ ). While less than 1% of TCGA-STAD patients  
252 have a mutation in *GATA4* according to cBioPortal (Gao et al. 2013), *GATA4* CN alterations are  
253 more frequent in TCGA-STAD (Fig. 3B). Further, an interval of Chr8p containing *GATA4* is  
254 frequently deleted in human epithelial cancers (Cai and Sablina 2016). Thus, GATA CN  
255 alterations are a potential mechanism for *GATA4* expression changes (Fig. 2F) and GATA  
256 activity in gastric cancer cell lines and primary tumors (Fig. 1C). We observed an opposite  
257 pattern for *MAPK9* kinase which is amplified in the Mes-like cell lines compared to Epi-like (Fig.  
258 3D, 3E). MAPK and JNK family members are known to be regulators of AP-1 (Karin et al. 1997),  
259 and *MAPK9* amplification is associated with higher expression of AP-1 complex members  
260 (*FOSL1*, Fig. 2F) and higher AP-1 transcriptional activity in Mes-like cell lines and primary  
261 stomach cancer (Fig. 1C). Taken together, we see evidence for DNA copy number changes in  
262 transcription factors and signaling pathway components as potential drivers of transcriptional  
263 dysregulation in gastric cancer cell lines.

264 Mes-like GC cell lines have expression and regulatory signatures  
265 similar to fibroblasts, while Epi-like GC cell lines retain  
266 gastrointestinal signatures

267 Because we detected the activation of non-stomach TFs among many GC cell lines, we next  
268 sought to characterize their cell states by comparing to known cell types. We first compared the  
269 GC cell line expression profiles to known normal cell types by calculating their correlation with  
270 Genotype-Tissue Expression (GTEx) RNA-seq data (Lonsdale et al. 2013) (Methods, Fig.4A).  
271 Mes-like GC cell line expression was more similar to fibroblasts ( $r=0.41$ ) than to stomach  
272 ( $r=0.29$ ), while Epi-like and Intermediate were most similar to stomach and esophagus mucosa,  
273 respectively (Fig. 4A). We further confirmed that the gene expression in Mes-like GC cell lines is  
274 more similar to fibroblast-derived cells than stomach in two additional independent RNA  
275 datasets from ENCODE and TCGA-STAD (Methods, Fig. 4B) (The ENCODE Consortium 2012;  
276 Luo et al. 2020). While TCGA-Normal is more similar to normal stomach than TCGA-STAD, the  
277 similarity of TCGA-STAD to fibroblast is weaker than the Mes-like cell lines (Fig. 4B).

278

279 We extended these comparisons to DNase-seq data and found consistent similarities between  
280 the chromatin accessibility of the Mes-like GC cell lines with ENCODE DNase-seq of fibroblast  
281 (Fig. 4C). We trained gkm-SVM sequence models on open chromatin regions for each GC cell  
282 line, TCGA-STAD, ENCODE fibroblast and stomach tissues, and compared their weight vectors  
283 (Methods). The chromatin accessibility of Mes-like GC is highly correlated with the fibroblast-  
284 derived cell lines HT1080 and ELR, as well as primary fibroblast tissue from ENCODE, as  
285 opposed to primary stomach. In contrast, the Epi-like cell lines and TCGA-STAD  
286 ATAC-seq were more correlated with the primary stomach.

287

288 In addition to these average correlations of expression and chromatin accessibility, we now  
289 show that differential TF activity among GC cell lines positions them along a continuum of cell  
290 states between normal stomach and fibroblast. Scatter plots of correlation to stomach and  
291 fibroblast (Fig. 4D, 4E, Supplemental Fig. S16, S17) show that TCGA-Normal had the highest  
292 correlation with ENCODE normal stomach and the lowest with fibroblast, while TCGA-STAD  
293 tumors have a higher correlation with fibroblast and a lower correlation with stomach (Fig. 4D).  
294 The Mes-like GC cell lines have a higher similarity to fibroblast than the primary tumors,  
295 correlations being as high as 0.8. Consistent with the RNA-seq profiles, gkm-SVM models  
296 trained on chromatin accessibility profiles also show that Mes-like cell line accessibility is more  
297 similar to fibroblast accessibility than normal stomach tissue (Fig. 4E, Supplemental Fig. S18,  
298 S19). While it is clear that GC cell lines are becoming more similar to fibroblasts by activation of  
299 fibroblast TFs, the bulk TCGA-STAD profiles may be moving toward a fibroblast cell state either  
300 by interaction with stroma through fibrosis or by the presence of cancer-associated fibroblasts  
301 (CAF) in the tumor microenvironment (Kalluri and Zeisberg 2006; Piersma et al. 2020).

302

303 We next addressed cell line similarity at the single cell level using stomach cancer scRNA-seq  
304 data from Kim et al. (Kim et al. 2022) (Methods). We identified 6 distinct populations of cells  
305 using UMAP (McInnes et al. 2020) (Fig. 5A, Supplemental Fig. S20), and based on marker gene  
306 expression (Methods, Supplemental Table S4), we will identify clusters 2 and 5 as  
307 gastrointestinal and cluster 6 as fibroblasts. We compared the expression of GC cell lines to  
308 these single-cell clusters and found high similarity between stomach clusters 2 and 5 with Epi-  
309 like GC cell lines (Fig. 5B) and high similarity between Mes-like GC cell lines and fibroblast  
310 cluster 6 (Fig. 5C). Expression of stomach and fibroblast specific genes from GTEx RNA-seq  
311 confirms our association of these clusters with these cell types (Fig. 5D, 5E, Methods).

312 Fibrotic gene expression in GC cell lines is driven by activation of  
313 flanking enhancers

314 While we observed similarity of ATAC-seq and RNA-seq profiles across Mes-like and Epi-like  
315 cell lines (Fig. 1A, 2A), we next show that their differential expression is consistent with direct  
316 induction by enhancers flanking Mes-like and Epi-like genes. We calculated the average ATAC-  
317 seq signal in all distal ATAC-seq peaks within 50kb of a Transcription Start Site (TSS) of  
318 differentially expressed genes between Epi-like or Mes-like GC cell lines. The accessibility of  
319 peaks flanking genes more highly expressed in Epi-like cell lines is 1.763 times higher in Epi-  
320 like cell lines compared to Mes-like, and the accessibility of peaks flanking genes more highly  
321 expressed in Mes-like cell lines is 1.756 times higher in Mes-like cell lines compared to Epi-like  
322 (Fig. 6A, 6B), providing a direct connection between altered chromatin state and expression  
323 responses.

324

325 Among the genes upregulated in Mes-like cell lines are the fibrotic gene *FGF5* and many  
326 collagen genes downstream of the TGF $\beta$ /SMAD pathway (Shin et al. 2019; Xie et al. 2020;  
327 Verrecchia et al. 2001). Their role in gastrointestinal cancer (Luo et al. 2019) and cancer-  
328 associated fibroblasts (Bordignon et al. 2019) has been well-established. Inhibition of FGF5  
329 decreases proliferation and metastasis in hepatocellular carcinoma (Fang et al. 2015). *COL1A1*  
330 has a similar effect on tumor behavior and proliferation (Nissen et al. 2019). *FGF5*, *COL1A2*,  
331 *COL6A3*, *COL5A1*, *COL12A1*, *COL1A1*, and *COL6A2* are upregulated in Mes-like GC and have  
332 1.5-fold to 6-fold higher accessibility in flanking ATAC peaks. The members of the fibrotic  
333 pathways *SMAD3*, *SMAD4*, *TGF $\beta$ 2*, *COL1A1*, *FGF1*, *FGF5*, and *FGF7* are differentially  
334 expressed (Fig. 6C, Mann–Whitney *U* test  $p < 0.05$ ) in Mes-like and Epi-like GC.

335

336 To understand the regulatory mechanism for the differential expression of *COL1A1* and *FGF5*,  
337 we compared the chromatin accessibility of their flanking *cis*-regulatory regions (Fig. 6D, 6E).  
338 For *COL1A1*, we found upstream distal peaks *COL1A1\_E1* and *COL1A1\_E2* differentially active  
339 in Mes-like and Epi-like GC cell lines and in primary TCGA-STAD tumors compared to normal  
340 stomach (Fig. 6D). They are both located within a loop-extrusion (LE) model predicted CTCF  
341 loop (Xi and Beer 2021). Similarly, *FGF5* is flanked by an upstream distal peak (*FGF5\_E*) (Fig.  
342 6E), which is active in primary TCGA-STAD cancer and inactive in healthy stomach tissue.  
343 Across cell lines, the expression of *COL1A1* is highly correlated with the activity of peaks  
344 *COL1A1\_E1* and *COL1A1\_E2* ( $r=0.64$  and  $0.58$ , Supplemental Fig. S21, S22). Similarly, the  
345 relatively higher expression of *FGF5* in 3 of the 5 Mes-like cell lines is correlated with the activity  
346 of the distal peak (*FGF5\_E*) ( $r=0.74$ , Supplemental Fig. S23). *COL1A1\_E1*, *COL1A1\_E2*, and  
347 *FGF5\_E* are all bound by RUNX2 and JUN in LPS141 (Bevill et al. 2023). Therefore, *COL1A1*  
348 and *FGF5* are dysregulated in both the Mes-like GC cell lines and in stomach cancer tissue,  
349 consistent with increased chromatin accessibility of newly identified flanking distal peaks which  
350 are likely acting as enhancers.

## 351 Mes-like GC cell lines have a distinct regulatory landscape 352 compared to fibroblasts

353 Since the Mes-like cell lines are epigenomically and transcriptionally similar to fibroblasts, we  
354 next sought to verify that they are not simply normal fibroblasts, or possibly cancer-associated  
355 fibroblasts (CAF) (Sahai et al. 2020). Normal fibroblasts would be distinguishable by lower  
356 mutation or CN alteration rates. However, we found that the Epi-like, Intermediate, and Mes-like  
357 GC cell lines had indistinguishable mutation rates in the coding regions (Fig. 7A,  $t$ -test  $p > 0.12$ ,  
358 Methods). Additionally, we calculated the average rate of CN alterations in each cell line (Fig.  
359 7B, Methods), and the three groups of cell lines had similar rates. While *GES1* is derived from

360 normal stomach epithelium, it was transformed with SV40 virus containing the T antigens, which  
361 downregulate TP53 and appear to have mimicked the high mutation burden of the other cell  
362 lines, which were not transformed with SV40 and do not have any ATAC-seq reads mapping to  
363 SV40. Because of these mutation and copy number profiles, the Mes-like cell lines are clearly  
364 distinct from normal fibroblasts.

365

366 The Mes-like cell lines are also distinct from normal fibroblasts in their chromatin accessibility  
367 profiles. PCA analysis on distal open-chromatin regions in GC cell lines, normal fibroblast and  
368 stomach ENCODE DNase-seq, showed differences between each group (Fig. 7C, Methods,  
369 analysis similar to Fig. 1A). The variation between Epi-like, Intermediate, and Mes-like GC cell  
370 lines is aligned with an axis that also explains most of the variance between stomach and  
371 fibroblast accessibility profiles. In addition, gkm-SVM sequence models trained on differentially-  
372 accessible ATAC distal peaks between Mes-like GC and normal fibroblast highlight that TFs  
373 binding AP-1, RUNX, GATA, RARA, and KLF motifs have higher activity in Mes-like GC, while  
374 fibroblasts have higher accessibility explained by TWIST and ZEB DNA-binding sites (Fig. 7D).  
375 By aggregating motif activity scores following pairwise comparisons between each Mes-like cell  
376 line and each ENCODE fibroblast, we rediscovered AP-1 and RUNX to be active in Mes-like GC  
377 as opposed to TWIST and ZEB being more active in normal fibroblasts (Fig. 7E).

378

379 Although the expression profiles of Mes-like GC cell lines are similar to fibroblasts, there are  
380 also systematic expression differences consistent with the above noted chromatin accessibility  
381 changes. TFs with higher expression levels in normal stomach such as *KLF5*, *FOXA1*, *HNF4A*,  
382 and *POU5F1* were among the TFs having a higher expression in Mes-like GC compared to  
383 normal ENCODE fibroblast (Fig. 7F, Methods,  $FDR < 0.01$ ,  $|\log_2FC| > 2$ ). The differential activity  
384 of KLF, TWIST, and AP-1 discussed in the motif analysis above is consistent with the  
385 upregulation of *KLF5* and *FOSL1* in Mes-like GC, compared to the higher expression of *TWIST2*

386 in normal ENCODE fibroblast (Fig. 7F). In addition, fibroblast marker genes such as *FN1*, *DCN*,  
387 *COL6A1*, *COL1A1*, and *FGF7* have higher expression in ENCODE fibroblast compared to Mes-  
388 like GC (Fig. 7G). Taken together, we see that despite similar chromatin accessibility and gene  
389 expression signatures between Mes-like GC and fibroblast tissues, the epigenomic state of  
390 Mes-like GC can be clearly distinguished from fibroblasts by the activity of AP-1, RUNX, and the  
391 stomach motif KLF, as well as lower activity of TWIST and ZEB compared to fibroblasts.

## 392 Discussion

393 We provide the first systematic map of gastric cancer TF activity inferred from an unbiased  
394 machine learning approach applied to ATAC-seq and RNA-seq data. Our sequence-based  
395 machine learning analysis revealed that most of the variation among GC cell lines is driven by a  
396 fibrotic vs. epithelial regulatory program differentially activated by a small set of TFs, most of  
397 which have not previously been directly associated with EMT in gastric cancer. We used a panel  
398 of gastric cancer-derived cell lines in combination with the TCGA stomach-adenocarcinoma  
399 (TCGA-STAD) cohort to identify heterogeneous regulatory programs in gastric cancer (GC) and  
400 their concomitant transcriptional responses. The inferred cell-line regulatory programs explained  
401 much of the variation among TCGA-STAD samples. Our direct analysis of GC cell lines allowed  
402 for the isolation of cancer cell regulatory programs without noise derived from intra-tumor  
403 heterogeneity and variation in tumor immune infiltration.

404

405 Using ATAC-seq and computational sequence models, along with RNA-seq and scRNA-seq, we  
406 identified the largely novel set of TFs involved in mesenchymal (Mes-like) GC (RUNX2, SNAI2,  
407 ZEB1, AP-1 dimer) as opposed to the less invasive epithelial state (Epi-like GC) and its  
408 regulators (GATA4, GATA6, KLF5, HNF4A, FOXA2, GRHL2) (Fig. 8A). Mutation and DNA copy  
409 number analysis identified genetic events associated with the activation of this regulatory

410 program. The Mes-like transcriptional state was often associated with *GATA4* DNA deletion and  
411 *MAPK9* amplification (Fig. 8B). This suggests copy number variation in TFs and signaling  
412 components can play a significant role as cancer driver mutations in addition to previously  
413 reported enhancer hijacking resulting from genomic structural rearrangement (Wang et al.  
414 2021).

415

416 EMT in cancer is a complex process (Lovisa 2021; Yang et al. 2020), but we found that  
417 individual GC cell lines activate the gene expression signature associated with EMT to varying  
418 degrees. DNase-seq and RNA-seq from ENCODE and GTEx provided evidence for the  
419 existence of a fibrotic phenotype in Mes-like GC through EMT (Fig. 8C & 8D). Downstream  
420 fibroblast genes such as *FGF5* and *COL1A1* were found to be upregulated in both Mes-like GC  
421 and TCGA-STAD, and we identified flanking *cis*-regulatory enhancer elements with increased  
422 activity in Mes-like GC vs. Epi-like GC as well as TCGA-STAD compared to normal stomach.  
423 Survival analysis shows that activation of Mes-like TFs and their target genes are predictive of  
424 disease severity and suggests both prognostic and therapeutic utility.

425

426 Taken together, we identified altered regulatory programs in gastric cancer along with their  
427 distinct transcriptional responses driven by differential enhancer activity. We identify copy  
428 number DNA alterations as genomic aberrations responsible for the dysregulation of these core  
429 TFs. Our findings suggest that activation of this small set of TFs driving the Mes-like GC  
430 regulatory program plays an important role in cancer progression and highlights new biology  
431 and potential therapeutic opportunities in gastric cancer.

432

## 433 Methods

### 434 ATAC-seq data processing and gkm-SVM training

435 The raw GC cell line ATAC-seq data analyzed in this study was reported in (Xu et al. 2023;  
436 Sheng et al. 2021; Xing et al. 2020; Ho et al. 2023) and uploaded to GSE264550. ATAC-seq  
437 paired reads were mapped to hg38 genome with Bowtie 2 version 2.2.5 (Langmead and  
438 Salzberg 2012). For gkm-SVM training, peaks were called by MACS2 (Zhang et al. 2008). Distal  
439 peaks were defined 300bp bins centered on the top 10k MACS2 peaks, after removing  
440 promoters (<2kb of a TSS) and peaks open in more than 30% of ENCODE Dnase-seq samples  
441 (mostly CTCF sites). Negative set genomic regions were selected randomly with matched GC  
442 and repeat content as described in (Beer et al. 2020; Shigaki et al. 2019). Training positive  
443 peaks vs. the negative sets generated robust models with a median cross-validation AUROC of  
444 0.925 (Supplemental Fig. S24). We used default settings and either the gkmSVM-R or lsgkm  
445 packages (Ghandi et al. 2014a, 2016; Lee 2016). This training procedure was shown to  
446 generate gkm-SVM models with a strong performance in both MPRA experiments and  
447 predictions of variant impact (Yan et al. 2021; Kreimer et al. 2017; Beer 2017). For ENCODE  
448 DNase-seq normal stomach we used ENCSR782SSS (Roadmap Epigenomics Consortium et  
449 al. 2015) downloaded from [encodeproject.org](https://encodeproject.org) and for TCGA stomach tumor ATAC-seq we used  
450 TCGA-BR-A4J6 (Supplemental Fig. S25, S26) (Corces et al. 2018). ENCODE DNase-seq and  
451 TCGA STAD ATAC-seq (Corces et al. 2018; Shigaki et al. 2019) were processed in a similar  
452 manner, as described in (Shigaki et al. 2019). We chose gkm-SVM for this task because of its  
453 robust performance on small (differential accessibility) datasets and because it has been shown  
454 that gapped *k*-mers are efficient representations of TFBS (Yan et al. 2021; Ghandi et al. 2014b)  
455 and protein motifs (Amanchy et al. 2011). Each gkm-SVM model generates a score function that

456 can be specified as weights for each gapped  $k$ -mer (Ghandi et al. 2014b), and the score for  
457 each sequence is the sum of weights for all gapped  $k$ -mers in a sequence. These weights  
458 quantify the contribution of each gapped  $k$ -mer to chromatin accessibility. To map these gapped  
459  $k$ -mer weights to interpretable TF activity, after training gkm-SVM models, we extracted TF  
460 models and their inferred activity using gkm-PWM, which infers the frequency of PWMs (TFBS  
461 position weight matrices) by minimizing the error between the observed and generated weight  
462 vectors. All of the 10 STAD ATAC samples with >10k distal peaks and AUROC > 0.9 detect  
463 some activation of RUNX or AP-1 when trained against normal stomach DNase-seq (ENCODE:  
464 ENCSR782SSS): 7 detect AP-1 (TCGA.CD.A48C, TCGA.VQ.A94O, TCGA.VQ.A8PJ,  
465 TCGA.BR.A4J6, TCGA.BR.A4CS, TCGA.HF.A5NB, TCGA.BR.A4J4), and 7 detect RUNX  
466 (TCGA.VQ.A94O, TCGA.VQ.A91W, TCGA.BR.A4J6, TCGA.BR.A4IY, TCGA.HF.A5NB,  
467 TCGA.BR.A4J4, TCGA.CD.A486) (Supplemental Fig. S5). For the GC cell lines, to generate a  
468 more compact set of TFBS motifs, we trained each set of positive ATAC distal peaks vs. each  
469 other (300 pairs of experiments,  $n=2000$  differentially active sequences in each of the two cell  
470 types, median AUROC = 0.922, Supplemental Fig. S27) and used the 19 most commonly  
471 detected motifs as PWMs to extract TF activity (AP-1, RUNX, AP2, NFKB, ZEB (AREB6), NFI,  
472 ONECUT, EBOX, EHF, TCF/LEF, TEAD, GRHL, HNF4A, HNF1B, HOXB13, FOX, KLF, SOX,  
473 GATA) (Fig. 1E, Supplemental Fig. S28). Although many of these motifs are detected in other  
474 cell types (McClymont et al. 2018), the combinations are specific to endodermal lineages (Luo  
475 et al. 2023). Clustering in the inferred motif activity space generates a PCA plot with the same  
476 grouping detected in the ATAC-seq signal (Fig. 1A).

## 477 ATAC-seq analysis around differentially expressed genes

478 We calculated the ATAC-seq signal in peaks within 50k around the TSS of differentially  
479 expressed genes in Epi and Mes (higher in Epi, and higher in Mes). The average activity in

480 flanking peaks averaged over all Epi, Mes, and Intermediate cell lines is shown in Fig. 6A.

481 Similar results were found using peaks within 100k, 150k, 200k, or 250k from TSS.

482

## 483 **ChIP-seq data**

484 ChIP-seq samples from ENCODE were used in the genome browser tracks: ZEB1 (ENCODE

485 ENCSR000BVN) and SMAD4 (ENCODE ENCSR826YMT). ChIP-seq data for RUNX2 and JUN

486 in Fig. 1F are from (Bevill et al. 2023). ChIP-seq data for KLF5 in AGS are from (Liu et al. 2020)

487 and ChIP-seq data for GATA4 and GATA6 in AGS are from (Chia et al. 2015) with raw data

488 uploaded to GSE51705.

## 489 **RNA-seq data processing**

### 490 **GC cell lines RNA-seq**

491 The raw GC cell line RNA-seq data analyzed in this study was reported in refs (Xu et al. 2023;

492 Sheng et al. 2021; Xing et al. 2020; Ho et al. 2023) and uploaded to GSE266159, GSE157750,

493 and GSE85465. To normalize FPKM values, we divided them by the sample's upper-quartile

494 (75th percentile) expression value and multiplied them by the average upper-quartile values

495 across all samples to scale (upper-quartile normalization) and then  $\log_2$ -transformed.

### 496 **TCGA-STAD and TCGA-Normal RNA-seq**

497 TCGA-STAD (stomach adenocarcinoma) and TCGA-Normal (normal adjacent stomach tissue)

498 RNA-seq HTSeq counts were downloaded from The Cancer Genome Atlas (TCGA) GDC portal

499 (<https://portal.gdc.cancer.gov/>) (Anders et al. 2015). 356 tumor (TCGA-STAD) and 32 normal

500 stomach (TCGA-Normal) samples were initially obtained. TCGA-STAD and TCGA-Normal

501 samples with a higher correlation to GTEx esophagus than to GTEx stomach tissue RNA-seq  
502 signature were removed (with a method similar to the GTEx correlation analysis described  
503 below). The remaining 322 TCGA-STAD and 17 TCGA-Normal samples were upper-quartile  
504 normalized over the set of protein-coding genes as labeled by GENCODE V35 annotation  
505 (Frankish et al. 2019).

506

## 507 Chromatin accessibility analysis of Mes-like GC and normal 508 ENCODE fibroblast and stomach

509 To compare chromatin accessibility profiles, we generated a union set of all distal peaks of  
510 ENCODE stomach (n=12) and primary fibroblast (n=30), and 25 GC cell lines to perform PCA.  
511 We trained gkm-SVM models on the 2k most differentially accessible peaks in all pairs of 5  
512 Mes-like GC cell lines and 30 ENCODE primary fibroblasts and ranked non-similar motifs by  
513 their average gkm-PWM Z-score.

## 514 RNA-seq analysis

### 515 Tissue-specific protein-coding genes

516 To reduce noise, RNA-seq analyses were performed using 11312 “tissue-specific” genes, which  
517 are derived from GTEx RNA-seq profiles of 54 healthy human tissues (GTEx portal:  
518 [GTEx\\_Analysis\\_2017-06-05\\_v8\\_RNASeQCv1.1.9\\_gene\\_median\\_tpm.gct.gz](http://GTEx_Analysis_2017-06-05_v8_RNASeQCv1.1.9_gene_median_tpm.gct.gz)). To define tissue-  
519 specific genes, we calculated the number of GTEx tissues, where a gene has a TPM expression  
520 above the median expression in that tissue (min = 0, max = 54). If this number is 0, the gene is  
521 lowly expressed in all tissues and if this number is 54, the gene is ubiquitously expressed (e.g.  
522 ribosomal genes). Tissue-specific genes are then defined as the subset of protein-coding genes

523 that are neither lowly expressed nor ubiquitously expressed in all GTEx healthy human tissues.

524 Protein-coding genes are found according to GENCODE V35 annotation.

## 525 PCA of RNA-seq

526 Principal Component Analysis (PCA) was performed using `prcomp()` and drawn with `ggplot2` in

527 R (R Core Team 2023). For the joint cell line and primary tissue PCA, GC cell line, and TCGA-

528 STAD FPKM expression values were upper-quartile normalized together over the set of tissue-

529 specific genes. Genes having the highest (positive or negative) weights in the PCA were

530 obtained using `R prcomp$rotation[, 1]` and `prcomp$rotation[, 2]` for PC1 and PC2, where the top

531 100 genes were used for gene set analysis. We used upper-quartile normalized RNA-seq from

532 LPS141 (Chen et al. 2019) to add this cell line to the PCA in Supplemental Fig. S6, which

533 shows that the transcriptional state of LPS141 is very similar to the Mes-like GC cell lines.

## 534 Differentially expressed genes in Mes-like vs. Epi-like cell lines

535 Differentially expressed genes between the Mes-like and Epi-like GC cell lines were found using

536 the `edgeR` (Robinson et al. 2010) package in R. The list of genes was limited to 11312 tissue-

537 specific protein-coding genes, as described above. 920 genes (including 72 TFs) were found to

538 be upregulated in Mes-like GC cell lines ( $|\log_2FC| > 2$ ,  $FDR < 0.05$ , Supplemental Table S1),

539 while 850 genes (including 51 TFs) have a higher expression in Epi-like cell lines (Supplemental

540 Table S2). The list of all 1639 human TFs was retrieved from

541 <http://humantfs.ccb.utoronto.ca/download.php> (Version 1.01) (Lambert et al. 2018).

## 542 Survival analysis in ACRG and TCGA-STAD

543 Gene expression values for 300 ACRG gastric cancer patients were downloaded from NCBI-

544 GEO GSE62254. Patients were stratified into Mes-like (n=150) and Epi-like (n=150) groups

545 based on the median difference in the mean expression values of Mes and Epi DE TFs.  
546 Similarly, TCGA-STAD patients (n=315) were split into two groups using the median expression  
547 difference value. Plots and p-values were produced using survminer  
548 (<https://github.com/kassambara/survminer>), ggsurvfit  
549 (<https://github.com/pharmaverse/ggsurvfit>), and survival (<https://github.com/therneau/survival>) R  
550 packages.

## 551 Gene set analysis

552 Gene Set Enrichment Analysis (GSEA) software (Mootha et al. 2003; Subramanian et al. 2005)  
553 with hallmark gene sets (Liberzon et al. 2015) was used to identify the biological function of  
554 differentially expressed genes in Mes-like vs. Epi-like GC (default GSEA parameters).  
555 Additionally, MSigDB (Liberzon et al. 2011) gene sets H (hallmark), C2 (curated gene sets), C5  
556 (ontology gene sets), and C6 (oncogenic gene sets) were used with top 100 PC\$rotation (see  
557 PCA of RNA-seq in Methods) to find enriched gene sets. Hypergeometric test calculated the p-  
558 values, which were then adjusted using `p.adjust(method = "fdr")` in R.

## 559 Differentially expressed genes in Mes-like GC vs. ENCODE fibroblast

560 RNA-seq data for normal human fibroblast was downloaded from the ENCODE portal and was  
561 upper-quartile normalized with the GC cell line RNA-seq profiles. Normalization and the DE  
562 gene analysis were performed similar to what described above (FDR < 0.01,  $\log_2FC > 2$ ) where  
563 982 and 5611 genes are upregulated in Mes-like GC and ENCODE fibroblast samples,  
564 respectively. 96 of the upregulated genes in the Mes-like group and 45 in the ENCODE  
565 fibroblast group are TFs. To compare the expression of these DE TFs in normal stomach and  
566 normal fibroblast, we used GTEx normal stomach and GTEx normal fibroblast to rank the

567 expression values of these DE TFs (between Mes vs. ENCODE-fibroblast) in the independent  
568 GTEx stomach and fibroblast expression profiles.

## 569 Copy number and mutation analysis

570 TCGA-STAD mutation rates were downloaded from cBioPortal (Gao et al. 2013; Cerami et al.  
571 2012) [https://www.cbioportal.org/study/summary?id=stad\\_tcga\\_pan\\_can\\_atlas\\_2018](https://www.cbioportal.org/study/summary?id=stad_tcga_pan_can_atlas_2018)). Of the  
572 TCGA-STAD samples for which the RNA-seq data was available, 255 samples are  
573 microsatellite stable (MSS) and were used for further analysis, while the microsatellite instable  
574 (MSI) samples were removed because of the possibility of having more passenger mutations.  
575 The genes selected in the plot are TFs differentially expressed between Mes-like and Epi-like  
576 GC cell lines. The x-axis is the  $\log_2$  fold-change of FPKM-UQ values between the 14 samples  
577 with a CDH1 nonsynonymous or splice-site mutation and the 241 samples without. Y-axis p-  
578 values were calculated using the *t*-test.

579

580 Copy number alteration rates for all human TFs in TCGA-STAD (Pan-Cancer) samples were  
581 downloaded from cBioPortal, and include copy number deletion and copy number amplification  
582 events as calculated by GISTIC (Mermel et al. 2011). GISTIC assigns an integer “score”  
583 between -2 and +2 to each gene depending on the number of DNA fragments aligned to it in the  
584 sequencing data. GISTIC “scores” “-2” and “-1” represent deep and shallow DNA deletions  
585 (possibly homozygous and heterozygous deletions), while “+2” and “+1” indicate deep and  
586 shallow DNA amplifications, respectively. A GISTIC “score” “0” for a gene in a sample means  
587 that the gene DNA copy number is not changed and is normal (i.e. the copy number is 2 in a  
588 healthy diploid human genome).

589

590 Whole genome sequencing (WGS) data of GC cell lines were processed as previously  
591 described (Xing et al. 2020). Briefly, copy number variations were identified using CNVkit with  
592 default parameters (bcbio-nextgen v0.9.3) (Talevich et al. 2016). As GC cell lines have no  
593 matched germline samples, CNVs were call against a non-matched normal sample. We filtered  
594 high-confidence copy number predictions (n=2061) by removing protein-coding genes whose  
595 expression correlation with copy number across cell lines is less than 0.5. We calculated copy  
596 number differences between the 15 Epi-like GC cell lines and the combined 5 Mes-like and 5  
597 Intermediate cell lines, and plotted the average copy number differences for each gene in the  
598 two groups, i.e.  $\text{mean}(\text{NonEpi CN}) - \text{mean}(\text{Epi CN})$ ). The y-axis is the p-value for the difference  
599 between Epi and NonEpi copy number scores using the Mann–Whitney *U* test. We also  
600 analyzed SNVs in the GC cell lines with SmuRF (Huang et al. 2020) and found mutations in  
601 *GATA6* and *ZFPM2* ("Friend of GATA family member 2") in 4/5 of the Mes-like cell lines, but  
602 their frequency was close to the threshold for common SNPs, and after filtering, the p-values  
603 were not significant.

604

605 To compare the non-synonymous mutation burden across groups of GC cell lines, performed  
606 variant calling on the WGS of GC cell lines using SMuRF (Huang et al. 2020). As no matched  
607 germline is available for these cell lines, we filtered out common population variants from the  
608 dbSNP database ( $AF > 1\%$ ). Then we calculated the number of protein-altering mutations  
609 (missense, truncating, start and stop codon gain/loss, and splice region variants) in each cell  
610 line (Supplemental Table S5).

611

612 To calculate the copy number burden in GC cell lines, we converted the CNVkit output to a  
613 GISTIC-like integer score (ranging from -2 to +2, described above). For each cell line, the  
614 absolute values of gene CN scores are averaged to report the cell line's overall CN status.

## 615 RNA-seq and ATAC-seq correlations

### 616 Correlation analysis using GTEx tissues RNA-seq

617 GTEx RNA-seq TPM values for 54 healthy human tissues were retrieved as described above.  
618 To compare GC cell line RNA-seq with GTEx, the cell line's FPKM values were converted to  
619 TPM. To preserve space in the figures, the TPM expression values for all GTEx brain tissues  
620 were aggregated (averaged). Mean expression for each tissue-specific gene (described above)  
621 for the three cell line groups (Mes, Epi, Intermediate) was calculated and this resulted in a  
622  $11312 * 3$  TPM expression matrix for the cell lines. Then Pearson's correlation was calculated  
623 between the GC cell line groups and GTEx tissues, using TPM expression values.

### 624 Correlation heatmap with ENCODE fibroblast and stomach RNA-seq

625 The following fibroblast and stomach RNA-seq profiles were downloaded from the ENCODE  
626 portal: HT1080 (ENCFF754UAP), bronchus fibroblast of lung (ENCFF716LRF), fibroblast skin of  
627 abdomen (ENCFF010QUB), fibroblast skin of scalp (ENCFF385POO), stomach 3-yr child  
628 (ENCFF299YCQ), stomach 37-yr adult (ENCFF683JSC), stomach 34-yr adult (ENCFF547FBP).  
629 Upper-quartile normalized FPKM values of ENCODE samples, GC cell lines, and TCGA-STAD  
630 were then compared in a heatmap (Pearson's correlation). The 11312 tissue-specific genes  
631 (described above) were used to calculate the correlation. In the scatter plot, correlation values  
632 for each individual sample are shown.

### 633 Correlation analysis with ENCODE fibroblast and stomach DNase-seq

634 The following fibroblast and stomach DNase-seq profiles were downloaded from the ENCODE  
635 portal and processed as described above: cardiac fibroblast (ENCSR000ENI), lung fibroblast  
636 (ENCSR000EPR), ELR fibroblast cell line (ENCSR240TPI), HT1080 fibrosarcoma cell line

637 (ENCSR000FDI), stomach 34-yr adult (ENCSR782SSS), stomach 54-yr adult  
638 (ENCSR163PKT), stomach 3-yr child (ENCSR246PXX). gkm-SVM models were trained on  
639 ENCODE DNase-seq, GC cell line ATAC-seq, and TCGA-STAD ATAC-seq profiles, as  
640 described above. The gkm-SVM output for each sample is a weight vector corresponding to  
641  $4^{11/2}$  *k*-mers (A/T/C/G 11-mers), where the weight value and its sign indicate how  
642 overrepresented or underrepresented the *k*-mer (TF binding site) is in the chromatin  
643 accessibility data (ATAC or DNase-seq). To find the similarities between the chromatin  
644 accessibility profiles, a heatmap (Pearson's correlation) was drawn using the ENCODE, cell  
645 lines, and TCGA-STAD samples. To calculate the correlation for each group (i.e Mes-like, Epi-  
646 like, TCGA-STAD), their mean gkm-SVM weight vectors were used for each group of samples.  
647 In the scatter plot, correlation values for each individual sample are shown. In Fig. 4D and Fig.  
648 4E, lung fibroblast (ENCSR000EPR) and adult stomach (ENCSR782SSS) were used.

## 649 scRNA-seq data processing and analysis

650 scRNA-seq raw counts for 24 patients (gastric tumor and normal adjacent tissue) were  
651 downloaded from Kim et al. study (Kim et al. 2022) on NCBI GEO (GSE150290). Cells having 1)  
652 less than 500 detected transcripts, 2) more than 20000 transcripts, 3) less than 500 detected  
653 genes, or 4) more than 10% mitochondrial or hemoglobin genes were filtered out for being low-  
654 quality. This resulted in an expression profile of a total of about 117,000 cells derived from  
655 tumor and normal adjacent tissue of 24 GC patients. Log-normalization was performed using  
656 Seurat (Hao et al. 2021) and UMAP was used for dimensionality reduction. Clusters were  
657 defined by running a Gaussian Mixture Model (GMM) over the UMAP space using the R uwot  
658 package (`n_components = 2`, `n_neighbors = 30`) (McInnes et al. 2020). To find the correlation  
659 and similarities between scRNA-seq populations (clusters) and GC cell line groups, the  
660 correlation of scRNA and bulk RNA was calculated for each cell. The average RNA-seq

661 expression for the Epi-like or Mes-like groups of cell lines over their differentially expressed (DE)  
662 genes (1770 genes) was calculated. Then for each cell, the Pearson's correlation of the scRNA-  
663 seq expression values (for 1770 DE genes) and the average bulk RNA-seq expression of Mes-  
664 like and Epi-like cell lines was calculated.

665

666 Kim et al. (Kim et al. 2022) discovered several cell populations including immune cells,  
667 fibroblasts, and various gastric cell types. We saw a similar pattern of marker genes in our  
668 analysis where clusters 1, 3, and 4 express immune-related markers such as *CD74*, *CD83*,  
669 *PTPRC (CD45)*, *IL32*, and various MHC/HLA class II variants. Clusters 2 and 5 express  
670 gastrointestinal gene markers such as *TFF1*, *TFF2*, *TFF3*, *PGC*, and *REG4* similar to what was  
671 discovered in (Kim et al. 2022), whereas cluster 6 expresses fibrotic markers including *COL1A1*,  
672 *COL1A2*, *COL3A1*, *COL6A1*, *COL6A2*, *PDGFRA*, *FN1*, and *MMP2*. This indicates that single-  
673 cell clusters 2 and 5 include stomach cells, whereas cluster 6 is comprised of fibroblasts. To  
674 define a more systematic list of highly-expressed gene markers for fibroblast (or stomach)  
675 tissue, we used the same GTEx 54 healthy tissue RNA-seq profiles as described above. We  
676 sorted the TPM expression values of 11312 tissue-specific genes (described above) in GTEx  
677 fibroblast (or stomach). We used the top 100 highly-expressed genes as the gene markers for  
678 fibroblast (or stomach). To find the biological interpretation of each single cell cluster, we used  
679 these gene markers. For every single cell, the average expression over the top 100 genes  
680 (fibroblast or stomach) was calculated and the boxplot (Fig. 5D, 5E) is drawn based on the  
681 single cell clusters, where clusters 2 and 5 were combined as they were both similarly  
682 correlated with the same group of GC cell lines. The average expression values across each  
683 cluster (in the boxplots) are compared using *t*-test, where  $p < 10^{-8}$ . Thus, gene expression in  
684 Epi-like cell lines is highly correlated to that of stomach cells, unlike Mes-like GC cell lines which  
685 have a high correlation with the fibroblast single-cell population.

686

687

## 688 Loop Extrusion Model Predictions

689 CTCF loop predictions shown in Fig. 3C, 6D, and 6E use Loop Extrusion model predictions (Xi  
690 and Beer 2021) using CTCF ChIP-seq in endoderm (Luo et al. 2023). The loop extrusion model  
691 uses orientation and binding strength information from ChIP-seq to predict loop probability, is  
692 consistent with CTCF ChIA-PET loop measurements, and generalizes better than more  
693 complicated CTCF loop prediction methods (Xi and Beer 2018).

694

## 695 Data access

696 All raw and processed sequencing data generated in this study have been submitted to the  
697 NCBI Gene Expression Omnibus (GEO; <https://www.ncbi.nlm.nih.gov/geo/>) under accession  
698 numbers GSE264550 and GSE266159. Source files for differential accessibility and expression  
699 analysis are included as Supplemental\_Code.zip and both code and gkm-SVM models are  
700 available at <http://beerlab.org/gccl>.

## 701 Competing interest statement

702 PT has stock and other ownership interests in Tempus Healthcare, previous research funding  
703 from Kyowa Hakko Kirin and Thermo Fisher Scientific, and patents/other intellectual property  
704 through the Agency for Science and Technology Research, Singapore (all outside the submitted  
705 work). All remaining authors have declared no conflicts of interest.

## 706 Acknowledgments

707 This work was supported by the following grants: MRM, DS, and MAB were supported by  
708 National Institutes of Health – National Human Genome Research Institute grant R01  
709 HG012367 to MAB. WH, YAG, SWTH, AJS and PT supported by the National Medical  
710 Research Council grant MOH-000967-00 (MOH-STaR21jun-0001) and by the National  
711 Research Foundation, Singapore, and Singapore Ministry of Health’s National Medical  
712 Research Council under its Open Fund-Large Collaborative Grant (“OF-LCG”) (MOH-  
713 OFLCG18May-0003). The cell line GES1 was a gift from Alfred Cheng from the Chinese  
714 University of Hong Kong. We thank Salvador Casan-Galdn and Brad Bernstein for providing  
715 the ChIP-seq bigwigs from (Bevill et al. 2023). We thank Roger Sik Yin Foo and Chukwuemeka  
716 George Anene-Nzeli for assisting with ATAC-seq library preparation and sequencing.  
717 *Author contributions:* Most analysis performed and designed by MRM and MAB, with  
718 contributions from AJS, WH, YAG, and DS. Overall study design conceived by MAB and AJS.  
719 All authors provided editorial and intellectual input.

720

721

722

723 

## References

- 724 Aghdassi A, Sendler M, Guenther A, Mayerle J, Behn C-O, Heidecke C-D, Friess H, Büchler M,  
725 Evert M, Lerch MM, et al. 2012. Recruitment of histone deacetylases HDAC1 and  
726 HDAC2 by the transcriptional repressor ZEB1 downregulates E-cadherin expression in  
727 pancreatic cancer. *Gut* **61**: 439–448.
- 728 Ajani JA, Lee J, Sano T, Janjigian YY, Fan D, Song S. 2017. Gastric adenocarcinoma. *Nat Rev*  
729 *Dis Primer* **3**: 1–19.
- 730 Amanchy R, Kandasamy K, Mathivanan S, Periaswamy B, Reddy R, Yoon W-H, Joore J, Beer  
731 MA, Cope L, Pandey A. 2011. Identification of Novel Phosphorylation Motifs Through an  
732 Integrative Computational and Experimental Analysis of the Human Phosphoproteome. *J*  
733 *Proteomics Bioinform* **4**: 22–35.
- 734 Anders S, Pyl PT, Huber W. 2015. HTSeq—a Python framework to work with high-throughput  
735 sequencing data. *Bioinformatics* **31**: 166–169.
- 736 Aran D, Sirota M, Butte AJ. 2015. Systematic pan-cancer analysis of tumour purity. *Nat*  
737 *Commun* **6**: 8971.
- 738 Badgwell B, Das P, Ajani J. 2017. Treatment of localized gastric and gastroesophageal  
739 adenocarcinoma: the role of accurate staging and preoperative therapy. *J Hematol*  
740 *OncolJ Hematol Oncol* **10**: 149.
- 741 Bakiri L, Macho-Maschler S, Custic I, Niemiec J, Guío-Carrión A, Hasenfuss SC, Eger A, Müller  
742 M, Beug H, Wagner EF. 2015. Fra-1/AP-1 induces EMT in mammary epithelial cells by  
743 modulating Zeb1/2 and TGF $\beta$  expression. *Cell Death Differ* **22**: 336–350.
- 744 Bass AJ, Thorsson V, Shmulevich I, Reynolds SM, Miller M, Bernard B, Hinoue T, Laird PW,  
745 Curtis C, Shen H, et al. 2014. Comprehensive molecular characterization of gastric  
746 adenocarcinoma. *Nature* **513**: 202–209.
- 747 Beer MA. 2017. Predicting enhancer activity and variant impact using gkm-SVM. *Hum Mutat* **38**:  
748 1251–1258.
- 749 Beer MA, Shigaki D, Huangfu D. 2020. Enhancer Predictions and Genome-Wide Regulatory  
750 Circuits. *Annu Rev Genomics Hum Genet* **21**: 37–54.
- 751 Berlth F, Bollschweiler E, Drebber U, Hoelscher AH, Moenig S. 2014. Pathohistological  
752 classification systems in gastric cancer: Diagnostic relevance and prognostic value.  
753 *World J Gastroenterol WJG* **20**: 5679–5684.
- 754 Bevill SM, Casaní-Galdón S, El Farran CA, Cytrynbaum EG, Macias KA, Oldeman SE, Oliveira  
755 KJ, Moore MM, Hegazi E, Adriaens C, et al. 2023. Impact of supraphysiologic MDM2  
756 expression on chromatin networks and therapeutic responses in sarcoma. *Cell*  
757 *Genomics* **3**: 100321.

- 758 Bordignon P, Bottoni G, Xu X, Popescu AS, Truan Z, Guenova E, Kofler L, Jafari P, Ostano P,  
759 Röcken M, et al. 2019. Dualism of FGF and TGF- $\beta$  Signaling in Heterogeneous Cancer-  
760 Associated Fibroblast Activation with ETV1 as a Critical Determinant. *Cell Rep* **28**: 2358-  
761 2372.e6.
- 762 Brabletz T, Kalluri R, Nieto MA, Weinberg RA. 2018. EMT in cancer. *Nat Rev Cancer* **18**: 128-  
763 134.
- 764 Byles V, Zhu L, Lovaas JD, Chmielewski LK, Wang J, Faller DV, Dai Y. 2012. SIRT1 induces  
765 EMT by cooperating with EMT transcription factors and enhances prostate cancer cell  
766 migration and metastasis. *Oncogene* **31**: 4619–4629.
- 767 Cai Y, Sablina AA. 2016. Cancer-associated chromosomal deletions: Size makes a difference.  
768 *Cell Cycle* **15**: 2850–2851.
- 769 Cerami E, Gao J, Dogrusoz U, Gross BE, Sumer SO, Aksoy BA, Jacobsen A, Byrne CJ, Heuer  
770 ML, Larsson E, et al. 2012. The cBio Cancer Genomics Portal: An Open Platform for  
771 Exploring Multidimensional Cancer Genomics Data. *Cancer Discov* **2**: 401–404.
- 772 Charafe-Jauffret E, Ginestier C, Monville F, Finetti P, Adélaïde J, Cervera N, Fekairi S, Xerri L,  
773 Jacquemier J, Birnbaum D, et al. 2006. Gene expression profiling of breast cell lines  
774 identifies potential new basal markers. *Oncogene* **25**: 2273–2284.
- 775 Chen Y, Xu L, Mayakonda A, Huang M-L, Kanojia D, Tan TZ, Dakle P, Lin RY-T, Ke X-Y, Said  
776 JW, et al. 2019. Bromodomain and extraterminal proteins foster the core transcriptional  
777 regulatory programs and confer vulnerability in liposarcoma. *Nat Commun* **10**: 1353.
- 778 Chia N-Y, Deng N, Das K, Huang D, Hu L, Zhu Y, Lim KH, Lee M-H, Wu J, Sam XX, et al. 2015.  
779 Regulatory crosstalk between lineage-survival oncogenes KLF5, GATA4 and GATA6  
780 cooperatively promotes gastric cancer development. *Gut* **64**: 707–719.
- 781 Corces MR, Granja JM, Shams S, Louie BH, Seoane JA, Zhou W, Silva TC, Groeneveld C,  
782 Wong CK, Cho SW, et al. 2018. The chromatin accessibility landscape of primary human  
783 cancers. *Science* **362**. <https://science.sciencemag.org/content/362/6413/eaav1898>.
- 784 Cristescu R, Lee J, Nebozhyn M, Kim K-M, Ting JC, Wong SS, Liu J, Yue YG, Wang J, Yu K, et  
785 al. 2015. Molecular analysis of gastric cancer identifies subtypes associated with distinct  
786 clinical outcomes. *Nat Med* **21**: 449–456.
- 787 Fang F, Chang R, Yu L, Lei X, Xiao S, Yang H, Yang L-Y. 2015. MicroRNA-188-5p suppresses  
788 tumor cell proliferation and metastasis by directly targeting FGF5 in hepatocellular  
789 carcinoma. *J Hepatol* **63**: 874–885.
- 790 Fodde R. 2002. The APC gene in colorectal cancer. *Eur J Cancer* **38**: 867–871.
- 791 Frankish A, Diekhans M, Ferreira A-M, Johnson R, Jungreis I, Loveland J, Mudge JM, Sisu C,  
792 Wright J, Armstrong J, et al. 2019. GENCODE reference annotation for the human and  
793 mouse genomes. *Nucleic Acids Res* **47**: D766–D773.

- 794 Gao J, Aksoy BA, Dogrusoz U, Dresdner G, Gross B, Sumer SO, Sun Y, Jacobsen A, Sinha R,  
795 Larsson E, et al. 2013. Integrative Analysis of Complex Cancer Genomics and Clinical  
796 Profiles Using the cBioPortal. *Sci Signal* **6**: pl1–pl1.
- 797 Gate RE, Cheng CS, Aiden AP, Siba A, Tabaka M, Lituiev D, Machol I, Gordon MG,  
798 Subramaniam M, Shamim M, et al. 2018. Genetic determinants of co-accessible  
799 chromatin regions in activated T cells across humans. *Nat Genet* **50**: 1140.
- 800 Ghandi M, Lee D, Mohammad-Noori M, Beer MA. 2014a. Enhanced Regulatory Sequence  
801 Prediction Using Gapped k-mer Features. *PLoS Comput Biol* **10**: e1003711.
- 802 Ghandi M, Mohammad-Noori M, Beer MA. 2014b. Robust k-mer frequency estimation using  
803 gapped k-mers. *J Math Biol* **69**: 469–500.
- 804 Ghandi M, Mohammad-Noori M, Ghareghani N, Lee D, Garraway L, Beer MA. 2016. gkmSVM:  
805 an R package for gapped-kmer SVM. *Bioinformatics* **32**: 2205–2207.
- 806 Guneri-Sozeri PY, Özden-Yılmaz G, Kisim A, Cakiroglu E, Eray A, Uzuner H, Karakulah G,  
807 Pesen-Okvur D, Senturk S, Erkek-Ozhan S. 2023. FLI1 and FRA1 transcription factors  
808 drive the transcriptional regulatory networks characterizing muscle invasive bladder  
809 cancer. *Commun Biol* **6**: 1–17.
- 810 Hao Y, Hao S, Andersen-Nissen E, Mauck WM, Zheng S, Butler A, Lee MJ, Wilk AJ, Darby C,  
811 Zager M, et al. 2021. Integrated analysis of multimodal single-cell data. *Cell* **184**: 3573-  
812 3587.e29.
- 813 Ho SWT, Sheng T, Xing M, Ooi WF, Xu C, Sundar R, Huang KK, Li Z, Kumar V, Ramnarayanan  
814 K. 2023. Regulatory enhancer profiling of mesenchymal-type gastric cancer reveals  
815 subtype-specific epigenomic landscapes and targetable vulnerabilities. *Gut* **72**: 226–241.
- 816 Huang W, Guo YA, Chang MM, Skanderup AJ. 2020. Ensemble-Based Somatic Mutation  
817 Calling in Cancer Genomes. In *Bioinformatics for Cancer Immunotherapy: Methods and*  
818 *Protocols* (ed. S. Boegel), *Methods in Molecular Biology*, pp. 37–46, Springer US, New  
819 York, NY [https://doi.org/10.1007/978-1-0716-0327-7\\_3](https://doi.org/10.1007/978-1-0716-0327-7_3).
- 820 Kalluri R, Zeisberg M. 2006. Fibroblasts in cancer. *Nat Rev Cancer* **6**: 392–401.
- 821 Karin M, Liu Z, Zandi E. 1997. AP-1 function and regulation. *Curr Opin Cell Biol* **9**: 240–246.
- 822 Kim J, Park C, Kim KH, Kim EH, Kim H, Woo JK, Seong JK, Nam KT, Lee YC, Cho SY. 2022.  
823 Single-cell analysis of gastric pre-cancerous and cancer lesions reveals cell lineage  
824 diversity and intratumoral heterogeneity. *Npj Precis Oncol* **6**: 1–11.
- 825 Kreimer A, Zeng H, Edwards MD, Guo Y, Tian K, Shin S, Welch R, Wainberg M, Mohan R,  
826 Sinnott-Armstrong NA, et al. 2017. Predicting gene expression in massively parallel  
827 reporter assays: A comparative study. *Hum Mutat* **38**: 1240–1250.
- 828 Lambert SA, Jolma A, Campitelli LF, Das PK, Yin Y, Albu M, Chen X, Taipale J, Hughes TR,  
829 Weirauch MT. 2018. The Human Transcription Factors. *Cell* **172**: 650–665.

- 830 Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nat Methods* **9**:  
831 nmeth.1923.
- 832 Lee D. 2016. LS-GKM: a new gkm-SVM for large-scale datasets. *Bioinformatics* **32**: 2196–2198.
- 833 Lee D, Gorkin DU, Baker M, Strober BJ, Asoni AL, McCallion AS, Beer MA. 2015. A method to  
834 predict the impact of regulatory variants from DNA sequence. *Nat Genet* **47**: 955–961.
- 835 Lei Z, Tan IB, Das K, Deng N, Zouridis H, Pattison S, Chua C, Feng Z, Guan YK, Ooi CH, et al.  
836 2013. Identification of Molecular Subtypes of Gastric Cancer With Different Responses  
837 to PI3-Kinase Inhibitors and 5-Fluorouracil. *Gastroenterology* **145**: 554–565.
- 838 Li QV, Dixon G, Verma N, Rosen BP, Gordillo M, Luo R, Xu C, Wang Q, Soh C-L, Yang D, et al.  
839 2019. Genome-scale screens identify JNK–JUN signaling as a barrier for pluripotency  
840 exit and endoderm differentiation. *Nat Genet* **51**: 999–1010.
- 841 Liberzon A, Birger C, Thorvaldsdóttir H, Ghandi M, Mesirov JP, Tamayo P. 2015. The Molecular  
842 Signatures Database Hallmark Gene Set Collection. *Cell Syst* **1**: 417–425.
- 843 Liberzon A, Subramanian A, Pinchback R, Thorvaldsdóttir H, Tamayo P, Mesirov JP. 2011.  
844 Molecular signatures database (MSigDB) 3.0. *Bioinformatics* **27**: 1739–1740.
- 845 Liu Y, Guo B, Aguilera-Jimenez E, Chu VS, Zhou J, Wu Z, Francis JM, Yang X, Choi PS, Bailey  
846 SD, et al. 2020. Chromatin Looping Shapes KLF5-Dependent Transcriptional Programs  
847 in Human Epithelial Cancers. *Cancer Res* **80**: 5464–5477.
- 848 Lonsdale J, Thomas J, Salvatore M, Phillips R, Lo E, Shad S, Hasz R, Walters G, Garcia F,  
849 Young N, et al. 2013. The Genotype-Tissue Expression (GTEx) project. *Nat Genet* **45**:  
850 580–585.
- 851 Lovisa S. 2021. Epithelial-to-Mesenchymal Transition in Fibrosis: Concepts and Targeting  
852 Strategies. *Front Pharmacol* **12**.  
853 <https://www.frontiersin.org/journals/pharmacology/articles/10.3389/fphar.2021.737570/full>  
854 l.
- 855 Lui FH, Tuan B, Swenson SL, Wong RJ. 2014. Ethnic Disparities in Gastric Cancer Incidence  
856 and Survival in the USA: An Updated Analysis of 1992–2009 SEER Data. *Dig Dis Sci*  
857 **59**: 3027–3034.
- 858 Luo J, Chen X-Q, Li P. 2019. The Role of TGF- $\beta$  and Its Receptors in Gastrointestinal Cancers.  
859 *Transl Oncol* **12**: 475–484.
- 860 Luo R, Yan J, Oh JW, Xi W, Shigaki D, Wong W, Cho HS, Murphy D, Cutler R, Rosen BP, et al.  
861 2023. Dynamic network-guided CRISPRi screen identifies CTCF-loop-constrained  
862 nonlinear enhancer gene regulatory activity during cell state transitions. *Nat Genet* **55**:  
863 1336–1346.
- 864 Luo W, Fedda F, Lynch P, Tan D. 2018. CDH1 Gene and Hereditary Diffuse Gastric Cancer  
865 Syndrome: Molecular and Histological Alterations and Implications for Diagnosis And  
866 Treatment. *Front Pharmacol* **9**.  
867 <https://www.frontiersin.org/journals/pharmacology/articles/10.3389/fphar.2018.01421>.

- 868 Luo Y, Hitz BC, Gabdank I, Hilton JA, Kagda MS, Lam B, Myers Z, Sud P, Jou J, Lin K, et al.  
869 2020. New developments on the Encyclopedia of DNA Elements (ENCODE) data portal.  
870 *Nucleic Acids Res* **48**: D882–D889.
- 871 Martinelli P, Pau ECS, Cox T, Sainz B, Duseti N, Greenhalf W, Rinaldi L, Costello E, Ghaneh  
872 P, Malats N, et al. 2017. GATA6 regulates EMT and tumour dissemination, and is a  
873 marker of response to adjuvant chemotherapy in pancreatic cancer. *Gut* **66**: 1665–1676.
- 874 Martínez-Jiménez F, Muiños F, Sentís I, Deu-Pons J, Reyes-Salazar I, Arnedo-Pac C, Mularoni  
875 L, Pich O, Bonet J, Kranas H, et al. 2020. A compendium of mutational cancer driver  
876 genes. *Nat Rev Cancer* **20**: 555–572.
- 877 McClymont SA, Hook PW, Soto AI, Reed X, Law WD, Kerans SJ, Waite EL, Briceno NJ, Thole  
878 JF, Heckman MG, et al. 2018. Parkinson-Associated SNCA Enhancer Variants  
879 Revealed by Open Chromatin in Mouse Dopamine Neurons. *Am J Hum Genet* **103**:  
880 874–892.
- 881 McInnes L, Healy J, Melville J. 2020. UMAP: Uniform Manifold Approximation and Projection for  
882 Dimension Reduction. <http://arxiv.org/abs/1802.03426>.
- 883 Mermel CH, Schumacher SE, Hill B, Meyerson ML, Beroukhi R, Getz G. 2011. GISTIC2.0  
884 facilitates sensitive and confident localization of the targets of focal somatic copy-  
885 number alteration in human cancers. *Genome Biol* **12**: R41.
- 886 Mo A, Luo C, Davis FP, Mukamel EA, Henry GL, Nery JR, Urich MA, Picard S, Lister R, Eddy  
887 SR, et al. 2016. Epigenomic landscapes of retinal rods and cones. *eLife* **5**: e11613.
- 888 Mootha VK, Lindgren CM, Eriksson K-F, Subramanian A, Sihag S, Lehar J, Puigserver P,  
889 Carlsson E, Ridderstråle M, Laurila E, et al. 2003. PGC-1 $\alpha$ -responsive genes involved in  
890 oxidative phosphorylation are coordinately downregulated in human diabetes. *Nat Genet*  
891 **34**: 267–273.
- 892 Nissen NI, Karsdal M, Willumsen N. 2019. Collagens and Cancer associated fibroblasts in the  
893 reactive stroma and its relation to Cancer biology. *J Exp Clin Cancer Res* **38**: 115.
- 894 Niu D-F, Kondo T, Nakazawa T, Oishi N, Kawasaki T, Mochizuki K, Yamane T, Katoh R. 2012.  
895 Transcription factor Runx2 is a regulator of epithelial–mesenchymal transition and  
896 invasion in thyroid carcinomas. *Lab Invest* **92**: 1181–1190.
- 897 Oh SC, Sohn BH, Cheong J-H, Kim S-B, Lee JE, Park KC, Lee SH, Park J-L, Park Y-Y, Lee H-  
898 S, et al. 2018. Clinical and genomic landscape of gastric cancer with a mesenchymal  
899 phenotype. *Nat Commun* **9**: 1777.
- 900 Park S, Lee J, Kim YH, Park J, Shin J-W, Nam S. 2016. Clinical Relevance and Molecular  
901 Phenotypes in Gastric Cancer, of TP53 Mutations and Gene Expressions, in  
902 Combination With Other Gene Mutations. *Sci Rep* **6**: 34822.
- 903 Piersma B, Hayward M-K, Weaver VM. 2020. Fibrosis and cancer: A strained relationship.  
904 *Biochim Biophys Acta BBA - Rev Cancer* **1873**: 188356.

- 905 R Core Team. 2023. R: A Language and Environment for Statistical Computing. [https://www.R-](https://www.R-project.org)  
906 [project.org](https://www.R-project.org).
- 907 Roadmap Epigenomics Consortium, Kundaje A, Meuleman W, Ernst J, Bilenky M, Yen A,  
908 Heravi-Moussavi A, Kheradpour P, Zhang Z, Wang J, et al. 2015. Integrative analysis of  
909 111 reference human epigenomes. *Nature* **518**: 317–330.
- 910 Robinson MD, McCarthy DJ, Smyth GK. 2010. edgeR: a Bioconductor package for differential  
911 expression analysis of digital gene expression data. *Bioinformatics* **26**: 139–140.
- 912 Roche J. 2018. The Epithelial-to-Mesenchymal Transition in Cancer. *Cancers* **10**: 52.
- 913 Sahai E, Astsaturov I, Cukierman E, DeNardo DG, Egeblad M, Evans RM, Fearon D, Greten  
914 FR, Hingorani SR, Hunter T, et al. 2020. A framework for advancing our understanding  
915 of cancer-associated fibroblasts. *Nat Rev Cancer* **20**: 174–186.
- 916 Schuetz CS, Bonin M, Clare SE, Nieselt K, Sotlar K, Walter M, Fehm T, Solomayer E, Riess O,  
917 Wallwiener D, et al. 2006. Progression-Specific Genes Identified by Expression Profiling  
918 of Matched Ductal Carcinomas In situ and Invasive Breast Tumors, Combining Laser  
919 Capture Microdissection and Oligonucleotide Microarray Analysis. *Cancer Res* **66**:  
920 5278–5286.
- 921 Sheng T, Ho SWT, Ooi WF, Xu C, Xing M, Padmanabhan N, Huang KK, Ma L, Ray M, Guo YA,  
922 et al. 2021. Integrative epigenomic and high-throughput functional enhancer profiling  
923 reveals determinants of enhancer heterogeneity in gastric cancer. *Genome Med* **13**: 158.
- 924 Shigaki D, Adato O, Adhikar AN, Dong S, Hawkins-Hooker A, Inoue F, Juven-Gershon T,  
925 Kenlay H, Martin B, Patra A, et al. 2019. Integration of Multiple Epigenomic Marks  
926 Improves Prediction of Variant Impact in Saturation Mutagenesis Reporter Assay. *Hum*  
927 *Mutat* **40**: 1280–1291.
- 928 Shin JY, Beckett JD, Bagirzadeh R, Creamer TJ, Shah AA, McMahan Z, Paik JJ, Sampedro  
929 MM, MacFarlane EG, Beer MA, et al. 2019. Epigenetic activation and memory at a  
930 TGFB2 enhancer in systemic sclerosis. *Sci Transl Med* **11**: eaaw0790.
- 931 Sohn BH, Hwang J-E, Jang H-J, Lee H-S, Oh SC, Shim J-J, Lee K-W, Kim EH, Yim SY, Lee  
932 SH, et al. 2017. Clinical Significance of Four Molecular Subtypes of Gastric Cancer  
933 Identified by The Cancer Genome Atlas Project. *Clin Cancer Res* **23**: 4441–4449.
- 934 Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, Paulovich A,  
935 Pomeroy SL, Golub TR, Lander ES, et al. 2005. Gene set enrichment analysis: A  
936 knowledge-based approach for interpreting genome-wide expression profiles. *Proc Natl*  
937 *Acad Sci* **102**: 15545–15550.
- 938 Talevich E, Shain AH, Botton T, Bastian BC. 2016. CNVkit: Genome-Wide Copy Number  
939 Detection and Visualization from Targeted DNA Sequencing. *PLoS Comput Biol* **12**:  
940 e1004873.
- 941 Tan IB, Ivanova T, Lim KH, Ong CW, Deng N, Lee J, Tan SH, Wu J, Lee MH, Ooi CH, et al.  
942 2011. Intrinsic Subtypes of Gastric Cancer, Based on Gene Expression Pattern, Predict  
943 Survival and Respond Differently to Chemotherapy. *Gastroenterology* **141**: 476-485.e11.

- 944 Taylor VM, Ko LK, Hwang JH, Sin M-K, Inadomi JM. 2015. Gastric Cancer in Asian American  
945 Populations: a Neglected Health Disparity. *Asian Pac J Cancer Prev* **15**: 10565–10571.
- 946 The ENCODE Consortium. 2012. An integrated encyclopedia of DNA elements in the human  
947 genome. *Nature* **489**: 57–74.
- 948 Vecchi M, Nuciforo P, Romagnoli S, Confalonieri S, Pellegrini C, Serio G, Quarto M, Capra M,  
949 Roviario GC, Contessini Avesani E, et al. 2007. Gene expression analysis of early and  
950 advanced gastric cancers. *Oncogene* **26**: 4284–4294.
- 951 Verrecchia F, Chu M-L, Mauviel A. 2001. Identification of Novel TGF- $\beta$ /Smad Gene Targets in  
952 Dermal Fibroblasts using a Combined cDNA Microarray/Promoter Transactivation  
953 Approach. *J Biol Chem* **276**: 17058–17062.
- 954 Wang X, Xu J, Zhang B, Hou Y, Song F, Lyu H, Yue F. 2021. Genome-wide detection of  
955 enhancer-hijacking events from chromatin interaction data in rearranged genomes. *Nat*  
956 *Methods* **18**: 661–668.
- 957 Xi W, Beer MA. 2018. Local epigenomic state cannot discriminate interacting and non-  
958 interacting enhancer–promoter pairs with high accuracy. *PLoS Comput Biol* **14**:  
959 e1006625.
- 960 Xi W, Beer MA. 2021. Loop competition and extrusion model predicts CTCF interaction  
961 specificity. *Nat Commun* **12**: 1046.
- 962 Xie Y, Su N, Yang J, Tan Q, Huang S, Jin M, Ni Z, Zhang B, Zhang D, Luo F, et al. 2020.  
963 FGF/FGFR signaling in health and disease. *Signal Transduct Target Ther* **5**: 1–38.
- 964 Xing M, Ooi WF, Tan J, Qamra A, Lee P-H, Li Z, Xu C, Padmanabhan N, Lim JQ, Guo YA, et al.  
965 2020. Genomic and epigenomic EBF1 alterations modulate TERT expression in gastric  
966 cancer. *J Clin Invest* **130**. <https://doi.org/10.1172/JCI126726>.
- 967 Xu C, Huang KK, Law JH, Chua JS, Sheng T, Flores NM, Pizzi MP, Okabe A, Tan ALK, Zhu F,  
968 et al. 2023. Comprehensive molecular phenotyping of ARID1A-deficient gastric cancer  
969 reveals pervasive epigenomic reprogramming and therapeutic opportunities. *Gut* **72**:  
970 1651–1663.
- 971 Yan J, Qiu Y, Ribeiro dos Santos AM, Yin Y, Li YE, Vinckier N, Nariai N, Benaglio P, Raman A,  
972 Li X, et al. 2021. Systematic analysis of binding of transcription factors to noncoding  
973 variants. *Nature* **591**: 147–151.
- 974 Yang J, Antin P, Berx G, Blanpain C, Brabletz T, Bronner M, Campbell K, Cano A, Casanova J,  
975 Christofori G, et al. 2020. Guidelines and definitions for research on epithelial–  
976 mesenchymal transition. *Nat Rev Mol Cell Biol* **21**: 341–352.
- 977 Yu K, Chen B, Aran D, Charalel J, Yau C, Wolf DM, van 't Veer LJ, Butte AJ, Goldstein T, Sirota  
978 M. 2019. Comprehensive transcriptomic analysis of cell lines as models of primary  
979 tumors across 22 tumor types. *Nat Commun* **10**: 3574.

980 Zhang Y, Liu T, Meyer CA, Eeckhoute J, Johnson DS, Bernstein BE, Nussbaum C, Myers RM,  
981 Brown M, Li W, et al. 2008. Model-based Analysis of CHIP-Seq (MACS). *Genome Biol* **9**:  
982 R137.

983

## Figure Legends

**Figure 1) Variation in chromatin accessibility profiles across GC cell lines is consistent with differential TF activity inferred from sequence-based machine learning.** (A) PCA of ATAC-seq detects variable accessibility clustered roughly into three groups of GC cell lines (Mes-like, Intermediate, Epi-like). (B) gkm-SVM training produces similar weight vectors in Mes-like, Intermediate, and Epi-like GC cell lines, resulting from the activation of similar sets of TFs. (C) Training gkm-SVM on differentially active peaks detects RUNX, AP-1, and TEAD activation in Mes-like peaks, and KLF, GATA, GRHL, FOXA, HNF4A activation in Epi-like peaks. Similar TF activation is detected when training on primary GC tumor (TCGA-STAD) vs. normal stomach. ZEB is a transcriptional repressor and the absence of its motif from open-chromatin regions shows its activity and hence, the flipped sign (red color) for ZEB activity score. (D) TFBS PWM logos (E) gkm-SVM inferred activity (dot size) of these TFs across all samples detects common patterns of activation, but additional heterogeneity within each group. (F) ChIP-seq validation experiments show that RUNX2 and AP-1 bind to Mes-high distal peaks and do not bind to Epi-high peaks in LPS141, consistent with the machine learning predictions of RUNX and AP-1 activity. On the other hand, GATA4, GATA6 and KLF5 bind to Epi-high peaks and do not bind to Mes-high peaks in AGS. AGS is one of the Epi-like GC cell lines and LPS141 is a mesenchymal liposarcoma cell line with a very similar transcriptional profile to Mes-like GC cell lines (see Supplemental Fig. S6).

**Figure 2) RNA-seq expression of GC cell lines is concordant with EMT and advanced-vs-early GC expression signatures and is consistent with TF activity inferred from chromatin accessibility.** (A) PCA of gene expression profiles is consistent with PCA from ATAC-seq. (B) Combined PCA of GC cell lines with TCGA-STAD (tumor) and TCGA-normal shows shared variation between Mes-like vs. Epi-like and tumor vs. normal. (C,D) The genes explaining the variation (top 100 PCA gene weights) in the direction of PC1 in (A) and PC2 in (B) both are most enriched for genes highly expressed in EMT and advanced-vs-early GC. (E) Examples of marker gene expression for EMT in cancer (higher in Mes, red; higher in Epi, green), with likelihood ratio test,  $FDR < 0.05$ , and  $|\log_2FC| > 2$ . (F) Differentially expressed TFs between Mes-like and Epi-like GC cell lines, with likelihood ratio test  $FDR < 0.05$  and  $|\log_2FC| > 2$ . (G) Survival analysis based on the expression of Mes and Epi DE TFs in the ACRG cohort.

**Figure 3) Mutations and copy number variants are associated with distinct Mes-like and Epi-like TF activation in GC.** (A) Nonsynonymous *CDH1* mutation status separates TFs with differential expression between Mes-like and Epi-like cell lines. Shown is TCGA-STAD expression in tumors with and without a *CDH1* nonsynonymous mutation ( $t$ -test  $p < 0.05$  and  $|\log_2FC| > 0.25$ ). (B) Human TFs ranked in order of DNA copy number (CN) alteration rate in TCGA-STAD primary tumors. (C) Chromatin accessibility landscape and TF binding of ZEB1 and SMAD4 in the *GATA6* locus. *GATA6\_E1* and *GATA6\_E2* are putative Epi-like and normal stomach enhancers, where *GATA6* binds in *GATA6* ChIP-seq in AGS GC cell line. (D) *GATA4* has reduced CN in non-Epi and *MAPK9* has amplified CN in non-Epi cell lines. Average CN differences between Epi GC cell lines and non-Epi (Mes and Intermediate) GC cell lines are shown for protein-coding genes whose expression and copy number correlate. CN values are obtained by DNA sequencing, where “-2” marks deep DNA deletion, “0” shows a normal diploid genome, and “+2” is for deep DNA amplification in the gene locus. Mann–Whitney  $U$  test  $p < 0.05$ . (E) CN variation is consistent with altered *GATA4* and *MAPK9* expression across samples.  $r$  is Pearson’s correlation between CN scores and gene expression.

**Figure 4) In both expression and chromatin accessibility, Mes-like GC cell lines are more similar to fibroblasts, and Epi-like GC cell lines are more similar to stomach.** (A) Mes-like expression is most similar to GTEx fibroblasts. Correlation of 54 healthy GTEx gene expression profiles for each group of GC cell lines in ~11,300 tissue-specific genes. (B) Correlation of ENCODE stomach and fibroblast gene expression profiles with average gene expression for each group of GC cell lines. HT1080 is a fibrosarcoma cell line. (C) Mes-like chromatin accessibility is also most similar to ENCODE fibroblast, and TCGA-STAD and Epi-like GC cell lines are most similar to ENCODE stomach. We measured the similarity of chromatin accessibility by correlation of gkm-SVM weight vectors trained on distal enhancers of each sample. ELR is a fibroblast-derived cell line. (D) Correlation of GC cell line gene expression profiles with a healthy ENCODE lung fibroblast and a healthy stomach sample. (E) Correlation of GC cell line chromatin accessibility profiles (gkm-SVM weights trained on cell line ATAC-seq) with healthy ENCODE lung fibroblast and stomach (gkm-SVM trained using DNase-seq).

**Figure 5) Mes-like GC is most similar in expression profile to the single-cell fibroblast population of both tumor and normal stomach tissue.** (A) UMAP of scRNA-seq of gastric tumor and adjacent normal stomach tissue from Kim et al 2022. Six cell populations are detected; 2 and 5 are stomach and 6 is fibroblast. (B) Epi-like cell line gene expression profiles are most highly correlated with the stomach clusters 2 and 5. (C) Mes-like gene expression profiles are most highly correlated with the fibroblast cluster, 6. (D) Highly expressed stomach-specific genes are expressed only in stomach clusters 2 and 5, combined due to similarity. ( $t$ -test  $p < 10^{-8}$  for *cluster 2+5* vs. other clusters). (E) Highly expressed fibroblast-specific genes are expressed only in fibroblast cluster 6 ( $t$ -test  $p < 10^{-9}$  for *cluster 6* vs. other clusters.)

**Figure 6) Mes and Epi upregulated genes are flanked by Mes-active and Epi-active distal enhancers.** (A) Highly expressed genes in Epi-like GC cell lines, are flanked by distal peaks with increased accessibility in Epi-like cell lines. (B) Highly expressed genes in Mes-like GC cell lines, are flanked by distal peaks with increased accessibility in Mes-like cell lines. All peaks within 50kb of TSS of DE genes were averaged. Many of the upregulated genes in Mes cell lines are members of the TGF $\beta$ /SMAD and FGF pathways, and their expression is shown in (C) (Mann–Whitney  $U$  test  $p < 0.05$ ). (D) In the *COL1A1* locus, *COL1A1\_E1* and *COL1A1\_E2* are putative distal enhancers with increased ATAC-seq signal in Mes-like GC cell lines and TCGA-STAD relative to Epi-like and normal stomach, consistent with higher *COL1A1* expression in Mes-like GC cell lines and TCGA STAD. ChIP-seq experiments in the mesenchymal liposarcoma cell line LPS141 identified AP-1 and RUNX2 binding in *COL1A1\_E1* and *COL1A1\_E2*. (E) Similarly, in the *FGF5* locus, *FGF5\_E* is a putative distal enhancer with increased ATAC-seq signal in Mes-like GC cell lines and TCGA-STAD relative to Epi-like and normal stomach, consistent with higher *FGF5* expression in Mes-like GC cell lines and TCGA STAD. *FGF5\_E* is a RUNX2 and AP-1 binding site in LPS141. *FGF5\_E*, *COL1A1\_E1* and *COL1A1\_E2* are contained in CTCF loops enclosing the target gene promoter.

**Figure 7) Despite some similarity in gene expression and enhancer activity, Mes-like GC and fibroblast have distinct TF responses.** All groups of cell lines have similar levels of (A) somatic mutations (A) and (B) copy number alterations. (C) PCA of chromatin accessibility of GC cell lines, stomach, and fibroblasts shows that while all groups are distinct, there is an axis along which differences in accessibility between stomach and fibroblasts are shared between Epi-like and Mes-like cell lines. Epi-like cell lines are more similar to stomach, and Mes-like cell lines are more similar to but still distinct from, primary fibroblasts. (D) Mes-like vs fibroblast accessibility differences are driven by differential activity of a small set of TFs. (E) Rank plot of the average TF motif weights in all pairwise comparisons between Mes-like cell lines and ENCODE fibroblast samples. (F) Differential expression of TFs between stomach and fibroblast shows that Mes-like cell lines retain high expression of stomach-specific TFs and do not upregulate all fibroblast TFs (especially *TWIST2*) to the same degree. Ranking of DE TFs between Mes-like and ENCODE-fibroblast (FDR < 0.01) in normal GTEx stomach (x-axis) and normal GTEx fibroblast profiles (y-axis) are shown. Some of the DE TFs upregulated in Mes-like GC (*FOXA1*, *KLF5*, etc.), have a higher expression (lower ranking) in GTEx normal stomach, while *TWIST2* has a higher expression in ENCODE fibroblast and a higher expression (lower ranking) in GTEx fibroblast. (G) Inferred TF motif activities are supported by differential expression of stomach TFs and fibroblast genes in Mes-like GC, normal ENCODE-fibroblast, and ENCODE-stomach (\*\* FDR < 0.01, \* FDR < 0.05, likelihood ratio test).

**Figure 8) Summary of TFs driving GC epigenomic heterogeneity and EMT.** (A) Differentially active TFs in Mes-like and Epi-like GC cell lines. (B) DNA copy number changes in Mes-like and Epi-like GC cell lines. (C) Mes-like GC has a high correlation of expression and distal enhancer activity with fibroblast, whereas Epi-like GC is more similar to stomach. (D) Mes-like GC TF response and chromatin profile are different from that of fibroblast, despite the shared patterns of expression and chromatin accessibility.

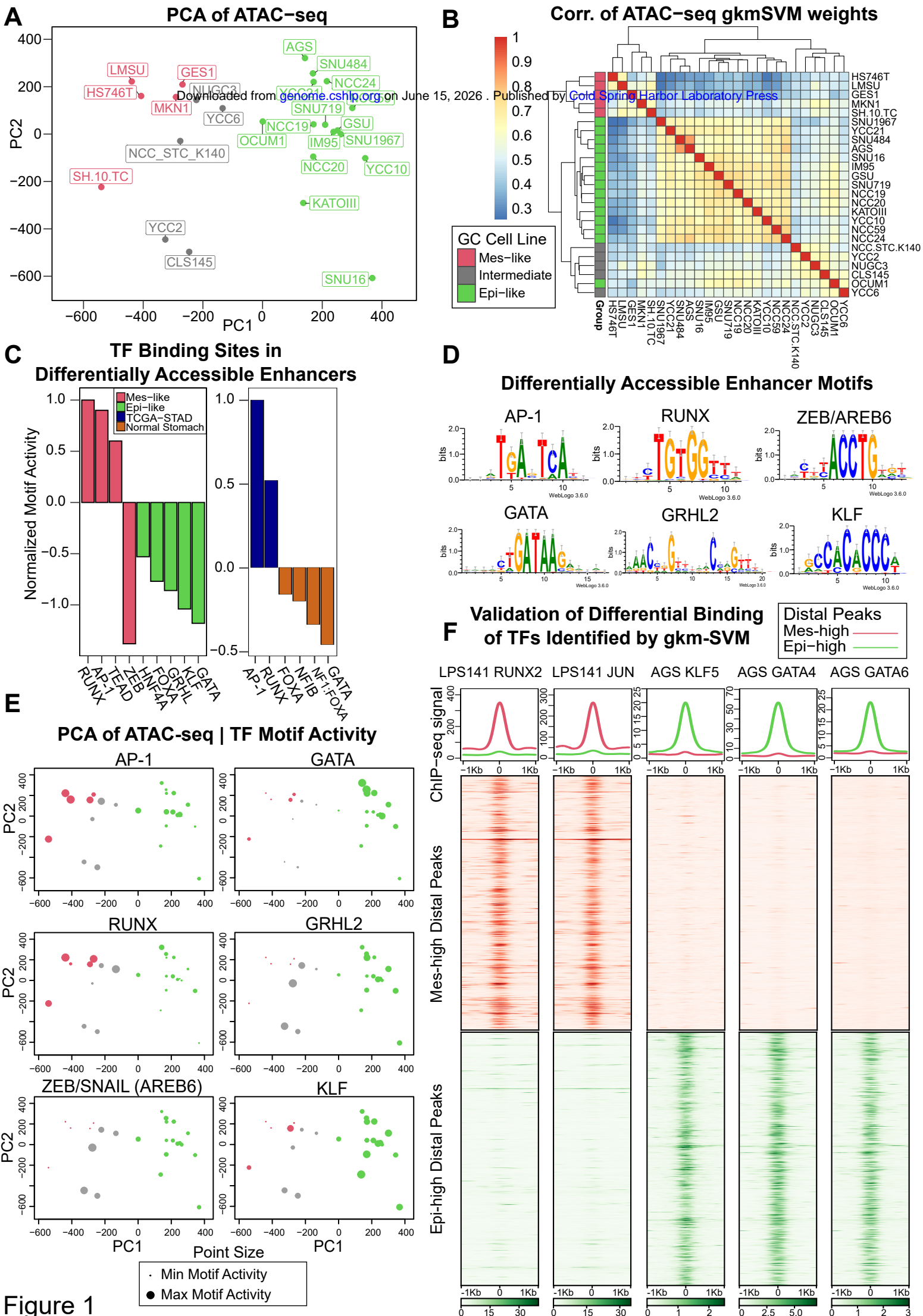


Figure 1

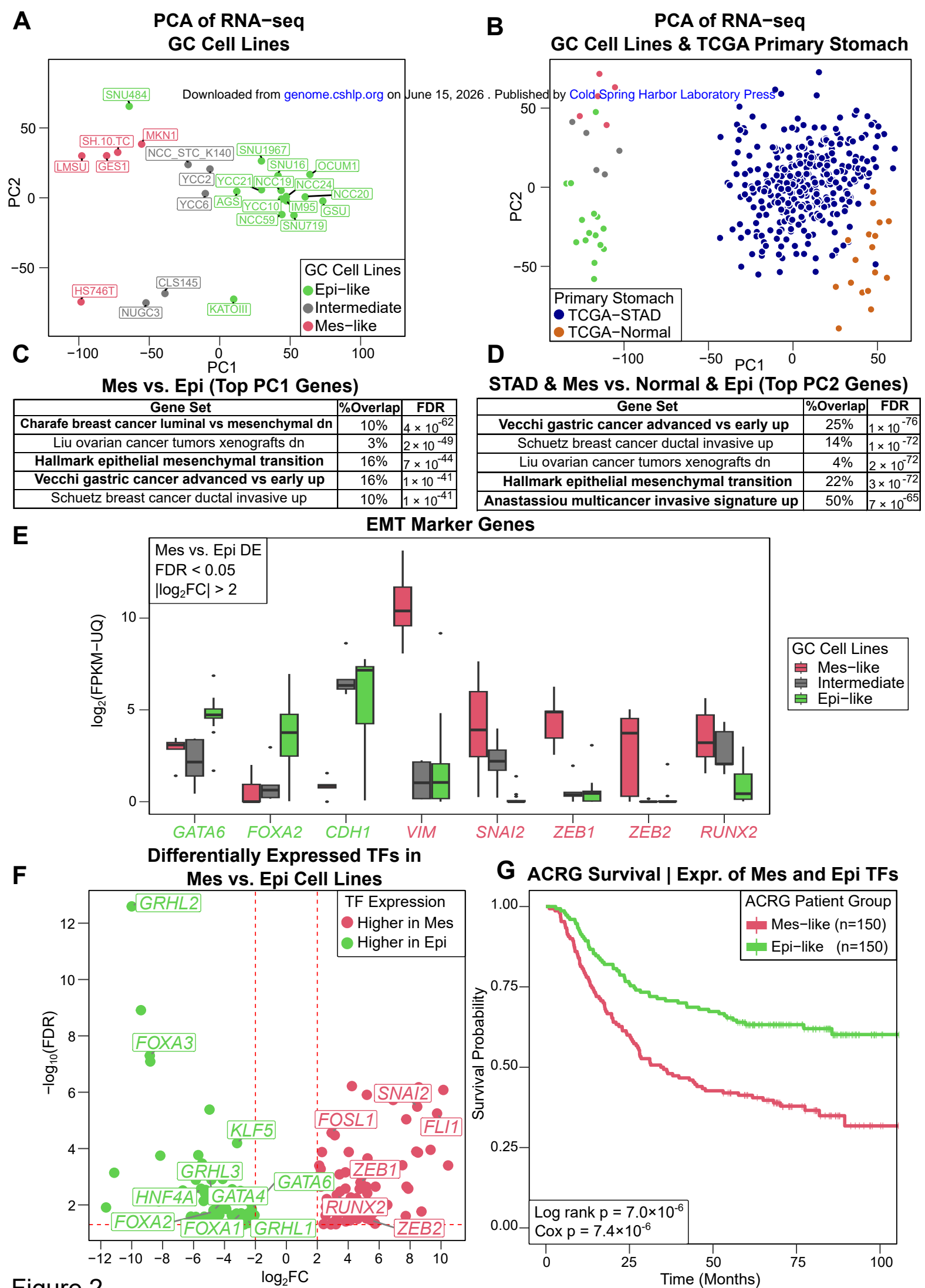


Figure 2

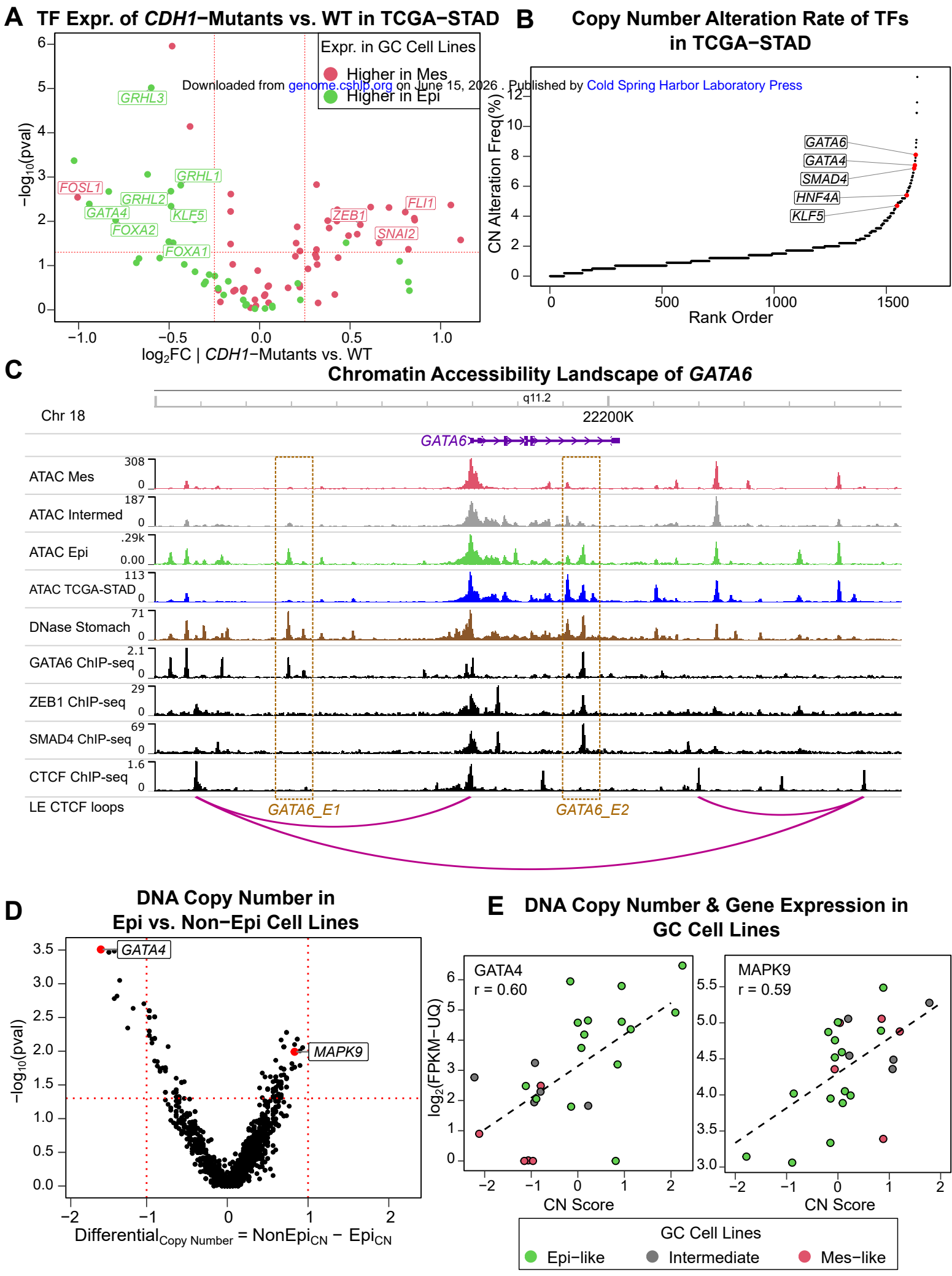


Figure 3

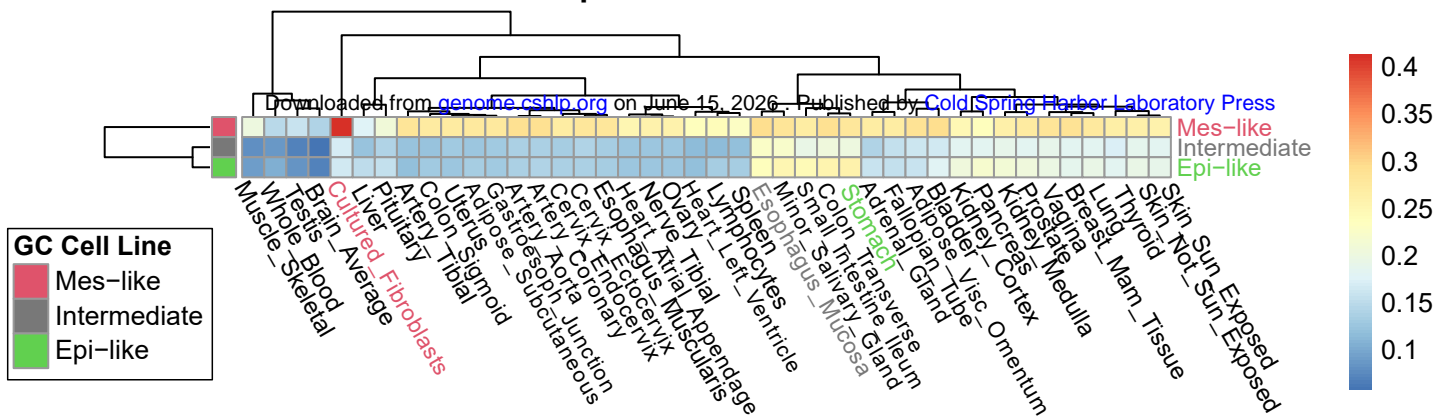
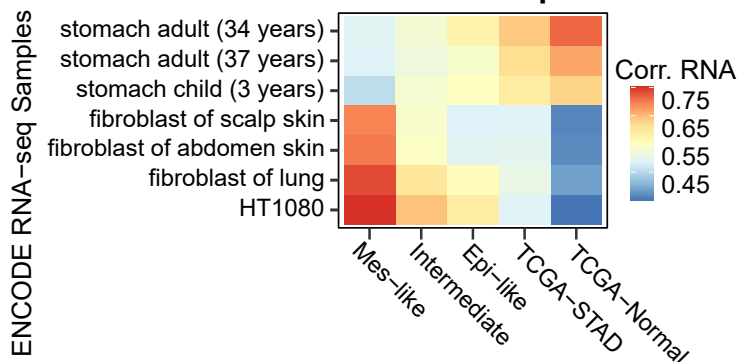
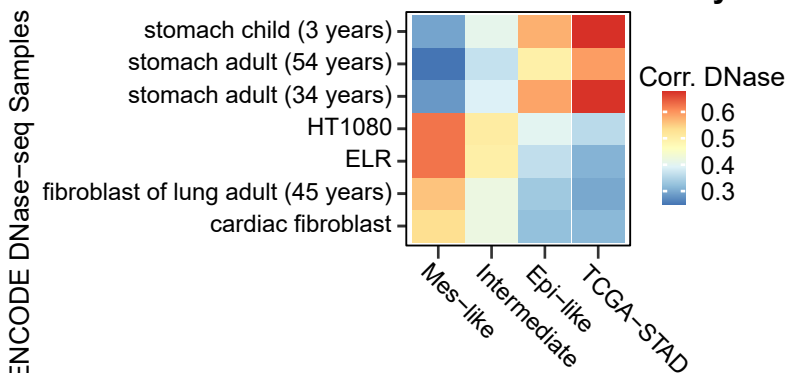
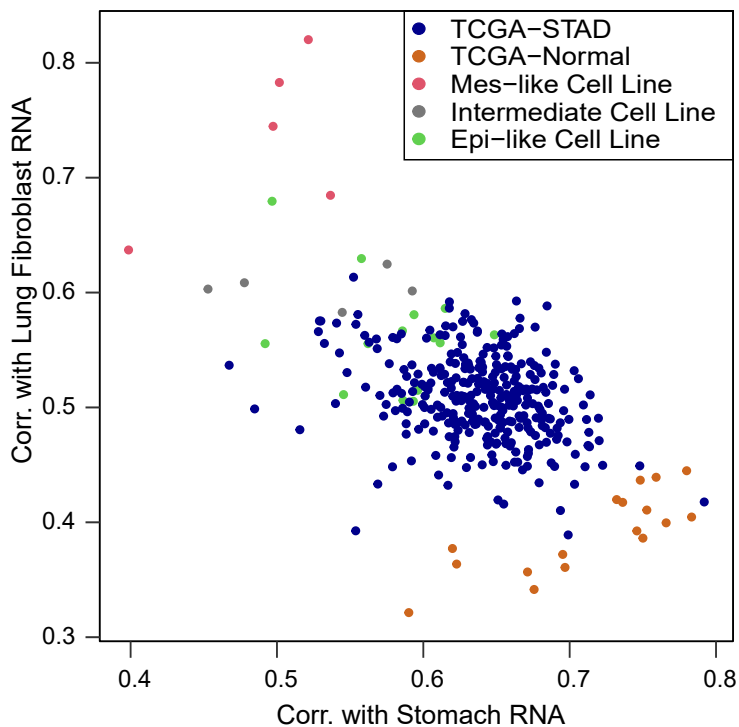
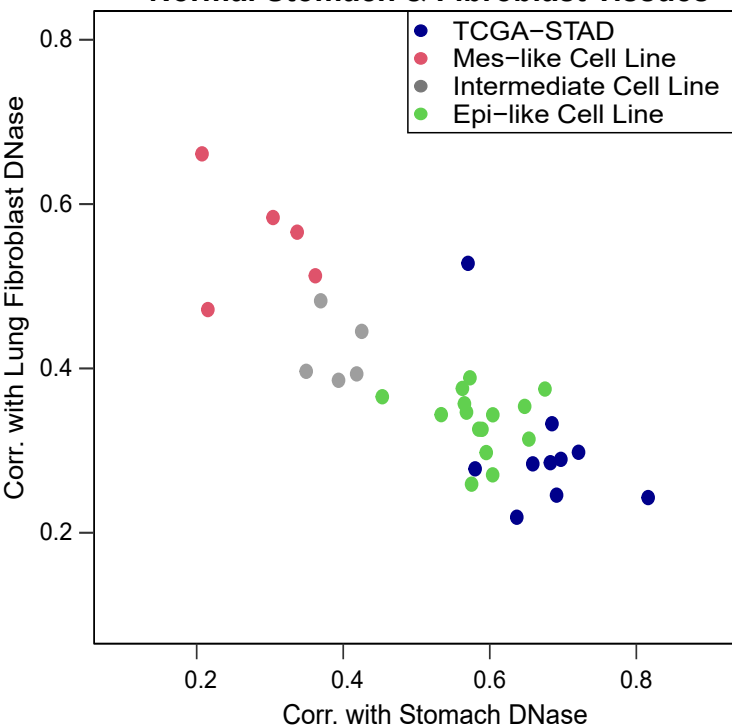
**A****Correlation of Gene Expression in GC Cell Lines and GTEx Tissues****B****Correlation of Gene Expression****C****Correlation of Chromatin Accessibility****D****Correlation of Gene Expression with Normal Stomach & Fibroblast Tissues****E****Correlation of Chromatin Accessibility with Normal Stomach & Fibroblast Tissues**

Figure 4

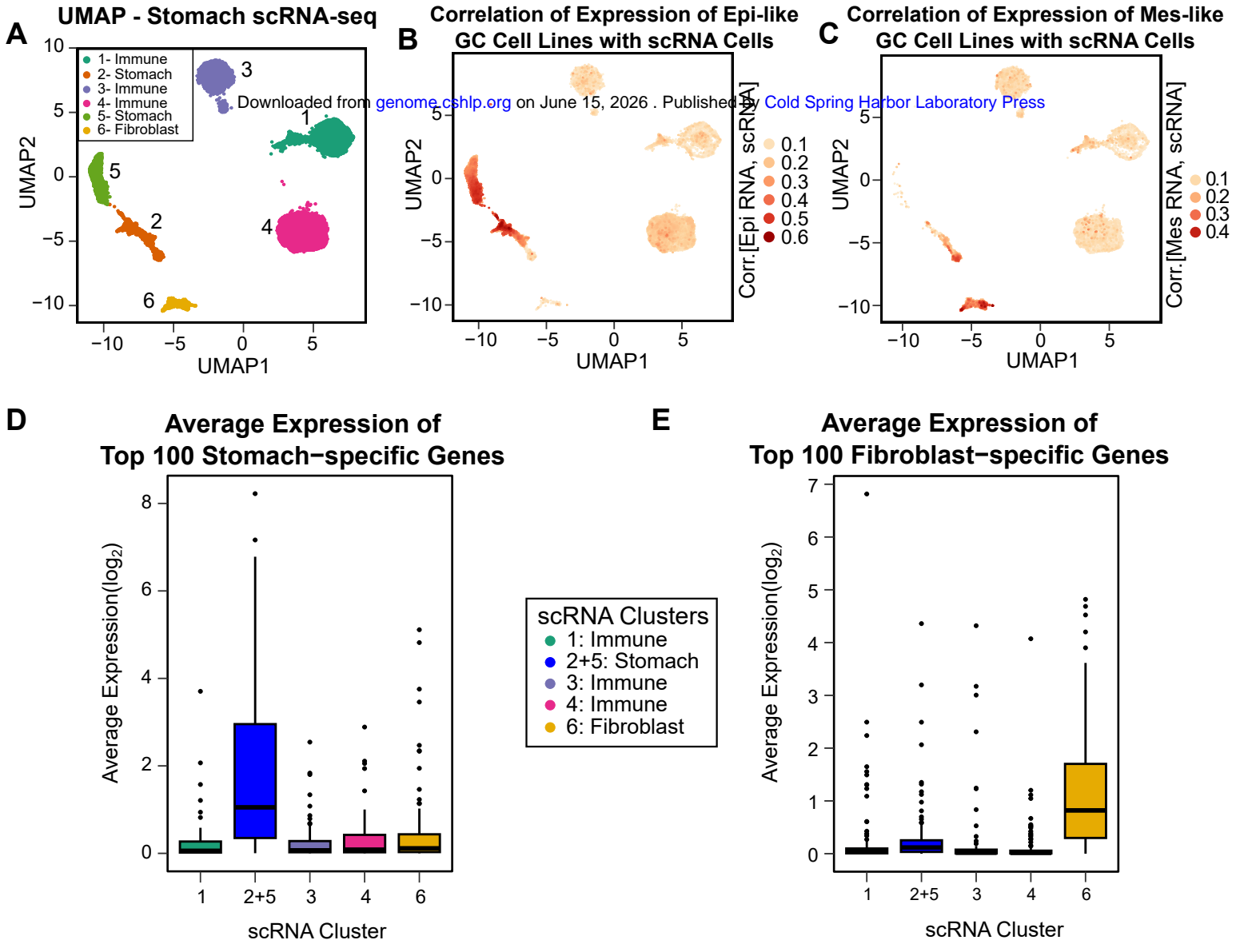
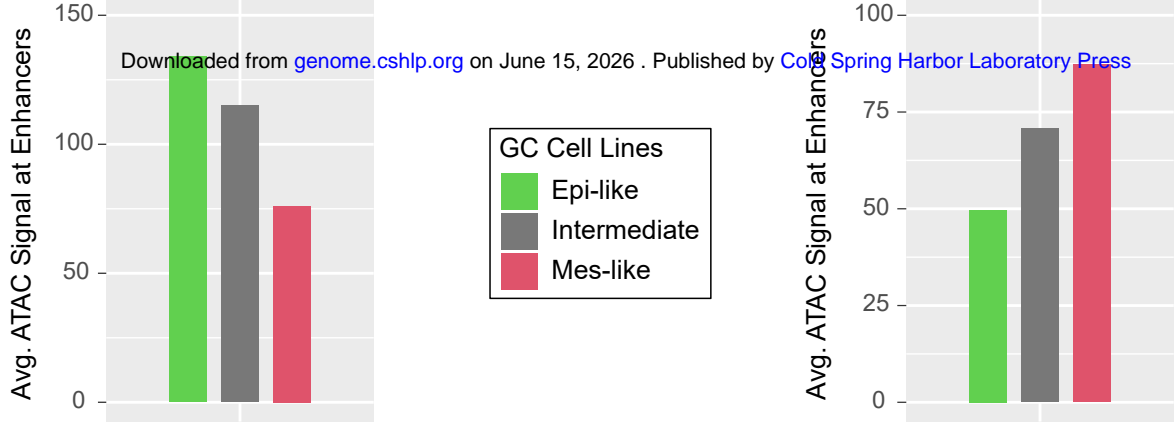


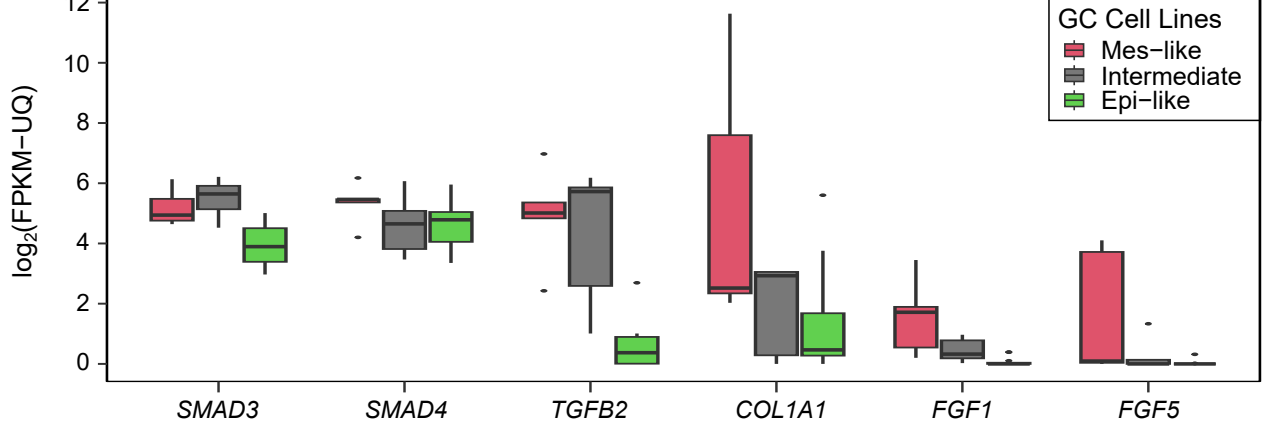
Figure 5

**A** Chromatin Accessibility in Distal Enhancers of Highly-expressed Genes in Epi-like GC

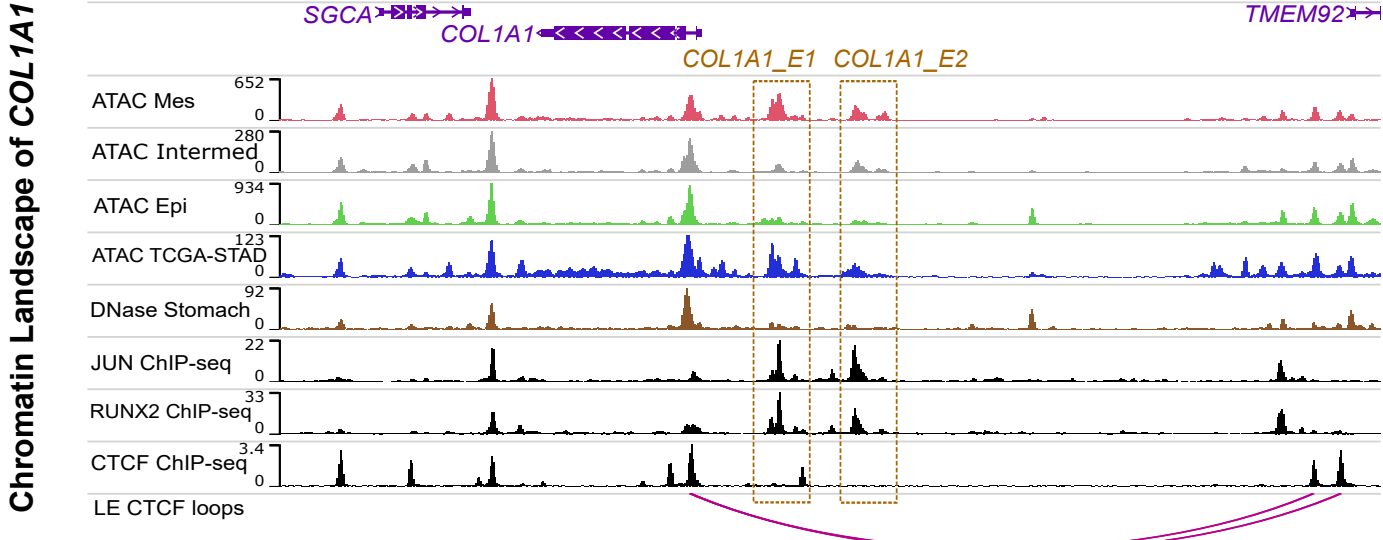
**B** Chromatin Accessibility in Distal Enhancers of Highly-expressed Genes in Mes-like GC



**C** Upstream and Downstream Genes in TGFB/SMAD and FGF Pathways



**D** Chromatin Landscape of COL1A1



**E** Chromatin Landscape of FGF5

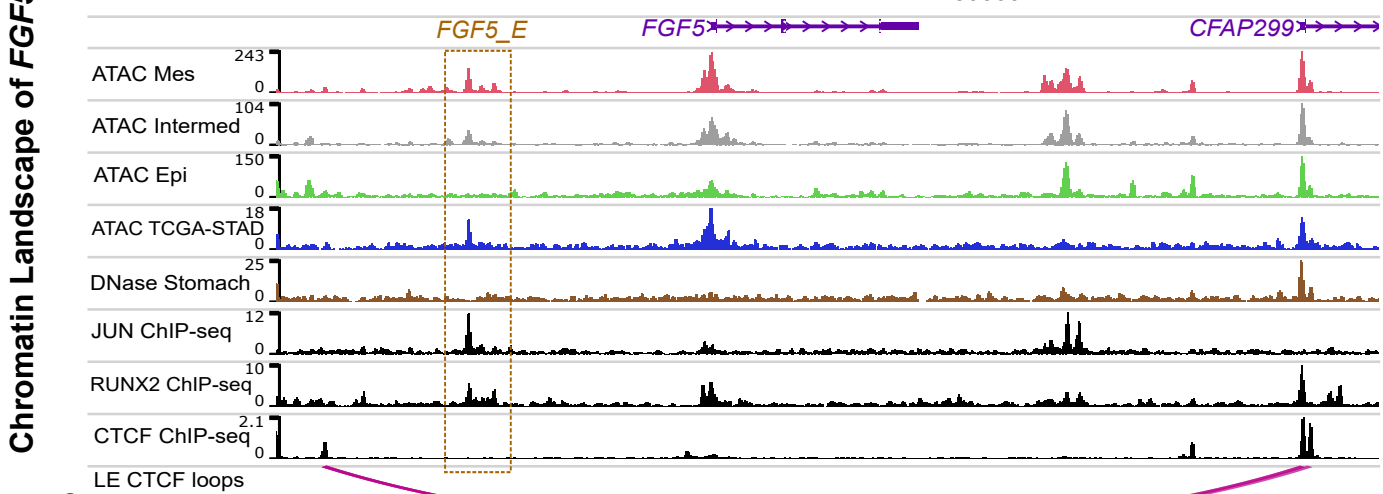
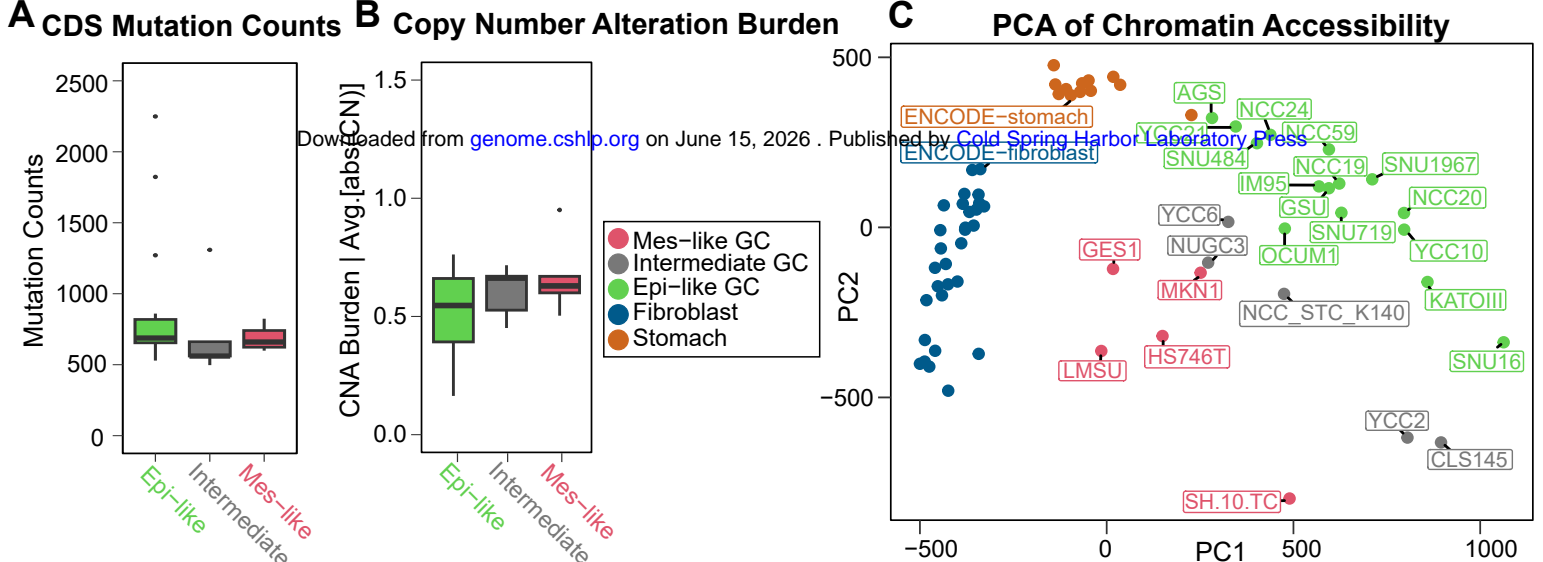


Figure 6

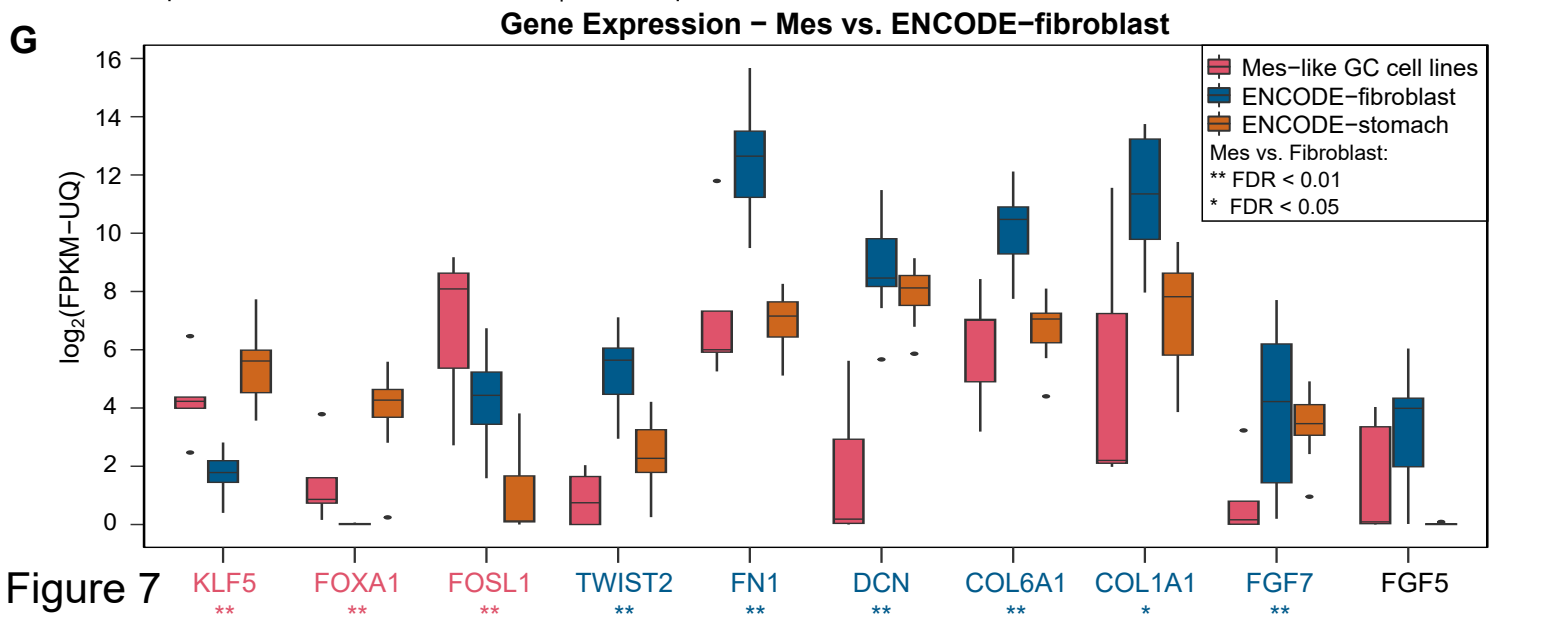
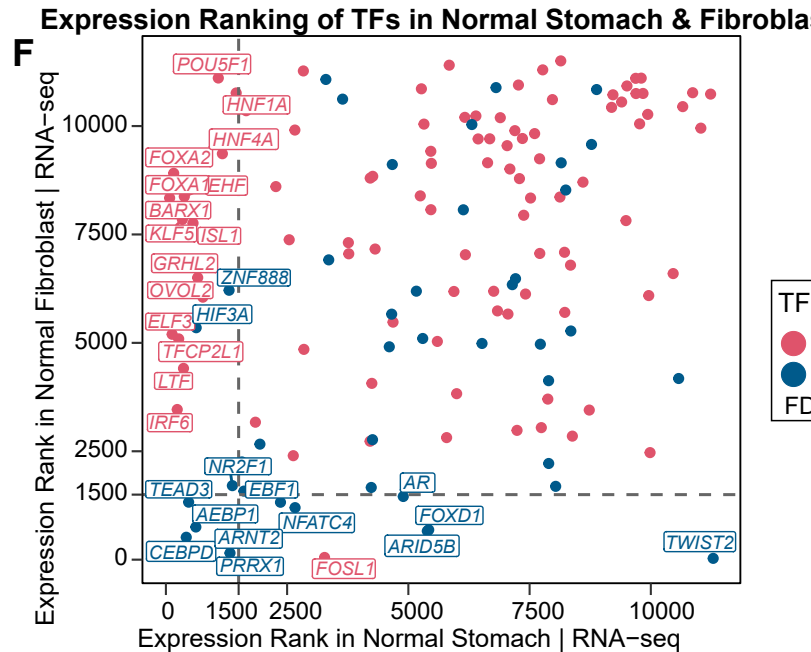
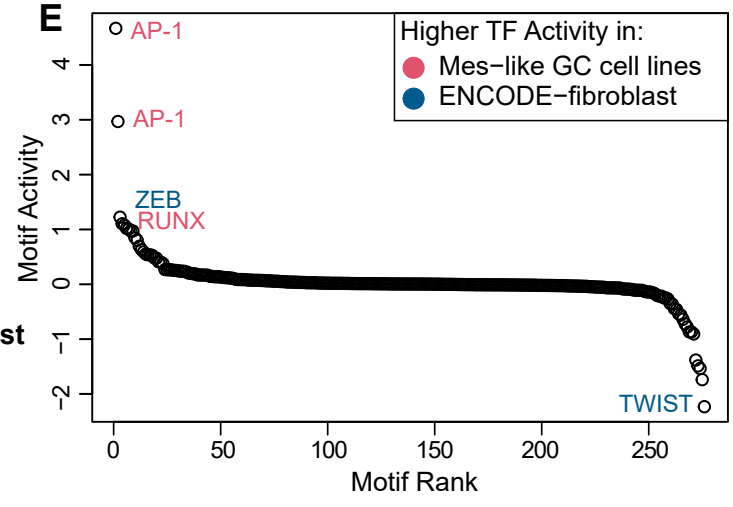


**TFBS Enriched in Differentially Accessible Distal Enhancers** **TF Motif Ranks in Mes GC vs. Normal Fibroblast**

**D Mes (ATAC) vs. Fibroblast (DNase)**

N	Motif	Z	Z
1	AP-1	4.67	1.00
2	ZEB	1.22	0.26
3	RUNX	1.10	0.24
4	GATA	1.08	0.23
5	KLF	0.97	0.21
6	RARA	0.85	0.18
7	AP-2g	0.69	0.15
8	TWIST	-2.23	0.48

Higher Accessibility in Mes-like (top 6 motifs)  
Higher Accessibility in ENCODE Fibroblast (motif 8)



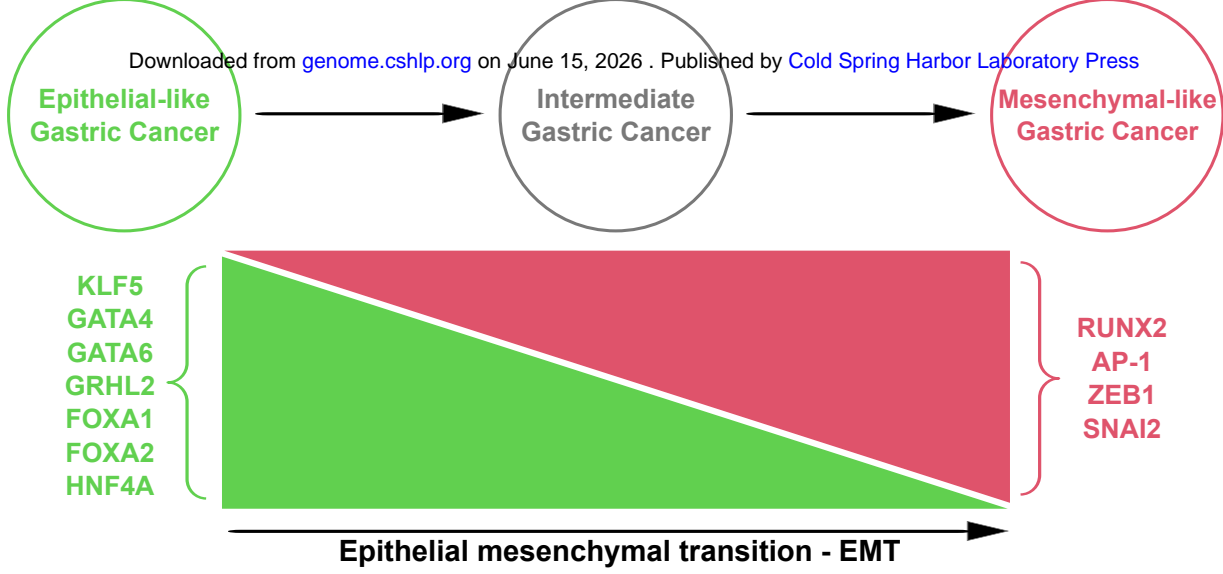
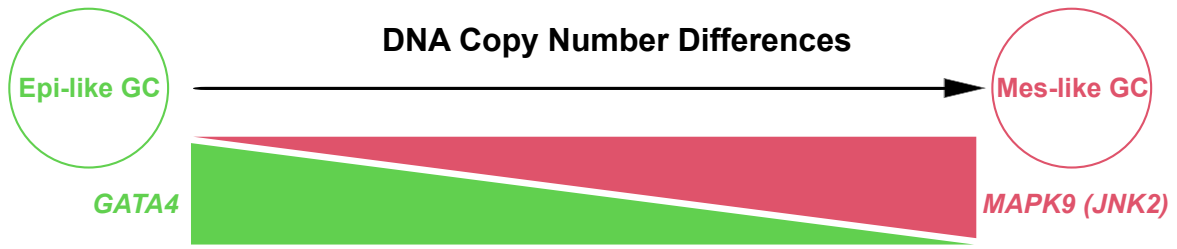
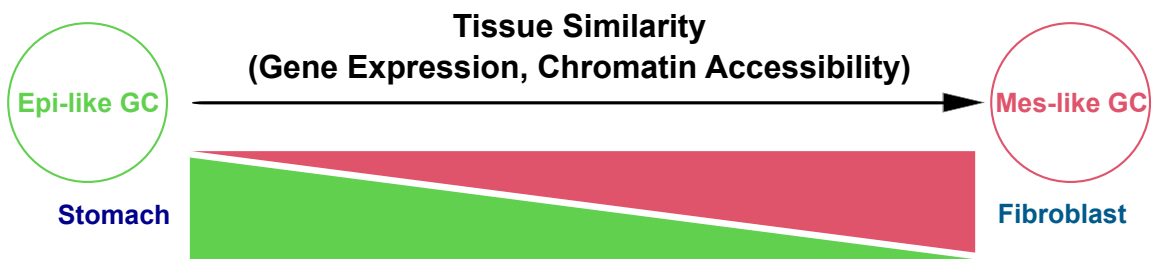
**A****Transcription Factor Activity Driving EMT in Gastric Cancer****B****C****D**

Figure 8