

1 **TITLE PAGE**

2 Full Title: Epigenetic misactivation of a distal developmental enhancer cluster drives
3 SOX2 overexpression in breast and lung cancer

4

5 Author List and Affiliations

6 Luis E. Abatti¹, Patricia Lado-Fernández^{2,3}, Linh Huynh⁴, Manuel Collado², Michael M.
7 Hoffman^{4,5,6,7}, and Jennifer A. Mitchell¹

8

9 ¹Department of Cell and Systems Biology, University of Toronto, Toronto, Ontario,
10 Canada

11 ²Laboratory of Cell Senescence, Cancer and Aging, Health Research Institute of
12 Santiago de Compostela (IDIS), Xerencia de Xestión Integrada de Santiago
13 (XXIS/SERGAS), Santiago de Compostela, Spain

14 ³Department of Physiology and Center for Research in Molecular Medicine and Chronic
15 Diseases (CiMUS), Universidade de Santiago de Compostela, Santiago de
16 Compostela, Spain

17 ⁴Princess Margaret Cancer Centre, University Health Network, Toronto, Ontario,
18 Canada

19 ⁵Department of Medical Biophysics, University of Toronto, Toronto, Ontario, Canada

20 ⁶Department of Computer Science, University of Toronto, Toronto, Ontario, Canada

21 ⁷Vector Institute for Artificial Intelligence, Toronto, Ontario, Canada

22

23 Corresponding authors:

24 ja.mitchell@utoronto.ca, luis.abatti@mail.utoronto.ca

25 Current address:

26 Department of Cell and Systems Biology, University of Toronto, Toronto, Canada

27 **ABSTRACT**

28 Enhancer reprogramming has been proposed as a major source of gene expression
29 dysregulation during tumorigenesis. Here, we identify *SOX2* developmental enhancers
30 that are misactivated in breast and lung carcinoma. Deletion of the SRR124–134
31 enhancer cluster disrupts *SOX2* transcription and genome-wide chromatin accessibility
32 in cancer cells. ATAC- and RNA-seq analysis of primary tumors shows that chromatin
33 accessibility at this cluster is correlated with *SOX2* overexpression in breast and lung
34 cancer. We further identify *FOXA1* as an activator and *NFIB* as a repressor of SRR124–
35 134 activity and *SOX2* transcription. Notably, the conserved SRR124 and SRR134
36 regions are essential during mouse development, where homozygous deletion results in
37 the lethal failure of esophageal-tracheal separation. Our findings indicate that the
38 SRR124–134 enhancer cluster drives *SOX2* expression during development. In breast
39 and lung cancer, *FOXA1*-induced aberrant activity of the SRR124–134 cluster drives
40 *SOX2* overexpression, demonstrating how developmental enhancers can be
41 recommissioned during tumorigenesis. These results highlight the importance of
42 understanding enhancer dynamics during development and disease while also
43 providing new opportunities for therapeutic intervention by targeting aberrantly activated
44 developmental enhancers.

45 INTRODUCTION

46 Multicellular organisms have a deeply organized hierarchy of distinct cell lineages
47 originally derived from the same embryonic progenitor. The differentiation of early
48 embryonic cells involves the activation of specific transcriptional programs that drive
49 their commitment toward distinct cell types. This process is mediated by key
50 developmental-associated transcription factors ^{reviewed in 1}, which interact genome-wide
51 with cis-regulatory regions and are responsible for progressively restricting the
52 epigenome, repressing regulatory regions associated with pluripotency ^{2,3} and activating
53 enhancers that control the expression of lineage-specific genes ⁴⁻⁶. This establishes an
54 epigenetic regulatory “memory” that maintains cells in their own lineage compartment,
55 reinforcing their transcriptional programs and repressing previous uncommitted states ⁷.
56 This epigenetic landscape, however, becomes profoundly disturbed during
57 tumorigenesis ⁸⁻¹⁰, causing cells to lose their identity and assume a dysfunctional state
58 that combines regulatory features from other cell lineages in a mechanism known as
59 “enhancer reprogramming” ^{7-9,11}. Although this has been proposed as one of the
60 sources of transcriptional dysregulation, it remains unclear if early developmental
61 enhancers that were decommissioned during lineage differentiation are reactivated in
62 the disease state.

63 SRY-box transcription factor 2 (SOX2) is a pioneer transcription factor required for
64 pluripotency maintenance in embryonic stem cells ^{12,13} and reprogramming to induced
65 pluripotent stem cells in mammals ¹⁴⁻¹⁶. Two proximal enhancers were once deemed
66 crucial for driving *Sox2* expression during early development: Sox2 Regulatory Region 1
67 (SRR1) and SRR2 ¹⁷⁻¹⁹. Deletion of SRR1 and SRR2, however, has no effect on *Sox2*

68 expression in mouse embryonic stem cells ²⁰. In contrast, deletion of a distal Sox2
69 Control Region (SCR), 106 kb downstream of the Sox2 promoter, causes a profound
70 loss of Sox2 expression in mouse embryonic stem cells ^{20,21} and in blastocysts, where
71 SCR deletion causes peri-implantation lethality ²². However, the contribution of these
72 regulatory regions in driving SOX2 expression in other contexts remains poorly
73 understood.

74 SOX2 is also involved in tissue morphogenesis and homeostasis of the brain ²³,
75 eyes ²⁴, esophagus ²⁵, inner ear ²⁶, lungs ²⁷, skin ²⁸, stomach ²⁹, taste buds ³⁰ and
76 trachea ³¹ in both human and mouse. In these tissues, SOX2 expression is regulated
77 precisely in space and time at critical stages of development. For example, proper
78 levels of Sox2 expression are required for the complete separation of the anterior
79 foregut into the esophagus and trachea in mice ^{22,25,32} and in humans ^{33,34}, as the
80 disruption of Sox2 expression leads to an abnormal developmental condition known as
81 esophageal atresia with distal tracheoesophageal fistula (EA/TEF) ^{reviewed in 35,36}. After
82 the anterior foregut is properly separated, Sox2 expression ranges from the esophagus
83 to the stomach in the gut ^{25,29}, and throughout the trachea, bronchi, and upper portion of
84 the lungs in the developing airways ³¹. Proper branching morphogenesis at the tip of the
85 lungs, however, requires temporary downregulation of Sox2, followed by reactivation
86 after lung bud establishment ²⁷. Sox2 also retains an essential function in multiple
87 mature epithelial tissues, where it is highly expressed in proliferative and self-renewing
88 adult stem cells necessary for maintaining and replacing terminally differentiated cells
89 within the epithelium of the brain, bronchi, esophagus, stomach, and trachea ^{29,31,37,38}.

90 Tumorigenesis of at least 25 different cancer types involves SOX2 overexpression
91 reviewed in ³⁹. This overexpression is linked to increased cellular replication rates, more
92 aggressive tumor grades, and poor patient outcomes in breast carcinoma (BRCA) ^{40–44};
93 colon adenocarcinoma (COAD) ^{45–48}; glioblastoma (GBM) ^{49–52}; liver hepatocellular
94 carcinoma (LIHC) ⁵³; lung adenocarcinoma (LUAD) ^{54–56}; and lung squamous cell
95 carcinoma (LUSC) ^{57,58}. These clinical and molecular characteristics arise from the
96 participation of SOX2 in the formation and maintenance of tumor-initiating cells that
97 resemble tissue progenitor cells, as evidenced by BRCA ^{44,59,60}, GBM ^{51,61–63}, LUAD ⁶⁴,
98 and LUSC ⁶⁵ studies. SOX2 knockdown, on the other hand, often results in diminished
99 levels of cell replication, invasion, and treatment resistance in these cancer types
100 ^{40,44,54,56,57,66–68}. Despite the involvement of SOX2 in the progression of multiple cancer
101 types, little is known about the mechanisms that cause SOX2 overexpression in cancer.

102 Here, we identify a novel enhancer cluster misactivated in breast and lung cancer.
103 This cluster contains two regions, located 124 and 134 kb downstream of the SOX2
104 promoter and referred to as SRR124–134, that drive transcription in BRCA, LUAD, and
105 LUSC. Deletion of this cluster results in significant SOX2 downregulation, leading to
106 genome-wide changes in chromatin accessibility and a globally disrupted transcriptome.
107 The SRR124–134 cluster is highly accessible in most breast and lung patient tumors,
108 where chromatin accessibility at these regions is correlated with SOX2 overexpression
109 and is regulated positively by FOXA1 and negatively by NFIB. Finally, we found that
110 both SRR124 and SRR134 are highly conserved in the mouse and are essential for
111 postnatal survival, as homozygous deletion of their homologous regions results in lethal
112 EA/TEF. These findings serve as a prime example of how cancer cells activate

- 113 enhancers that were decommissioned during development to drive the expression of
- 114 developmentally associated transcription factors during tumorigenesis.

115 RESULTS

116 Two regions downstream of *SOX2* gain enhancer features in cancer cells

117 *SOX2* overexpression occurs in multiple types of cancer ^{reviewed in 39}. To examine
118 which cancer types have the highest levels of *SOX2* upregulation, we performed
119 differential expression analysis by calculating the \log_2 fold change (\log_2 FC) of *SOX2*
120 transcription from 21 TCGA primary solid tumors (see Supplementary Table S1 for
121 cancer type abbreviations) compared to normal tissue samples ⁶⁹. We found that BRCA
122 (\log_2 FC = 3.31), COAD (\log_2 FC = 1.38), GBM (\log_2 FC = 2.05), LIHC (\log_2 FC = 3.22),
123 LUAD (\log_2 FC = 1.36), and LUSC (\log_2 FC = 4.91) tumors had the greatest *SOX2*
124 upregulation (\log_2 FC > 1; FDR-adjusted Q < 0.01; Figure 1A, Supplementary Table
125 S2). As a negative control, we ran this same analysis using the housekeeping gene
126 *PUM1* ⁷⁰ and found no cancer types with significant upregulation of this gene
127 (Supplementary Figure S1A, Supplementary Table S3).

128 Next, we divided BRCA, COAD, GBM, LIHC, LUAD, and LUSC patients (n =
129 3,064) into four groups according to their *SOX2* expression. Gene expression levels
130 were measured by RNA-seq counts normalized by library size and transformed to a \log_2
131 scale, hereinafter referred to as \log_2 counts. Cancer patients within the top group (25%
132 highest *SOX2* expression; \log_2 counts > 10.06) have a significantly ($P = 1.27 \times 10^{-23}$, log-
133 rank test) lower overall probability of survival compared to cancer patients within the
134 bottom group (25% lowest *SOX2* expression; \log_2 counts < 1.68) (Supplementary
135 Figure S1B, Supplementary Table S4). We also examined the relationship between
136 *SOX2* copy number and *SOX2* overexpression within these six tumor types. Although
137 previous studies have shown that *SOX2* is frequently amplified in squamous cell

138 carcinoma ^{57,58,71,72}, we found that most BRCA (88%), COAD (98%), GBM (91%), LIHC
139 (94%), and LUAD (92%) tumors were diploid for SOX2. In addition, BRCA ($P = 0.011$,
140 Holm-adjusted Dunn's test), GBM ($P = 1.18 \times 10^{-3}$), LIHC ($P = 0.016$), LUAD ($P = 0.012$),
141 and LUSC ($P = 2.72 \times 10^{-11}$) diploid tumors significantly overexpressed SOX2 compared
142 to normal tissue (Figure 1B, Supplementary Table S5). This indicates that gene
143 amplification is dispensable for driving SOX2 overexpression in most cancer types.

144 We investigated whether the SOX2 locus gains epigenetic features associated
145 with active enhancers in cancer. Enhancer features commonly include accessible
146 chromatin determined by either Assay for Transposase Accessible Chromatin with high-
147 throughput sequencing (ATAC-seq) ⁷³ or DNase I hypersensitive sites sequencing
148 (DNase-seq) ⁷⁴, and histone modifications including histone H3 lysine 4
149 monomethylation (H3K4me1) and histone H3 lysine 27 acetylation (H3K27ac) ^{75,76}. To
150 study gains in enhancer features within the SOX2 locus, we initially focused our
151 analyses on luminal A breast cancer, the most common subtype of BRCA to
152 significantly ($P = 0.021$, Tukey's test) overexpress SOX2 (Supplementary Figure S1C)
153 ⁶⁹. MCF-7 cells are a widely used ER⁺/PR⁺/HER2⁻ luminal A breast adenocarcinoma
154 model ⁷⁷, which have been previously described to overexpress SOX2 ^{40,68,78,79}. After
155 confirming that SOX2 is one of the most upregulated genes in MCF-7 cells (\log_2 FC =
156 10.75; FDR-adjusted Q = 2.20×10^{-36} ; Supplementary Figure S1D, Supplementary Table
157 S6) compared to healthy breast epithelium ⁸⁰, we contrasted their chromatin
158 accessibility and histone modifications ⁸¹. By intersecting 1,500 bp regions that contain
159 at least 500 bp overlap between H3K27ac and ATAC-seq peaks, we found that 19
160 putative enhancers gained (\log_2 FC > 1) both these features within ± 1 Mb from the

161 SOX2 transcription start site (TSS) in MCF-7 cells (Figure 1C, Supplementary Table
162 S7). Besides the SOX2 promoter (pSOX2), we identified a downstream cluster
163 containing two regions that have gained the highest ATAC-seq and H3K27ac signal in
164 MCF-7 cells: SRR124 (124 kb downstream of pSOX2) and SRR134 (134 kb
165 downstream of pSOX2). The previously described SRR1, SRR2¹⁷⁻¹⁹, and the human
166 ortholog of the mouse SCR (hSCR)^{20,21}, however, lacked substantial gains in enhancer
167 features in MCF-7 cells.

168 Alongside gains in chromatin features, another characteristic of active enhancers
169 is the binding of numerous (> 10) transcription factors⁸²⁻⁸⁴. Chromatin
170 Immunoprecipitation Sequencing (ChIP-seq) data from ENCODE⁸¹ on 117 transcription
171 factors revealed 48 different factors present at the SRR124–134 cluster in MCF-7 cells,
172 with the majority (47) of these factors present at SRR134 (Figure 1D). Transcription
173 factors bound at both SRR124 and SRR134 include CEBPB, CREB1, FOXA1, FOXM1,
174 NFIB, NR2F2, TCF12, and ZNF217. An additional feature of distal enhancers is that
175 they contact their target genes through long-range chromatin interactions^{85,86}.
176 Chromatin Interaction Analysis by Paired-End-Tag sequencing (ChIA-PET)⁸⁷ showed
177 two interesting RNA Polymerase II (RNAPII)-mediated chromatin interactions in MCF-7
178 cells: one between the SOX2 gene and SRR134, and one between SRR124 and
179 SRR134 (Figure 1E). Beyond MCF-7 cells, we found that H520 (LUSC), PC-9 (LUAD),
180 and T47D (luminal A BRCA) cancer cell lines, which display varying levels of SOX2
181 expression (Supplementary Figure S1E), also gained substantial enhancer features at
182 SRR124 and SRR134 when compared to healthy tissue (Figure 1E)⁸⁸⁻⁹⁰. Together,

183 these data suggest SRR124 and SRR134 could be active enhancers driving SOX2
184 transcription in BRCA, LUAD, and LUSC.

185

186 **Figure 1: A cluster 124–134 kilobases downstream of SOX2 gains enhancer**

187 **features in cancer cells. (A)** Super-logarithmic RNA-seq volcano plot of SOX2

188 expression from 21 cancer types compared to normal tissue ⁶⁹. Cancer types with log₂

189 FC > 1 and FDR-adjusted Q < 0.01 were considered to significantly overexpress SOX2.

190 Error bars: standard error. **(B)** SOX2 log₂-normalized expression (log₂ counts)

191 associated with the SOX2 copy number from BRCA (n = 1174), COAD (n = 483), GBM

192 (n = 155), LIHC (n = 414), LUAD (n = 552), and LUSC (n = 546) patient tumors ⁶⁹. RNA-

193 seq reads were normalized to library size using DESeq2 ⁹¹. Error bars: standard

194 deviation. Significance analysis by Dunn's test ⁹² with Holm correction ⁹³. **(C)** 1,500 bp

195 genomic regions within ± 1 Mb from the SOX2 transcription start site (TSS) that gained

196 enhancer features in MCF-7 cells ⁸¹ compared to healthy breast epithelium ⁸⁰. Regions

197 that gained both ATAC-seq and H3K27ac ChIP-seq signal above our threshold (log₂ FC

198 > 1, dashed line) are highlighted in pink. Each region was labelled according to their

199 distance in kilobases (kb) to the SOX2 promoter (pSOX2, bolded). **(D)** ChIP-seq signal

200 for H3K4me1 and H3K27ac, ATAC-seq signal, and transcription factor ChIP-seq peaks

201 at the SRR124–134 cluster in MCF-7 cells. Datasets from ENCODE ⁸¹. **(E)** UCSC

202 Genome Browser ⁹⁴ display of H3K4me1 and H3K27ac ChIP-seq signal, DNase-seq

203 and ATAC-seq chromatin accessibility signal, and ChIA-PET RNA Polymerase II

204 (RNAPII) interactions around the SOX2 gene within breast (normal tissue and 2 BRCA

205 cancer cell lines) and lung (normal tissue, one LUAD, and one LUSC cancer cell lines)

206 samples^{81,88–90}. Relevant RNAPII interactions (between SRR124 and SRR134, and
207 between SRR134 and pSOX2) are highlighted in maroon.

208

209 **The SRR124–134 cluster is essential for SOX2 expression in BRCA and LUAD** 210 **cells**

211 To assess SRR124 and SRR134 enhancer activity alongside the embryonic-
212 associated SRR1, SRR2, and hSCR regions, we used a reporter vector containing the
213 firefly luciferase gene under the control of a minimal promoter (minP, pGL4.23). We
214 transfected each enhancer construct into the BRCA (MCF-7, T47D), LUAD (PC-9), and
215 LUSC (H520) cell lines and measured luciferase activity as a relative fold change (FC)
216 compared to the empty minP vector. SRR134 demonstrated the strongest enhancer
217 activity, with the MCF-7 (FC = 6.42; $P < 2 \times 10^{-16}$, Dunnett's test), T47D (FC = 3.36; $P =$
218 9.34×10^{-10}), H520 (FC = 2.37; $P = 1.22 \times 10^{-6}$), and PC-9 (FC = 2.03; $P = 9.79 \times 10^{-5}$) cell
219 lines displaying a significant increase in luciferase activity compared to minP (Figure
220 2A). SRR124 also showed a modest, significant increase in luciferase activity compared
221 to minP in the MCF-7 (FC = 1.53; $P = 4.27 \times 10^{-2}$), T47D (FC = 1.80; $P = 4.57 \times 10^{-2}$), and
222 PC-9 (FC = 1.60; $P = 4.27 \times 10^{-2}$) cell lines. The embryonic-associated enhancers SRR1,
223 SRR2, and hSCR, however, showed no significant enhancer activity ($P > 0.05$) in all
224 four cell lines.

225 Although reporter assays can be used to assess enhancer activity, enhancer
226 knockout approaches remain the current gold standard method for enhancer validation
227^{95,96}. To investigate whether the SRR124–134 cluster drives SOX2 expression in cancer
228 cells, we used CRISPR/Cas9 to delete this cluster from the H520, MCF-7, PC-9, and

229 T47D cell lines. RT-qPCR showed that homozygous SRR124–134 deletion ($\Delta\text{ENH}^{-/-}$)
230 causes a profound ($> 99.5\%$) and significant ($P < 0.001$, Dunnett's test) loss of *SOX2*
231 expression in both the MCF-7 and PC-9 cell lines (Figure 2B). Immunoblot analysis
232 confirmed the depletion of the *SOX2* protein in $\Delta\text{ENH}^{-/-}$ MCF-7 cells (Supplementary
233 Figure S2A). Heterozygous SRR124–134 deletion ($\Delta\text{ENH}^{+/-}$) also significantly ($P <$
234 0.001) reduced *SOX2* expression by $\sim 60\%$ in both MCF-7 and PC-9 cells (Figure 2B).
235 Although we were unable to isolate a homozygous deletion clone from T47D cells,
236 multiple independent heterozygous $\Delta\text{ENH}^{+/-}$ T47D clonal isolates showed a significant
237 downregulation ($>50\%$; $P < 0.001$) in *SOX2* expression (Figure 2C). Interestingly, we
238 did not find a significant ($P > 0.05$) impact on *SOX2* expression in $\Delta\text{ENH}^{+/-}$ or $\Delta\text{ENH}^{-/-}$
239 H520 cells (Supplementary Figure S2B), which indicates that *SOX2* transcription is
240 sustained by a different mechanism in these cells. To assess the impact of the loss of
241 *SOX2* expression in the tumor initiation capacity of enhancer-deleted cells, we
242 performed a colony formation assay with MCF-7 and PC-9 $\Delta\text{ENH}^{-/-}$ cells. We found that
243 both MCF-7 ($P = 3.53 \times 10^{-4}$, t-test) and PC-9 ($P = 1.26 \times 10^{-5}$) $\Delta\text{ENH}^{-/-}$ cells showed a
244 significant decrease ($> 50\%$) in their ability to form colonies compared to WT cells
245 (Figure 2D), further suggesting that SRR124–134-driven *SOX2* overexpression is
246 required to sustain high tumor initiation capacity in BRCA and LUAD.

247 Next, we performed total RNA sequencing (RNA-seq) to measure changes in the
248 transcriptome of $\Delta\text{ENH}^{-/-}$ MCF-7 cells compared to WT MCF-7 cells. As expected, all
249 three replicates of each genotype clustered together (Supplementary Figure S2C). In
250 addition to *SOX2* downregulation (Figure 2E), differential expression analysis showed a
251 total of 529 genes differentially ($|\log_2 \text{FC}| > 1$; FDR-adjusted $Q < 0.01$) expressed in

252 Δ ENH^{-/-} MCF-7 cells (Figure 2F, Supplementary Table S8). From these, 312 genes
253 significantly lost expression (59%), whereas 217 (41%) genes significantly gained
254 expression in Δ ENH^{-/-} MCF-7 cells compared to WT MCF-7 cells (Supplementary
255 Figure S2D). *SOX2* was the gene with the highest loss in expression (\log_2 FC = -10.24;
256 $Q = 1.23 \times 10^{-43}$) in Δ ENH^{-/-} MCF-7 cells, followed by *CT83* (\log_2 FC = -8.43; $Q =$
257 1.07×10^{-8}), and *GUCY1A1* (\log_2 FC = -6.96; $Q = 5.09 \times 10^{-15}$). On the other hand, genes
258 with the most significant gain in expression within Δ ENH^{-/-} MCF-7 cells included the
259 protocadherins *PCDH7* (\log_2 FC = 5.34; $Q < 1 \times 10^{-200}$), *PCDH10* (\log_2 FC = 5.29; $Q <$
260 1×10^{-200}), and *PCDH11X* (\log_2 FC = 4.73; $Q = 9.29 \times 10^{-110}$). In addition, deletion of the
261 SRR124–134 cluster reduced *SOX2* expression back to the levels found in healthy
262 breast epithelium ($P = 0.48$, Tukey's test)^{80,81} (Figure 2G). Together, these data confirm
263 that the SRR124–134 cluster drives *SOX2* overexpression in BRCA and LUAD.

264

265 **Figure 2: The SRR124–134 cluster drives *SOX2* overexpression in MCF-7, T47D,**
266 **and PC-9 cells. (A)** Enhancer reporter assay comparing luciferase activity driven by the
267 SRR1, SRR2, SRR124, SRR134, and hSCR regions to an empty vector containing only
268 a minimal promoter (minP). Enhancer constructs were assayed in the BRCA (MCF-7,
269 T47D), LUAD (PC-9), and LUSC (H520) cell lines. Dashed line: average activity of
270 minP. Error bars: standard deviation. Significance analysis by Dunnett's test ($n = 5$; * P
271 < 0.05 , *** $P < 0.001$, ns: not significant)⁹⁷. **(B)** RT-qPCR analysis of *SOX2* transcript
272 levels in SRR124–134 heterozygous- (Δ ENH^{+/-}) and homozygous- (Δ ENH^{-/-}) deleted
273 MCF-7 (BRCA) and PC-9 (LUAD) clones compared to WT cells. Error bars: standard
274 deviation. Significance analysis by Dunnett's test ($n = 3$; *** $P < 0.001$). **(C)** RT-qPCR

275 analysis of *SOX2* transcript levels in three independent SRR124–134 heterozygous-
276 deleted ($\Delta\text{ENH}^{+/-}$) T47D clonal isolates compared to WT cells. Error bars: standard
277 deviation. Significance analysis by Dunnett's test ($n = 4$; *** $P < 0.001$). **(D)** Crystal violet
278 absorbance (570 nm) from a colony formation assay with WT and $\Delta\text{ENH}^{-/-}$ MCF-7 and
279 PC-9 cells. Total absorbance was normalized to the average absorbance from WT cells
280 within each cell line. Significance analysis by t-test with Holm correction ($n = 5$; *** $P <$
281 0.001). **(E)** UCSC Genome Browser ⁹⁴ view of the SRR124–134 cluster deletion in
282 $\Delta\text{ENH}^{-/-}$ MCF-7 cells with RNA-seq tracks from normal breast epithelium ⁸⁰, WT and
283 $\Delta\text{ENH}^{-/-}$ MCF-7 cells. Arrow: reduction in RNA-seq signal at the *SOX2* gene in $\Delta\text{ENH}^{-/-}$
284 MCF-7 cells. **(F)** Volcano plot with DESeq2 ⁹¹ differential expression analysis between
285 $\Delta\text{ENH}^{-/-}$ and WT MCF-7 cells. Blue: 312 genes that significantly lost expression (\log_2
286 $\text{FC} < -1$; FDR-adjusted $Q < 0.01$) in $\Delta\text{ENH}^{-/-}$ MCF-7 cells. Pink: 217 genes that
287 significantly gained expression ($\log_2 \text{FC} > 1$; $Q < 0.01$) in $\Delta\text{ENH}^{-/-}$ MCF-7 cells. Grey:
288 35,891 genes that maintained similar ($-1 \leq \log_2 \text{FC} \leq 1$) expression between $\Delta\text{ENH}^{-/-}$
289 and WT MCF-7 cells. **(G)** Comparison of *SOX2* transcript levels between WT MCF-7
290 and either $\Delta\text{ENH}^{-/-}$ MCF-7 or healthy breast epithelium cells ⁸⁰, and between $\Delta\text{ENH}^{-/-}$
291 MCF-7 and healthy breast epithelium cells. RNA-seq reads were normalized to library
292 size using DESeq2 ⁹¹. Error bars: standard deviation. Significance analysis by Tukey's
293 test (*** $P < 0.001$, ns: not significant) ⁹⁸.

294

295 **SOX2 regulates pathways associated with epithelium development in luminal A**
296 **BRCA**

297 Because SOX2 regulates cell proliferation and differentiation pathways in other
298 epithelial cells^{38,99}, we decided to further investigate the molecular function of SOX2 in
299 luminal A BRCA cells by utilizing our SOX2-depleted $\Delta\text{ENH}^{-/-}$ MCF-7 cell model. Gene
300 Set Enrichment Analysis (GSEA) showed a significant (FDR-adjusted $Q < 0.05$)
301 depletion of multiple epithelium-associated processes within the transcriptome of
302 $\Delta\text{ENH}^{-/-}$ MCF-7 cells, as indicated by normalized enrichment score [NES] < 1
303 (Supplementary Table S9). These processes included epidermis development (NES = -
304 1.93; $Q = 0.001$; Figure 3A), epithelial cell differentiation (NES = -1.67; $Q = 0.007$;
305 Figure 3B), and cornification (NES = -2.11; $Q = 0.006$; Figure 3C). This suggests that
306 SOX2 regulates epithelial development and differentiation in luminal A BRCA cells.

307 SOX2 is a pioneer transcription factor that associates with its motif in
308 heterochromatin¹⁰⁰ and recruits chromatin-modifying complexes¹⁰¹ in embryonic and
309 reprogrammed stem cells. We performed ATAC-seq in $\Delta\text{ENH}^{-/-}$ MCF-7 cells and
310 compared chromatin accessibility to WT MCF-7 cells to identify genome-wide loci that
311 are dependent on SOX2 to remain accessible in luminal A BRCA. As expected, the
312 ATAC-seq signal from all replicates was highly enriched around gene TSS
313 (Supplementary Figure S3A), with both WT and $\Delta\text{ENH}^{-/-}$ samples having higher
314 chromatin accessibility at the TSS of highly expressed genes (Supplementary Figure
315 S3B and Supplementary Figure S3C). Correlation analysis also confirmed the clustering
316 of all three replicates from each genotype (Supplementary Figure S3D). Together with
317 the SRR124–134 cluster and pSOX2 (Figure 3D), a total of 3,076 500-bp regions had
318 significant ($|\log_2 \text{FC}| > 1$; FDR-adjusted $Q < 0.01$) changes in chromatin accessibility in
319 $\Delta\text{ENH}^{-/-}$ compared to WT MCF-7 cells (Figure 3E, Supplementary Table S10). Most

320 regions (86%, 2,636 regions) significantly lost chromatin accessibility in $\Delta\text{ENH}^{-/-}$ MCF-7
321 cells and 76% (2,024 regions) of these regions also gained chromatin accessibility in
322 WT MCF-7 compared to healthy breast epithelium⁸⁰ (Supplementary Table S11).
323 Together, this indicates that SOX2 has an important role in regulating the chromatin
324 accessibility changes acquired in luminal A BRCA.

325 We used TOBIAS¹⁰² to analyze changes in transcription factor footprints within
326 ATAC-seq peaks in $\Delta\text{ENH}^{-/-}$ compared to WT MCF-7 cells. From 841 vertebrate motifs
327¹⁰³, we found a total of 281 motifs with a significant ($|\log_2 \text{FC}| > 0.1$; FDR-adjusted $Q <$
328 0.01) differential binding score (Figure 3F, Supplementary Table S12). Most of these
329 motifs (97%, 272 motifs) were underrepresented within ATAC-seq peaks in $\Delta\text{ENH}^{-/-}$
330 compared to WT MCF-7 cells, indicating that reduced SOX2 expression affects the
331 binding of multiple other transcription factors. Among them, the GRHL1 ($\log_2 \text{FC} = -$
332 0.519; $Q = 3 \times 10^{-179}$), TFCP2 ($\log_2 \text{FC} = -0.462$; $Q = 1.03 \times 10^{-172}$), RUNX2 ($\log_2 \text{FC} = -$
333 0.352; $Q = 8.02 \times 10^{-164}$), GRHL2 ($\log_2 \text{FC} = -0.343$; $Q = 4.43 \times 10^{-174}$), TEAD3 ($\log_2 \text{FC} =$
334 -0.235 ; $Q = 9.74 \times 10^{-155}$), and SOX4 ($\log_2 \text{FC} = -0.232$; $Q = 5.33 \times 10^{-167}$) motifs (Figure
335 3G) had the most reduced binding score in $\Delta\text{ENH}^{-/-}$ MCF-7 cells compared to WT MCF-
336 7 cells. These factors belong to three main motif clusters: GRHL/TFCP (cluster 33;
337 aaAACAGGTTtcAgtt), RUNX (cluster 60; ttctTGtGGTTttt), TEAD (cluster 2;
338 tccAcATTCCAggcCTTta), and SOX (cluster 8; acggaACAATGgaagTGTT)¹⁰³. The SOX
339 cluster also included the SOX2 ($\log_2 \text{FC} = -0.175$; $Q = 6.61 \times 10^{-139}$) motif.

340 Next, we aimed to analyze ChIP-seq data from transcription factors within these
341 motif clusters in MCF-7 cells. We utilized two published datasets: GRHL2⁸⁹ and RUNX2
342¹⁰⁴. Regions that lost ($\log_2 \text{FC} < -1$; $Q < 0.01$) chromatin accessibility in $\Delta\text{ENH}^{-/-}$ MCF-7

343 cells significantly overlapped with regions with binding of either of these transcription
344 factors ($P < 2 \times 10^{-16}$, hypergeometric test). Among the 2,636 regions that lost chromatin
345 accessibility, 40% (750 regions) also show GRHL2 binding (Supplementary Figure
346 S3E), whereas 21% (552 regions) share RUNX2 binding (Supplementary Figure S3F).
347 We found multiple SOX motifs significantly (FDR-adjusted $Q < 0.001$) enriched within
348 peaks from both GRHL2 (Supplementary Table S13) and RUNX2 (Supplementary
349 Table S14) ChIP-seq data, further suggesting that SOX2 collaborates with GRHL2 and
350 RUNX2 to maintain chromatin accessibility in luminal A BRCA. Expression levels of
351 either *GRHL2* or *RUNX2*, however, were not significantly affected by *SOX2*
352 downregulation in $\Delta ENH^{-/-}$ MCF-7 cells ($-1 \leq \log_2 FC \leq 1$; Supplementary Table S8),
353 indicating that they are not directly regulated by *SOX2* at the transcriptional level but
354 may interact at the protein level.

355

356 **Figure 3: SOX2 downregulation impacts chromatin accessibility in luminal A**
357 **BRCA. (A – C)** Gene Set Enrichment Analysis (GSEA) in the transcriptome of $\Delta ENH^{-/-}$
358 compared to WT MCF-7 cells. Genes were ranked according to their change in
359 expression ($\log_2 FC$). A subset of GO terms significantly enriched among
360 downregulated genes in $\Delta ENH^{-/-}$ MCF-7 cells are displayed, indicated by the
361 normalized enrichment score (NES) < 1 : **(A)** epidermis development, **(B)** epithelial cell
362 differentiation, and **(C)** cornification. GSEA was performed using clusterProfiler¹⁰⁵ with
363 an FDR-adjusted $Q < 0.05$ threshold. Green line: running enrichment score. **(D)** UCSC
364 Genome Browser⁹⁴ view of the SRR124–134 deletion in $\Delta ENH^{-/-}$ MCF-7 cells with
365 ATAC-seq tracks from breast epithelium⁸⁰, WT, and $\Delta ENH^{-/-}$ MCF-7 cells. **(E)** Volcano

366 plot with differential ATAC-seq analysis of $\Delta\text{ENH}^{-/-}$ MCF-7 cells compared to WT. Blue:
367 2,638 regions that lost ($\log_2 \text{FC} < -1$; FDR-adjusted $Q < 0.01$) chromatin accessibility in
368 $\Delta\text{ENH}^{-/-}$ MCF-7 cells. Pink: 440 regions that gained ($\log_2 \text{FC} > 1$; $Q < 0.01$) chromatin
369 accessibility in $\Delta\text{ENH}^{-/-}$ MCF-7 cells. Grey: 132,726 regions that retained chromatin
370 accessibility in $\Delta\text{ENH}^{-/-}$ MCF-7 cells ($-1 \leq \log_2 \text{FC} \leq 1$). Regions were labelled with their
371 closest gene within a ± 1 Mb distance threshold. Differential chromatin accessibility
372 analysis was performed using diffBind ¹⁰⁶. **(F)** Volcano plot with ATAC-seq footprint
373 analysis of differential transcription factor binding in $\Delta\text{ENH}^{-/-}$ MCF-7 cells compared to
374 WT. Blue: 272 underrepresented ($\log_2 \text{FC} < -0.1$; FDR-adjusted $Q < 0.01$) motifs in
375 ATAC-seq peaks from $\Delta\text{ENH}^{-/-}$ MCF-7 cells. Pink: 9 overrepresented ($\log_2 \text{FC} > 0.1$; Q
376 < 0.01) motifs in ATAC-seq peaks from $\Delta\text{ENH}^{-/-}$ MCF-7 cells. Grey: 560 motifs with no
377 representative change ($-0.1 \leq \log_2 \text{FC} \leq 0.1$) within ATAC-seq peaks from $\Delta\text{ENH}^{-/-}$
378 MCF-7 cells. **(G)** Sequence motifs of the top 6 transcription factors with the lowest
379 binding score in $\Delta\text{ENH}^{-/-}$ compared to WT MCF-7 cells: GRHL1, TFCP2, RUNX2,
380 GRHL2, TEAD3, SOX4. Footprint analysis was performed using TOBIAS ¹⁰² utilizing the
381 JASPAR 2022 motif database ¹⁰³.

382

383 **The SRR124–134 cluster is associated with SOX2 overexpression in primary**
384 **tumors**

385 With the confirmation that the SRR124–134 cluster drives SOX2 overexpression in
386 the BRCA and LUAD cell lines, we investigated chromatin accessibility at this enhancer
387 cluster within primary tumors isolated from cancer patients. By analyzing the pan-cancer
388 ATAC-seq dataset from TCGA ¹⁰⁷, we found that SRR124 and SRR134 are most

389 accessible within LUSC, LUAD, BRCA, bladder carcinoma (BLCA), stomach
390 adenocarcinoma (STAD), and uterine endometrial carcinoma (UCEC) patient tumors
391 (Figure 4A). We also quantified the ATAC-seq signal at six other regions (genomic
392 coordinates in Supplementary Table S15): the SOX2 embryonic-associated enhancers
393 (SRR1, SRR2, hSCR), pSOX2, a gene regulatory desert with no enhancer features
394 located between the SOX2 gene and the SRR124–134 cluster (desert), and the
395 promoter of the housekeeping gene *RAB7A* (pRAB7A, positive control). We then
396 compared the chromatin accessibility levels at each of these regions to the promoter of
397 the repressed olfactory gene *OR5K1* (pOR5K1, negative control). Both SRR124 and
398 SRR134 showed significantly increased ($P < 0.05$, Holm-adjusted Dunn's test)
399 chromatin accessibility when compared to pOR5K1 in BLCA (SRR124 $P = 0.014$;
400 SRR134 $P = 1.52 \times 10^{-3}$; Holm-adjusted Dunn's test), BRCA (SRR124 $P = 1.70 \times 10^{-20}$;
401 SRR134 $P = 1.03 \times 10^{-16}$), LUAD (SRR124 $P = 6.76 \times 10^{-7}$; SRR134 $P = 3.26 \times 10^{-6}$), LUSC
402 (SRR124 $P = 1.62 \times 10^{-6}$; SRR134 $P = 7.08 \times 10^{-4}$), STAD (SRR124 $P = 1.15 \times 10^{-4}$;
403 SRR134 $P = 1.96 \times 10^{-7}$), and UCEC (SRR124 $P = 3.15 \times 10^{-5}$; SRR134 $P = 0.025$)
404 patient tumors (Figure 4B).

405 One explanation for increased chromatin accessibility is locus amplification. While
406 LUSC had high levels of chromatin accessibility likely related to previously described
407 SOX2 amplifications^{57,58,71,72}, most patient tumors showed no evidence of locus
408 amplifications extending to the SRR124–134 cluster, as evidenced by the lack of
409 significant ($P > 0.05$) accessibility at the intermediate desert region. In contrast, the
410 SRR124–134 cluster displayed a consistent pattern of accessible chromatin across
411 multiple cancer types: BLCA, BRCA, LUAD, LUSC, STAD, and UCEC (Figure 4C).

412 GBM and LGG tumors lacked accessible chromatin at this cluster but displayed
413 increased chromatin accessibility at the SRR1 and SRR2 enhancers (Supplementary
414 Figure S4A, Supplementary Table S16), which is consistent with the evidence that
415 SRR1 and SRR2 drive *SOX2* expression in the neural lineage^{17,19,108}.

416 Next, we reasoned that an accessible SRR124–134 cluster drives subsequent
417 *SOX2* transcription within patient tumors. If this is the case, we expect to find positive
418 and significantly correlated chromatin accessibility between this enhancer cluster and
419 pSOX2. Indeed, we found that the majority of BRCA (58%), LUAD (82%), and LUSC
420 (69%) tumors have concurrent accessibility (\log_2 RPM > 0) at pSOX2, SRR124 and
421 SRR134. Patient tumors also showed a significant correlation (Pearson, R) between
422 accessible chromatin signal at pSOX2 and at both SRR124 and SRR134 in BRCA and
423 LUAD (Figure 4D). LUSC tumors showed a significant correlation between accessible
424 chromatin at pSOX2 and SRR124, but not at SRR134 (Figure 4D). As a negative
425 control, we measured the correlation between chromatin accessibility at pSOX2 and at
426 the *SOX2* desert region and found no significant ($P > 0.05$) correlation in any of these
427 cancer types (Supplementary Figure S4B). We also conducted a similar analysis after
428 segregating BRCA tumors into luminal A, luminal B, HER2⁺, and basal-like subtypes.
429 Interestingly, we found that both luminal A and luminal B tumors possess a significant
430 ($P < 0.05$) correlation between enhancer accessibility and pSOX2 accessibility, whereas
431 for HER2⁺ tumors the correlation was weaker (Supplementary Figure S4C). Basal-like
432 tumors, on the other hand, display no accessible chromatin at either SRR124 or
433 SRR134. In summary, a luminal-like BRCA phenotype correlates with increased
434 accessibility at the SRR124–134 cluster.

435 Finally, by separating BRCA, LUAD, and LUSC patient tumors according to their
436 chromatin accessibility at SRR124 and SRR134, we found that tumors with the most
437 accessible chromatin at each of these regions also significantly ($P < 0.05$, t-test)
438 overexpress *SOX2* compared to tumors with low chromatin accessibility at these
439 regions (Figure 4E, Supplementary Table S17). Together, these data are consistent
440 with a model in which increased chromatin accessibility at the SRR124–134 cluster
441 drives *SOX2* overexpression in BRCA, LUAD, and LUSC patient tumors.

442

443 **Figure 4: The SRR124–134 cluster is associated with *SOX2* overexpression in**
444 **cancer patient tumors. (A)** ATAC-seq signal (\log_2 RPM) at SRR124 and SRR134 for
445 294 patient tumors from 14 cancer types¹⁰⁷. Cancer types are sorted in descending
446 order by the median signal between all three regions. Dashed line: regions with a sum
447 of reads above our threshold (\log_2 RPM > 0) were considered “accessible”. Error bars:
448 standard deviation. Underscore: top 6 cancer types with the highest ATAC-seq median
449 signal. **(B)** ATAC-seq signal (\log_2 RPM) at the *RAB7A* promoter (pRAB7A), *SOX2*
450 promoter (pSOX2), SRR1, SRR2, SRR124, SRR134, hSCR, and a desert region within
451 the *SOX2* locus (desert) compared to the background signal at the repressed *OR5K1*
452 promoter (pOR5K1) in BLCA (n = 10), BRCA (n = 74), LUAD (n = 22), LUSC (n = 16),
453 STAD (n = 21), and UCEC (n = 13) patient tumors. Dashed line: regions with a sum of
454 reads above our threshold (\log_2 RPM > 0) were considered “accessible”. Error bars:
455 standard deviation. Significance analysis by Dunn’s test with Holm correction (* $P <$
456 0.05, ** $P < 0.01$, *** $P < 0.001$, ns: not significant). **(C)** UCSC Genome Browser⁹⁴
457 visualization of the *SOX2* region with ATAC-seq data from BLCA, BRCA, LUAD, LUSC,

458 STAD, and UCEC patient tumors (n = 5 in each cancer type)¹⁰⁷. ATAC-seq reads were
459 normalized by library size (RPM). Scale: 0 – 250 RPM. **(D)** ATAC-seq signal at SRR124
460 and SRR134 regions against ATAC-seq signal for the SOX2 promoter (pSOX2) from 74
461 BRCA, 22 LUAD, and 16 LUSC patient tumors. Correlation is shown for accessible
462 chromatin (\log_2 RPM > 0). Grey: tumors with closed chromatin (\log_2 RPM < 0) at either
463 region, not included in the correlation analysis. Significance analysis by Pearson
464 correlation. Bolded line: fitted linear regression model. Shaded area: 95% confidence
465 region for the regression fit. **(E)** Comparison of \log_2 -normalized SOX2 transcript levels
466 (\log_2 counts) between BRCA, LUAD, and LUSC patient tumors according to the
467 chromatin accessibility at SRR124 and SRR134 regions. Chromatin accessibility at
468 each region was considered “low” if \log_2 RPM < -1, or “high” if \log_2 RPM > 1. RNA-seq
469 reads were normalized to library size using DESeq2⁹¹. Error bars: standard deviation.
470 Significance analysis by a two-sided t-test with Holm correction.

471

472 **FOXA1 and NFIB are upstream regulators of the SRR124–134 cluster**

473 With the indication that the SRR124–134 cluster is driving SOX2 overexpression in
474 patient tumors, we investigated which transcription factors regulate this cluster in BRCA,
475 LUAD, and LUSC. From a list of 1622 human transcription factors¹⁰⁹, we found that the
476 expression of 115 transcription factors was significantly (FDR-adjusted Q < 0.05)
477 associated with chromatin accessibility levels at SRR124, whereas accessibility at
478 SRR134 was associated with the expression of 90 transcription factors (Figure 5A,
479 Supplementary Table S18). From this list, we focused our investigation on FOXA1 and
480 NFIB, which show binding at both SRR124 and SRR134 in MCF-7 cells⁸¹.

481 The expression of *FOXA1* is positively (Pearson correlation $R > 0$) and significantly
482 correlated to accessible chromatin at both SRR124 ($R = 0.39$; FDR-adjusted $Q =$
483 1.97×10^{-3}) and SRR134 ($R = 0.46$; $Q = 1.41 \times 10^{-4}$) (Figure 5B). By separating BRCA,
484 LUAD, and LUSC patient tumors according to the chromatin accessibility levels at each
485 region, we found that tumors with the most accessible chromatin within SRR124 ($P =$
486 2.38×10^{-4} , t-test) and SRR134 ($P = 1.53 \times 10^{-4}$) also significantly overexpress *FOXA1*
487 compared to tumors with low accessibility at these regions (Figure 5C, Supplementary
488 Table S19). On the other hand, we found the expression of *NFIB* to be negatively
489 (Pearson correlation $R < 0$) and significantly correlated with chromatin accessibility at
490 both SRR124 ($R = -0.49$; $Q = 4.12 \times 10^{-5}$) and SRR134 ($R = -0.51$; $Q = 1.32 \times 10^{-5}$)
491 (Figure 5D). Patient tumors with highly accessible chromatin within SRR124 ($P =$
492 1.46×10^{-6}) and SRR134 ($P = 1.24 \times 10^{-5}$) also display significantly downregulated *NFIB*
493 expression (Figure 5E, Supplementary Table S20). These data suggest that whereas
494 *FOXA1* could be inducing increased accessibility at the SRR124–134 cluster, *NFIB*
495 expression could counteract *FOXA1* by acting as a repressor.

496 To assess the contribution of these transcription factors to enhancer activity, we
497 overexpressed either *FOXA1* or *NFIB* in H520, MCF-7, PC-9, and T47D cells and
498 compared SRR124 and SRR134 activity to cells transfected with an empty vector
499 (mock) containing only a fluorescent marker. Although endogenous *FOXA1* and *NFIB*
500 expression levels are already high in both MCF-7 and T47D cells (Supplementary
501 Figure S5A), we found that overexpression of *FOXA1* significantly increased (\log_2 FC $>$
502 1; $P < 0.05$, Tukey's test) the enhancer activity of both SRR124 and SRR134 in the
503 H520, MCF-7, PC-9, and T47D cell lines, whereas *NFIB* overexpression led to a

504 significant decrease ($\log_2 FC < 1$; $P < 0.05$) in SRR124 and SRR134 enhancer activity
505 in the H520, MCF-7, and T47D cell lines (Figure 5F). This further indicates that *FOXA1*
506 overexpression increases SRR124–134 activity, whereas NFIB represses the activity of
507 this cluster.

508 To assess the importance of FOXA1 and NFIB motifs in modulating enhancer
509 activity, we analyzed the SRR134 sequence using the JASPAR2022 motif database ¹⁰³
510 and mutated FOXA1 (GTAAACA) or NFIB (TGGCAnnnnGCCAA) motifs (mutated
511 SRR134 sequences in Supplementary Table S21). We found that mutation of the
512 FOXA1 motif abolished SRR134 enhancer activity compared to WT SRR134 within
513 MCF-7 ($P = 1.53 \times 10^{-5}$, Tukey's test), PC-9 ($P = 1 \times 10^{-2}$), and T47D ($P = 4.48 \times 10^{-6}$) cells,
514 whereas no significant change ($P > 0.05$) in enhancer activity was found for the NFIB-
515 mutated construct (Figure 5G). These data indicate that the FOXA1 motif is crucial for
516 sustaining SRR134 activity, whereas the NFIB motif is dispensable in this context, as
517 would be expected for a negative regulator under conditions where the activity of the
518 target is high.

519 With the evidence that these transcription factors are modulating SRR124–134
520 activity, we investigated their transcriptional effects over *SOX2* expression. We used
521 CRISPR homology-directed repair (HDR) to create an MCF-7 cell line in which the
522 *SOX2* gene is tagged with a 2A self-cleaving peptide (P2A) followed by a blue
523 fluorescent protein (tagBFP). This cell line, MCF-7 *SOX2*-P2A-tagBFP, allows rapid
524 visualization of *SOX2* transcriptional changes by measuring tagBFP signal through
525 fluorescence-activated cell sorting (FACS). To validate this model, we sorted cells within
526 the top 10% (BFP^{+ve}) and bottom 10% (BFP^{-ve}) tagBFP signal (Supplementary Figure

527 S5B). We found that BFP^{+ve} cells showed a significant ($P = 4.25 \times 10^{-5}$, paired t-test)
528 increase in *SOX2* expression, and significantly upregulated transcription of enhancer
529 RNA (eRNA) at SRR124 ($P = 1.54 \times 10^{-4}$) and SRR134 ($P = 5.13 \times 10^{-5}$) compared to
530 BFP^{-ve} cells (Figure 5H). This confirms that the tagBFP signal is directly correlated to
531 *SOX2* transcription levels in MCF-7 *SOX2*-P2A-tagBFP cells.

532 Finally, we overexpressed *FOXA1* or *NFIB* in MCF-7 *SOX2*-P2A-tagBFP to assess
533 changes in *SOX2* transcription. Although overexpression of *FOXA1* did not significantly
534 (chi-squared $T(x)=63.70$) change tagBFP signal, we found that overexpression of *NFIB*
535 significantly (chi-squared $T(x)=1168.88$) reduced tagBFP signal compared to
536 transfection of an empty vector (mock) (Figure 5I). This confirms the repression effect of
537 *NFIB* over *SOX2* expression and illustrates a potential mechanism upstream of *SOX2*
538 that modulates chromatin accessibility at the SRR124–134 cluster and subsequent
539 control of *SOX2* transcription in cancer cells.

540

541 **Figure 5: FOXA1 and NFIB are upstream regulators of SRR124 and SRR134. (A)**
542 Heatmap of the Pearson correlation between transcription factor expression⁶⁹ and
543 chromatin accessibility¹⁰⁷ at SRR124 and SRR134 in BRCA, LUAD, and LUSC patient
544 tumors ($n = 111$). Transcription factors are ordered according to their correlation to
545 chromatin accessibility at each region. Red: transcription factors with a positive
546 correlation ($R > 0$; FDR-adjusted $Q < 0.05$) to chromatin accessibility. Blue: transcription
547 factors with a negative correlation ($R < 0$; $Q < 0.05$) to chromatin accessibility. Asterisk:
548 transcription factors that show binding at SRR124 or SRR134 by ChIP-seq⁸¹. **(B)**
549 Correlation analysis between *FOXA1* expression (\log_2 counts) and chromatin

550 accessibility (\log_2 RPM) at SRR124 and SRR134 regions in BRCA (n = 74), LUAD (n =
551 21), and LUSC (n = 16) tumors. RNA-seq reads were normalized to library size using
552 DESeq2⁹¹. Significance analysis by Pearson correlation (n = 111). Bolded line: fitted
553 linear regression model. Shaded area: 95% confidence region for the regression fit. **(C)**
554 Comparison of *FOXA1* expression (\log_2 counts) from BRCA, LUAD, and LUSC patient
555 tumors according to their chromatin accessibility at the SRR124 and SRR134 regions.
556 Chromatin accessibility at each region was considered “low” if \log_2 RPM < 1, or “high” if
557 \log_2 RPM > 1. RNA-seq reads were normalized to library size using DESeq2⁹¹. Error
558 bars: standard deviation. Significance analysis by a two-sided t-test with Holm
559 correction. **(D)** Correlation analysis between *NFIB* expression (\log_2 counts) and
560 chromatin accessibility (\log_2 RPM) at SRR124 and SRR134 regions in BRCA (n = 74),
561 LUAD (n = 21), and LUSC (n = 16) tumors. RNA-seq reads were normalized to library
562 size using DESeq2⁹¹. Significance analysis by Pearson correlation (n = 111). Bolded
563 line: fitted linear regression model. Shaded area: 95% confidence region for the
564 regression fit. **(E)** Comparison of *NFIB* expression (\log_2 counts) from BRCA, LUAD, and
565 LUSC patient tumors according to their chromatin accessibility at the SRR124 and
566 SRR134 regions. Chromatin accessibility at each region was considered “low” if \log_2
567 RPM < 1, or “high” if \log_2 RPM > 1. RNA-seq reads were normalized to library size using
568 DESeq2⁹¹. Error bars: standard deviation. Significance analysis by a two-sided t-test
569 with Holm correction. **(F)** Relative fold change (\log_2 FC) in luciferase activity driven by
570 SRR124 and SRR134 after overexpression of either *FOXA1* or *NFIB* compared to an
571 empty vector containing (mock negative control, miRFP670). Dashed line: average
572 activity of the mock control. Error bars: standard deviation. Significance analysis by

573 Tukey's test ($n = 5$; * $P < 0.05$, ** $P < 0.01$, *** $P < 0.001$, ns: not significant). **(G)**
574 Relative luciferase activity driven by WT, FOXA1-mutated, and NFIB-mutated SRR134
575 constructs compared to a minimal promoter (minP) vector in the MCF-7, PC-9, and
576 T47D cell lines. Dashed line: average activity of minP. Error bars: standard deviation.
577 Significance analysis by Tukey's test ($n = 5$; * $P < 0.05$, ** $P < 0.01$, *** $P < 0.001$, ns:
578 not significant). **(H)** RT-qPCR comparison of transcripts at SOX2, SRR124, and
579 SRR134 between sorted BFP^{-ve} and BFP^{+ve} MCF-7 cells normalized to unsorted
580 population. Error bars: standard deviation. Significance analysis by paired t-test with
581 Holm correction ($n = 6$; *** $P < 0.001$). **(I)** FACS density plot comparing tagBFP signal
582 between SOX2-P2A-tagBFP MCF-7 cells transfected with an empty vector (mock
583 negative control, miRFP670), FOXA1-T2A-miRFP670, or NFIB-T2A-miRFP670. tagBFP
584 signal was acquired from successfully transfected live cells (miRFP⁺/PI⁻) after 5 days
585 post-transfection. Significance analysis by FlowJo's chi-squared T(x) test. T(x) scores
586 above 1000 were considered "strongly significant" (*** $P < 0.001$), whereas T(x) scores
587 under 100 were considered "non-significant".

588

589 **SRR124 and SRR134 are conserved enhancers across mammals and are required**
590 **for the separation of the anterior foregut**

591 SOX2 is required for the proper development of multiple tissues³⁷, including the
592 digestive and respiratory systems in the mouse^{25,27,29,31,32,38} and in humans^{33,34}.
593 Therefore, we questioned whether the SRR124–134 cluster drives SOX2 expression in
594 additional contexts other than cancer. A compilation of chromatin accessibility data from
595 cardiac, digestive, embryonic, lymphoid, musculoskeletal, myeloid, neural, placental,

596 pulmonary, renal, skin, and vascular tissues^{80,81,110} showed that both SRR124 and
597 SRR134 display increased chromatin accessibility in digestive and respiratory tissues
598 alongside cancer samples (Figure 6A). By comparing DNase-seq signal from fetal lung
599 and stomach tissues⁸¹, we found that both SRR124 (lung $P = 1.25 \times 10^{-6}$; stomach $P =$
600 9.64×10^{-4} ; holm-adjusted Dunn's test) and SRR134 (lung $P = 1.14 \times 10^{-3}$; stomach $P =$
601 0.045), together with SRR2 (lung $P = 1.55 \times 10^{-3}$; stomach $P = 5.74 \times 10^{-5}$), are
602 significantly more accessible than pOR5K1 (Figure 6B, Supplementary Table S22). This
603 suggests that SRR124 and SRR134 are contributing to SOX2 expression during the
604 development of the digestive and respiratory systems.

605 Since critical developmental genes are often controlled by highly conserved
606 enhancers across species^{111,112}, we hypothesized that the SRR124–134 cluster might
607 regulate SOX2 expression during development in other species. By analyzing PhyloP
608 conservation scores^{94,113}, we discovered that both SRR124 and SRR134 contain a
609 highly conserved core sequence that is preserved across mammals, birds, reptiles, and
610 amphibians (Figure 6C). After aligning and comparing enhancer sequences between
611 humans and mice, we found the core sequence at both SRR124 and SRR134 are
612 highly conserved (> 80%) in the mouse genome (Supplementary Figure S6A). We
613 termed these homologous regions as mSRR96 (96 kb downstream of the mouse Sox2
614 promoter; homologous to the human SRR124) and mSRR102 (102 kb downstream of
615 the mouse Sox2 promoter; homologous to the human SRR134). Enhancer feature
616 analysis in the developing lung and stomach tissues in the mouse^{81,114} showed that
617 both mSRR96 and mSRR102 display increased chromatin accessibility and H3K27ac
618 signal throughout developmental days E14.5 to the 8th post-natal week (Figure 6D).

619 Interestingly, mSRR96 and mSRR102 display higher ATAC-seq and H3K27ac signal
620 towards the later stages of development in the lungs, but at early stages of development
621 in the stomach. This suggests a distinct spatiotemporal contribution of this homologous
622 cluster to *Sox2* expression during the development of these tissues in the mouse.
623 ATAC-seq quantification (genomic coordinates in Supplementary Table S23) showed
624 that both mSRR96 (lung $P = 5.54 \times 10^{-5}$; stomach $P = 2.37 \times 10^{-4}$; Holm-adjusted Dunn's
625 test) and mSRR102 (lung $P = 1.27 \times 10^{-3}$; stomach $P = 0.046$) are significantly more
626 accessible than the repressed promoter of the olfactory gene *Olf266* (pOlf266,
627 negative control) during the development of the lungs and stomach in the mouse
628 (Supplementary Figure S6B, Supplementary Table S24). Together, these results
629 suggest a conserved SOX2 regulatory mechanism across multiple species and support
630 a model in which the SRR124 and SRR134 enhancers and their homologs regulate
631 SOX2 expression during the development of the digestive and respiratory systems.

632 To assess the contribution of the mSRR96 and mSRR102 regions to the
633 development of the mouse, we generated a knockout containing a deletion spanning the
634 mSRR96–102 enhancer cluster (Δ mENH) (Figure 6E). We crossed animals carrying a
635 heterozygous mSRR96–102 deletion (Δ mENH^{+/-}) and determined the number of pups
636 alive at weaning (P21) from each genotype. We found a significant ($P = 3.92 \times 10^{-6}$, Chi-
637 squared test) deviation from the expected mendelian ratio, with no homozygous mice
638 (Δ mENH^{-/-}) alive at weaning (Figure 6F), demonstrating that the mSRR96–102
639 enhancer cluster is crucial for survival in the mouse. To investigate the resulting
640 phenotype in a homozygous mSRR96–102 enhancer deletion, we collected E18.5
641 embryos and prepared cross-sections at the thymus level from five animals of each

642 phenotype (WT, Δ mENH^{+/-}, and Δ mENH^{-/-}) (Figure 6G). Similar to other studies that
643 interfered with Sox2 expression during development^{22,25,32}, we found that all five
644 Δ mENH^{-/-} embryos developed EA/TEF, where the esophagus and trachea fail to
645 separate during embryonic development (Figure 6H). WT and Δ mENH^{+/-} embryos, on
646 the other hand, showed normal development of the esophageal and tracheal tissues.
647 Finally, immunohistochemistry showed the complete absence of the SOX2 protein
648 within the EA/TEF tissue in Δ mENH^{-/-} embryos, whereas WT and Δ mENH^{+/-} embryos
649 showed high levels of SOX2 protein within both the esophagus and tracheal tubes
650 (Figure 6I). Together, these results demonstrate that mSRR96 and mSRR102 are
651 imperative to drive Sox2 expression during the development of the esophagus and
652 trachea.

653

654 **Figure 6: The SRR124 and SRR134 enhancers are conserved across species and**
655 **are required for the separation of the esophagus and trachea in the mouse. (A)**

656 UCSC Genome Browser⁹⁴ view of the SOX2 region containing a compilation of
657 chromatin accessibility tracks of multiple human tissues^{80,81,110}. Arrow: increased
658 chromatin accessibility at the SRR124–134 cluster in cancer, digestive, and respiratory
659 tissues. **(B)** DNase-seq quantification (log₂ RPM) at the *RAB7A* promoter (pRAB7A),
660 *SOX2* promoter (pSOX2), SRR1, SRR2, SRR124, SRR134, human SCR (hSCR), and a
661 desert region within the *SOX2* locus (desert) compared to the background signal at the
662 repressed *OR5K1* promoter (pOR5K1) in lung and stomach embryonic tissues⁸¹.
663 Dashed line: Regions with a sum of reads above our threshold (log₂ RPM > 0) were
664 considered “accessible”. Error bars: standard deviation. Significance analysis by Dunn’s

665 test with Holm correction (* $P < 0.05$, ** $P < 0.01$, *** $P < 0.001$, ns: not significant). **(C)**
666 UCSC Genome Browser⁹⁴ with PhyloP conservation scores¹¹³ at the SRR124 and
667 SRR134 enhancers across mammals, birds, reptiles, and amphibians species. Black
668 lines: highly conserved sequences. Empty lines: variant sequences. **(D)** UCSC Genome
669 Browser⁹⁴ view of the *Sox2* region in the mouse. ATAC-seq and H3K27ac ChIP-seq
670 data from lung and stomach tissues throughout developmental days E14.5 to the 8th
671 post-natal week^{81,114}. mSRR96: homologous to SRR124. mSRR102: homologous to
672 SRR134. Reads were normalized to library size (RPM). **(E)** Illustration demonstrating
673 the mSRR96–102 enhancer cluster CRISPR deletion (Δ mENH) in C57BL/6J mouse
674 embryos. **(F)** Quantification and genotype of the C57BL/6J progeny from mSRR96–102-
675 deleted crossings (Δ mENH^{+/-}). Pups were counted and genotyped at weaning (P21).
676 The expected numbers of heterozygous and homozygous (Δ mENH^{-/-}) pups are twice
677 and equal to, respectively, the number of obtained WT animals. Significance analysis by
678 chi-squared test to measure the deviation in the number of obtained pups from the
679 expected mendelian ratio of 1:2:1 (WT : Δ mENH^{+/-} : Δ mENH^{-/-}). **(G)** Transverse cross-
680 section of fixed E18.5 embryos at the start of the thymus. **(H)** Embryo sections stained
681 with Hematoxylin and Eosin (H&E). The scale bar represents 500 μ m. Es: esophagus;
682 Tr: trachea; EA/TEF: esophageal atresia with distal tracheoesophageal fistula. **(I)**
683 Embryo sections stained with SOX2. The scale bar represents 500 μ m. Es: esophagus;
684 Tr: trachea; EA/TEF: esophageal atresia with distal tracheoesophageal fistula.

685 **DISCUSSION**

686 Our findings reveal that the SRR124–134 enhancer cluster is essential for Sox2
687 expression in the developing airway and digestive systems and is required for the
688 separation of the esophagus and trachea during mouse development. When
689 embryogenesis is complete, Sox2 expression is downregulated in most cell types as its
690 developmental enhancers are decommissioned. We propose that aberrant upregulation
691 of the pioneer factor *FOXA1* recommissions both SRR124 and SRR134 in tumor cells,
692 driving *SOX2* overexpression in breast and lung cancer. As *SOX2* is also a pioneer
693 transcription factor, increased levels of this protein further reprogram the chromatin
694 landscape of cancer cells, binding at multiple downstream regulatory regions, increasing
695 chromatin accessibility, and driving subsequent upregulation of genes associated with
696 epithelium development, ultimately supporting a tumor-initiating phenotype.

697 The observation that enhancers involved in the development of the airway and
698 digestive systems are recommissioned to support *SOX2* upregulation during
699 tumorigenesis is in line with observations that tumor-initiating cells acquire a less
700 differentiated phenotype^{115–118}. It is more surprising, however, that the *SOX2* gene is
701 regulated by common enhancers in both breast and lung cancer cells as enhancers are
702 usually highly tissue-specific^{7,111,112,119}. Our observation that *FOXA1* expression is
703 significantly correlated to chromatin accessibility at the SRR124–134 cluster and
704 increases the transcriptional output of the SRR124 and SRR134 enhancers provides a
705 mechanistic link between breast and lung developmental programs and cancer
706 progression. *FOXA1* is directly involved in the branching morphogenesis of the
707 epithelium in breast^{120,121} and lung^{122,123} tissues, where *SOX2* also plays an important

708 role^{27,59}. Overexpression of both *FOXA1*^{7–9,11,124–126} and *SOX2*^{54,65,127} have also been
709 individually linked to the activation of transcriptional programs associated with multiple
710 types of cancer. Therefore, we propose that *FOXA1* is one of the key players
711 responsible for misactivation of the SRR124–134 cluster in cancer, which then drives
712 *SOX2* overexpression in breast and lung tumors. As mutation of the *FOXA1* motif
713 disrupted SRR134 enhancer activity, and this motif is shared among other members of
714 the forkhead box (FOX) transcription factor family¹²⁸, it remains unclear if *FOXA1* alone
715 activates the SRR124–134 cluster, or whether other FOX proteins are involved in this
716 process. For example, *FOXM1* overexpression, which also showed binding at both
717 SRR124 and SRR134 in MCF-7 cells, has similarly been associated with poor patient
718 outcomes in multiple types of cancer¹²⁹.

719 In addition to the activating role of *FOXA1*, we identified *NFIB* as a negative
720 regulator of *SOX2* expression through inhibition of SRR124–134 activity. *NFIB* is
721 normally required for the development of multiple tissues^{reviewed in 130}, including the brain
722 and lungs^{131–133}, tissues in which *SOX2* expression is also tightly regulated^{27,134}. In the
723 lungs, *NFIB* is essential for promoting the maturation and differentiation of progenitor
724 cells^{131,132}. This is in stark contrast to *SOX2*, which inhibits the differentiation of lung
725 cells²⁷. Interestingly, *NFIB* seems to have paradoxical roles in cancer, acting both as a
726 tumor suppressor and as an oncogene in different tissues¹³⁵. Among its tumor
727 suppressor activity, *NFIB* acts as a barrier to skin carcinoma progression¹³⁶, and its
728 downregulation is associated with dedifferentiation and aggressiveness in LUAD¹³⁷. On
729 the other hand, *SOX2* promotes skin⁶⁵ and lung¹³⁸ cancer progression. As an
730 oncogene, *NFIB* promotes cell proliferation and metastasis in STAD¹³⁹, where *SOX2*

731 downregulation is associated with poor patient outcomes^{140–142}. With this contrasting
732 relationship between *SOX2* and *NFIB* across multiple tissues, we propose that *NFIB*
733 normally acts as a suppressor of *SRR124–134* activity and *SOX2* expression during the
734 differentiation of progenitor cells; downregulation of *NFIB* expression then results in
735 *SOX2* overexpression during tumorigenesis of the breast and lung.

736 We initially hypothesized that the neural enhancers *SRR1* and *SRR2*^{18,19,143},
737 and/or the pluripotency-associated *SCR*^{20,21} might be recommissioned during cancer
738 progression, as stem cell-related enhancers have been shown to acquire enhancer
739 features in tumorigenic cells¹⁴⁴. Although other studies have also proposed the
740 activation of either *SRR1*^{41,68} or *SRR2*^{145,146} as the main drivers of *SOX2*
741 overexpression in BRCA, we found no evidence of this mechanism and instead
742 identified the *SRR124–134* cluster as the main driver of *SOX2* expression in BRCA and
743 LUAD. Our patient tumor analysis did show that GBM and LGG were the only cancer
744 types that display a unique and consistent pattern of accessible chromatin at *SRR1* and
745 *SRR2*, which is likely related to glioma cells assuming a neural stem cell-like identity to
746 sustain high levels of cell proliferation in the brain⁶¹. In fact, *SRR2* deletion was shown
747 to downregulate *SOX2* and reduce cell proliferation in GBM cells¹⁴⁷, highlighting
748 enhancer specificity to different tumor types. In line with these findings, our observation
749 that PC-9 LUAD cells are dependent on *SRR124–134* for *SOX2* transcription, whereas
750 in H520 LUSC cells *SRR124–134* is dispensable, again highlights these tumor-type
751 specific regulatory mechanisms. LUSC tumors frequently amplify the *SOX2* locus
752^{57,58,71,72}, whereas LUAD tumors do not¹⁴⁸, indicating that different mechanisms are
753 involved in genome dysregulation in these two lung cancer subtypes. Interestingly, a

754 further downstream enhancer cluster located ~55 kb away from SRR124–134 is co-
755 amplified with *SOX2* in LUSC cell lines⁷², revealing an additional mechanism that could
756 sustain *SOX2* overexpression in the absence of the SRR124–134 cluster in certain
757 types of LUSC but not in LUAD.

758 Deletion of mSRR96–102, homolog of the human SRR124–134 cluster, resulted in
759 EA/TEF which is also observed in human cases with *SOX2* heterozygous mutations
760^{33,34}. Interestingly, a recent study showed that insertion of a CTCF insulation cluster
761 downstream of the *Sox2* gene, but upstream of mSRR96–102, disrupts *Sox2*
762 expression, impairs separation of the esophagus and trachea, and results in perinatal
763 lethality due to EA/TEF in the mouse²². This was of particular interest for understanding
764 enhancer functional nuances since the SCR, which is required for *Sox2* transcription at
765 implantation, can overcome the insulator effect of this insertion. The authors proposed
766 that enhancer density might explain the EA/TEF phenotype, as chromatin features
767 suggested that enhancers in the developing lung and stomach tissues might be spread
768 over a 400 kb domain²². The 6 kb deletion that removes the mSRR96–102 cluster
769 causing EA/TEF suggests this is not the case. Instead, we propose that the sensitivity of
770 each cell type to gene dosage is behind the differing ability of CTCF to block distal
771 enhancers. This is based on two observations: in humans, heterozygous *SOX2*
772 mutations are linked with the anophthalmia-esophageal-genital syndrome; in mice,
773 hypomorphic *Sox2* alleles display similar phenotypes in the eye²⁴ and EA/TEF^{25,32}.
774 This suggests that cells from the peri-implantation phase are less sensitive to lower
775 *Sox2* dosages compared to cells from the developing airways and digestive systems in
776 both species and explains the aberrant phenotypes observed at term.

777 Our findings illustrate how cis-regulatory regions can similarly drive gene
778 expression in both healthy and diseased contexts and serve as a prime example of how
779 developmental-associated enhancers may become misactivated in cancer. The fact that
780 we have found a digestive/respiratory-associated enhancer cluster driving gene
781 expression in a non-native context such as BRCA remains intriguing and reinforces a
782 model in which tumorigenic cells often revert to a progenitor-like state that combines
783 cis-regulatory features of progenitor cells from their own tissue compartments with those
784 of other developing lineages ⁷. This “dys-differentiation” mechanism seems to be
785 centered around the overexpression of a few key development-associated pioneer
786 transcription factors such as FOXA1. Identifying additional mechanisms that regulate
787 this enhancer recommissioning could lead to new approaches to target tumor-initiating
788 cells that depend on SOX2 overexpression.
789

790 **MATERIALS AND METHODS**

791 *Cell Culture*

792 MCF-7 cells were obtained from Eldad Zacksenhaus (Toronto General Hospital
793 Research Institute, Toronto, CA). H520 (HTB-182) and T47D (HTB-133) cells were
794 acquired from ATCC. PC-9 (90071810) cells were obtained from Sigma. Cell line
795 identities were confirmed by short tandem repeat profiling. MCF-7 and T47D cells were
796 grown in phenol red-free DMEM high glucose (Gibco), 10% FBS (Gibco), 1x Glutamax
797 (Gibco), 1x Sodium Pyruvate (Gibco), 1x Penicillin-Streptomycin (Gibco), 1x Non-
798 essential amino acids (Gibco), 25 mM HEPES (Gibco) and 0.01 mg/ml insulin (Sigma).
799 H520 and PC-9 cells were grown in phenol red-free RPMI-1640 (Gibco), 10% FBS
800 (Gibco), 1x Glutamax (Gibco), 1x Sodium Pyruvate (Gibco), 1x Penicillin-Streptomycin
801 (Gibco), 1x Non-essential amino acids (Gibco), and 25 mM HEPES (Gibco). Cells were
802 either passaged or had their medium replenished every three days.

803

804 *Genome editing*

805 Pairs of gRNA plasmids were constructed by inserting a 20 bp target sequence
806 (Supplementary Table S25) into an empty gRNA cloning vector (a gift from George
807 Church; Addgene plasmid # 41824; <http://n2t.net/addgene:41824>;
808 RRID:Addgene_41824)¹⁴⁹ containing either miRFP670 (Addgene plasmid #163748) or
809 tagBFP (Addgene plasmid #163747) fluorescent markers. Plasmids were sequenced to
810 confirm correct insertion. Both gRNA (1 µg each) vectors were co-transfected with 3 µg
811 of pCas9_GFP (a gift from Kiran Musunuru; Addgene plasmid #44719;
812 <http://n2t.net/addgene:44719>; RRID:Addgene_44719)¹⁵⁰ using Neon electroporation

813 (Life Technologies). After 72 hours of transfection, cells were FACS sorted to select
814 clones that contained all three plasmids. Sorted tagBFP⁺/GFP⁺/miRFP670⁺ cells were
815 grown in a bulk population and serially diluted into individual wells to generate isogenic
816 populations. Once fully grown, each well was screened by PCR to confirm the deletion.

817

818 *Gene tagging*

819 SOX2 was tagged with a P2A-tagBFP sequence in both alleles using CRISPR-mediated
820 homology-directed repair (HDR) ¹⁵¹. This strategy results in the expression of a single
821 transcript that is further translated into two separate proteins due to ribosomal skipping
822 ¹⁵². In summary, we designed a gRNA that targets the 3' end of the SOX2 stop codon
823 (Supplementary Table S25, Addgene plasmid #163752). We then amplified ~800 bp
824 homology arms upstream and downstream of the gRNA target sequence using high-
825 fidelity Phusion Polymerase. We purposely avoided amplification of the SOX2 promoter
826 sequence to reduce the likelihood of random integrations in the genome. Both homology
827 arms were then joined at each end of a P2A-tagBFP sequence using Gibson assembly.
828 Flanking primers containing the gRNA target sequence were used to reamplify SOX2-
829 P2A-tagBFP and add gRNA targets at both ends of the fragment; this approach allows
830 excision of the HDR sequence from the backbone plasmid once inside the cell ¹⁵³.
831 Finally, the full HDR sequence was inserted into a pJET1.2 (Thermo Scientific)
832 backbone, midipreped, and sequenced (Addgene #163751). 3µg of HDR template was
833 then co-transfected with 1µg of hCas9 (a gift from George Church; Addgene plasmid
834 #41815; <http://n2t.net/addgene:41815> ; RRID:Addgene_41815) ¹⁴⁹ and 1µg of gRNA
835 plasmid using Neon electroporation (Life Technologies). A week after transfection,

836 tagBFP⁺ cells were FACS sorted as a bulk population. Sorted cells were further grown
837 for two more weeks, and single tagBFP⁺ cells were isolated to generate isogenic
838 populations. Once fully grown, each clone was screened by PCR and sequenced to
839 confirm homozygous integration of P2A-tagBFP into the SOX2 locus.

840

841 *Luciferase assay*

842 Luciferase activity was measured using the dual-luciferase reporter assay (Promega
843 #E1960) that relies on the co-transfection of two plasmids: pGL4.23 (Firefly Luciferase,
844 *luc2*) and pGL4.75 (Renilla Luciferase). Assayed plasmids were constructed by
845 subcloning the empty pGL4.23 vector containing a minimal promoter (minP). SRR124,
846 SRR134, SRR1, SRR2, and hSCR were PCR-amplified (primers in Supplementary
847 Table S26) from MCF-7 genomic DNA using high-fidelity Phusion Polymerase and
848 inserted in the forward position downstream of the *luc2* gene at the NotI restriction site.
849 Constructs were sequenced to confirm correct insertions.

850 JASPAR2022¹⁰³ was used to find FOXA1 (GTAAACA) and NFIB
851 (TGGCAnnnnGCCAA) motifs in the SRR134 sequence. Only motifs with a score of 80%
852 or higher were further analyzed. Bases within each motif sequence were mutated until
853 the score was reduced below 80% without affecting co-occurring motifs or creating
854 novel binding sites. In total, four FOXA1 motifs and two NFIB motifs were mutated
855 (Supplementary Table S21). Engineered sequences were ordered as gene blocks
856 (Eurofins) and inserted into pGL4.23 in the forward position. Constructs were
857 sequenced to confirm correct insertions.

858 Cells were plated in 96-well plates with 4 technical replicates at $2 \cdot 10^4$ cells per well.
859 After 24 hours, a 200ng 50:1 mixture of enhancer vector and pGL4.75 was transfected
860 using Lipofectamine 3000 (0.05 μ l Lipofectamine:1 μ l Opti-mem). For transcription factor
861 overexpression analysis, a 200ng 50:10:1 mixture of enhancer vector, expression
862 plasmid, and pGL4.75 was transfected. After 48 hours of transfection, cells were lysed
863 in 1x Passive Lysis Buffer and stored at -80°C until all 5 biological replicates were
864 completed. Luciferase activity was measured in the Fluoroskan Ascent FL plate reader.
865 Enhancer activity was calculated by normalizing the firefly signal from pGL4.23 to the
866 *Renilla* signal from pGL4.75.

867

868 *Colony formation assay*

869 MCF-7 and PC-9 cells were seeded at low density (2,000 cells/well) into 6-well plates in
870 triplicate for each cell line. Culture media was renewed every 3 days. After 12 days,
871 cells were fixed with 3.7% paraformaldehyde for 10 minutes and stained with 0.5%
872 crystal violet for 20 minutes to quantify the number of colonies formed. Crystal violet
873 staining was then eluted with 10% acetic acid and absorbance was measured at 570
874 nm to evaluate cell proliferation. Each 6-well plate was considered one biological
875 replicate and the experiment was repeated five times for each cell line (n = 5).

876

877 *FACS analysis*

878 For analyzing the effects of *FOXA1* and *NFIB* overexpression, $2 \cdot 10^6$ SOX2-P2A-tagBFP
879 cells were transfected with 50nM of plasmid expressing either miRFP670 (a gift from
880 Vladislav Verkhusha; Addgene plasmid #79987; <http://n2t.net/addgene:79987>;

881 RRID:Addgene_79987), FOXA1-T2A-miRFP670 (Addgene plasmid #182335), or NFIB-
882 T2A-miRFP670 (Addgene plasmid #187222) in 5 replicates. Five days after
883 transfection, miRFP670, tagBFP, and propidium iodide (PI) (live/dead stain) signals
884 were acquired; the amount of tagBFP signal from miRFP670⁺/PI⁻ cells was compared
885 between each treatment across all replicates.

886 FlowJo's chi-squared T(x) test was used to contrast the effects of each treatment over
887 tagBFP expression; T(x) scores above 1000 were considered "strongly significant" (***),
888 whereas T(x) scores under 100 were considered "non-significant".

889

890 *Transcriptome analysis*

891 Total RNA was isolated from WT and enhancer-deleted (Δ ENH) cell lines using the
892 RNeasy kit. Genomic DNA was digested by Turbo DNase. 500-2,000ng of total RNA
893 was used in a reverse transcription reaction with random primers. cDNA was diluted in
894 H₂O and amplified in a qPCR reaction using SYBR Select Mix (primers in
895 Supplementary Table S27). Amplicons were sequenced to confirm primer specificity.
896 Gene expression was normalized to *PUM1*^{70,154,155}.

897 Total RNA was sent to The Centre for Applied Genomics (TCAG) for paired-end rRNA-
898 depleted total RNA-seq (Illumina 2500, 125 bp). Read quality was checked by fastQC,
899 trimmed using fastP¹⁵⁶ and mapped to the human genome (GRCh38/hg38) using
900 STAR 2.7¹⁵⁷. Healthy breast epithelium RNA-seq was obtained from ENCODE
901 (Supplementary Table S28)^{80,81}. Mapped reads were quantified using featureCounts¹⁵⁸
902 and imported into DESeq2⁹¹ for normalization and differential expression analysis.
903 Genes with a $|\log_2 \text{FC}| > 1$ and FDR-adjusted $Q < 0.01$ were considered significantly

904 changing. Differential gene expression was plotted using the EnhancedVolcano
905 package. Correlation and clustering heatmaps were plotted using the pheatmap R
906 package (<https://cran.r-project.org/web/packages/pheatmap/index.html>). Signal
907 enrichment plot was prepared using NGS.plot¹⁵⁹.

908 Cancer patient transcriptome data were obtained from TCGA⁶⁹ using the TCGAbiolinks
909 package¹⁶⁰. The overall survival KM-plot was calculated using clinical information from
910 TCGA¹⁶¹. Tumor transcriptome data were compared to healthy tissue using DESeq2.
911 RNA-seq reads were normalized to library size using DESeq2⁹¹ and transformed to a
912 log₂ scale [log₂ counts]. Differential gene expression was considered significant if |log₂
913 FC| > 1 and Q < 0.01.

914 Gene Set Enrichment Analysis (GSEA) was performed by ranking genes according to
915 their log₂ FC in Δ ENH^{-/-} versus WT MCF-7 cells. The ranking was then analyzed using
916 the GSEA function from the clusterProfiler package¹⁰⁵ with a threshold of FDR-adjusted
917 Q < 0.05 using the MSigDB GO term database (C5).

918

919 *Chromatin accessibility analysis*

920 Cells were grown in three separate wells (n = 3) and 50,000 cells were sent to Princess
921 Margaret Genomics Centre for ATAC-seq library preparation using the Omni-ATAC
922 protocol¹⁶². ATAC-seq libraries were sequenced using 50 bp paired-ended parameters
923 in the Illumina Novaseq 6000 platform. Read quality was checked by fastQC, trimmed
924 using fastP and mapped to the human genome (GRCh38/hg38) using STAR 2.7.
925 narrowPeaks were called using Genrich (<https://github.com/jsh58/Genrich>). Differential
926 chromatin accessibility analysis was performed using diffBind¹⁰⁶. ATAC-seq peaks with

927 a $|\log_2 \text{FC}| > 1$ and FDR-adjusted $Q < 0.01$ were considered significantly changing.
928 Correlation heatmaps were generated using diffBind. Signal enrichment plot was
929 prepared using NGS.plot¹⁵⁹. Genes were separated into three categories according to
930 their expression levels in our WT MCF-7 RNA-seq data.
931 Transcription factor footprint analysis was performed using TOBIAS¹⁰² with standard
932 settings. Motifs with a $|\log_2 \text{FC}| > 0.1$ and FDR-adjusted $Q < 0.01$ were considered
933 significantly enriched in each condition. Replicates ($n = 3$) were merged into a single
934 BAM file for each treatment. Motif enrichment at differential ATAC-seq peaks was
935 performed using HOMER¹⁶³. ATAC-seq peaks were assigned to their closest gene
936 within ± 1 Mb distance from their promoter using ChIPpeakAnno¹⁶⁴.
937 Cancer patient ATAC-seq data was obtained from TCGA¹⁰⁷. DNase-seq from human
938 developing tissues were obtained from ENCODE (Supplementary Table S28)^{80,81}.
939 Read quantification was calculated at the *RAB7a* (pRAB7a), *OR5K1* (pOR5K1), and
940 *SOX2* (pSOX2) promoters, together with SRR1, SRR2, SRR124, SRR134, hSCR, and
941 desert regions with a 1,500 bp window centered at the core of each region (genomic
942 coordinates of each region in Supplementary Table S15). Reads were normalized to
943 library size (RPM) and transformed to a \log_2 scale (\log_2 RPM) using a custom script
944 (<https://github.com/luisabatti/BAMquantify>). Each region's average \log_2 RPM was
945 compared to the *OR5K1* promoter for differential analysis using Dunn's test with Holm
946 correction. Correlations were calculated using Pearson's correlation test and considered
947 significant if FDR-adjusted $Q < 0.05$. Chromatin accessibility at SRR124 and SRR134
948 regions was considered low if $\log_2 \text{RPM} < -1$, medium if $-1 \leq \log_2 \text{RPM} \leq 1$, or high if \log_2
949 $\text{RPM} > 1$.

950 ATAC-seq from developing mouse lung and stomach tissues were obtained from
951 ENCODE (Supplementary Table S28)⁸¹ and others¹¹⁴. Conserved mouse regulatory
952 regions were lifted from the human build (GRCh38/hg38) to the mouse build
953 (GRCm38/mm10) using UCSC liftOver⁹⁴. The number of mapped reads was calculated
954 at the *Egf* (pEgf), *Olfir266* (pOlfir266), and *Sox2* (pSox2) promoters, together with the
955 mouse mSRR1, mSRR2, mSRR96, mSRR102, mSCR and desert regions with a 1,500
956 bp window at each location (genomic coordinates in Supplementary Table S23). Each
957 log₂-transformed region's reads per million (log₂ RPM) was compared to the negative
958 *Olfir266* promoter control for differential analysis using Dunn's test with Holm correction.

959

960 *Conservation analysis*

961 Cross-species evolutionary conservation was obtained using phyloP¹¹³. Pairwise
962 comparisons between human SRR124 and SRR134 (GRCh38/hg38) and mouse
963 mSRR96 and mSRR102 (GRCm38/mm10) sequences were plotted using FlexiDot¹⁶⁵
964 with an 80% conservation threshold.

965

966 *ChIP-seq analysis*

967 Transcription factor and histone modifications ChIP-seq were obtained from ENCODE⁸¹
968 (Supplementary Table S28) and others⁸⁸⁻⁹⁰ (Supplementary Table S29). H3K4me1 and
969 H3K27ac tracks were normalized to input and library size (log₂ RPM). ATAC-seq reads
970 were normalized to library size (RPM). Histone modification ChIP-seq tracks and
971 transcription factor ChIP-seq peaks were uploaded to the UCSC browser⁹⁴ for
972 visualization. Normalized H3K4me1, H3K27ac and ATAC-seq reads were quantified

973 and the difference in normalized signal was calculated using diffBind. Peaks with a $|\log_2$
974 $FC| > 1$ and $Q < 0.01$ were considered significantly changing.

975 Overlapping ChIP-seq and ATAC-seq peaks were analyzed using ChIPpeakAnno¹⁶⁴.
976 The hypergeometric test was performed by comparing the number of overlapping peaks
977 to the total size of the genome divided by the median peak size.

978

979 *Mouse line construction*

980 Our mSRR96–102 knockout mouse line (Δ mENH) was ordered from and generated by
981 The Centre for Phenogenomics (TCP) in Toronto, ON. The protocol for the generation
982 of the mouse line has been previously described¹⁶⁶. Briefly, C57BL/6J zygotes were
983 collected from superovulated, mated, and plugged female mice at 0.5-day post coitum.
984 Zygotes were electroporated with Cas9 RNP complexes (gRNA sequences in
985 Supplementary Table S25) and transferred into pseudopregnant female recipients
986 within 3-4 hours of electroporation. Born pups (founders) were screened by end-point
987 PCR and sequenced to confirm allelic mSRR96–102 deletions. Heterozygous
988 mSRR96–102 founders (Δ mENH^{+/-}) were then backcrossed to the parental strain to
989 confirm germline transmission. Once the mouse line was established and the mSRR96–
990 102 deletion was fully confirmed and sequenced in the N1 offspring, Δ mENH^{+/-} mice
991 were crossed and the number of live pups from each genotype (WT, Δ mENH^{+/-},
992 Δ mENH^{-/-}) was assessed at weaning (P21). The obtained number of live pups from
993 each genotype was then compared to the expected mendelian ratio of 1:2:1 (WT :
994 Δ mENH^{+/-} : Δ mENH^{-/-}) using a chi-squared test. Once the lethality of the homozygous
995 deletion was confirmed at weaning, E18.5 embryos generated from new Δ mENH^{+/-}

996 crosses were collected for further histological analyses. All procedures involving
997 animals were performed in compliance with the Animals for Research Act of Ontario
998 and the Guidelines of the Canadian Council on Animal Care. The TCP Animal Care
999 Committee reviewed and approved all procedures conducted on animals at the facility.

1000

1001 *Histological analyses*

1002 A total of 46 embryos were collected at E18.5 and fixed in 4% paraformaldehyde. Each
1003 embryo was genotyped, and a total of 15 embryos, 5 of each genotype (WT, Δ mENH^{+/-},
1004 Δ mENH^{-/-}), were randomly selected, processed, and embedded in paraffin for
1005 sectioning and further analysis. Tissue sections were collected at 4 μ m thickness roughly
1006 at the start of the thymus. Sections were prepared by the Pathology Core at TCP.

1007 Tissue sections were stained with Hematoxylin and Eosin (H&E) using an auto-stainer
1008 to ensure batch consistency. Slides were scanned using a Hamamatsu Nanozoomer
1009 slide scanner at 20X magnification. Images were then cropped and centered around the
1010 esophageal (Es) and tracheal (Tr) tissues.

1011 Embryo sections were submitted to heat-induced epitope retrieval with TRIS-EDTA (pH
1012 9.0) for 10 minutes, followed by quenching of endogenous peroxidase with Bloxall
1013 reagent (Vector). Non-specific antibody binding was blocked with 2.5 % normal horse
1014 serum (Vector), followed by incubation for 1 hour in Rabbit anti-SOX2 (Abcam,
1015 ab92494, 1:500). After washes, sections were incubated for 30 minutes with
1016 ImmPRESS Anti-Rabbit HRP (Vector) followed by DAB reagent, and counterstained in
1017 Mayer's hematoxylin.

1018 **SUPPORTING INFORMATION**

1019 **Supplementary Figure S1: (A)** Super-logarithmic volcano plot of *PUM1* expression
1020 from RNA-seq of 21 cancer types compared to normal tissue⁶⁹. Cancer types with log₂
1021 FC > 1 and FDR-adjusted Q < 0.01 were considered to significantly overexpress *PUM1*.
1022 Error bars: standard error. **(B)** Kaplan-Meier plot¹⁶⁷ of overall survival against time
1023 since diagnosis for 3,064 patients with BRCA (n = 1089), COAD (n = 453), GBM (n =
1024 153), LIHC (n = 370), LUAD (n = 504), and LUSC (n = 495) tumors¹⁶¹. We divided
1025 patients into four equal groups and compared two groups: high *SOX2* expression
1026 (range: 10.06–16.36 log₂ counts) and low *SOX2* expression (range: 0–1.67 log₂ counts).
1027 RNA-seq reads were normalized to library size using DESeq2⁹¹. Significance analysis
1028 by logrank test. The shadowed area represents the 95% confidence interval. **(C)**
1029 Comparison of *SOX2* expression (log₂ counts) between luminal A (n = 560), luminal B
1030 (n = 207), HER2+ (n = 82), basal-like (n = 190) breast cancer subtypes and normal
1031 mammary tissue (n = 152)⁶⁹. RNA-seq reads were normalized to library size using
1032 DESeq2⁹¹. Error bars: standard deviation. Significance analysis by Tukey's test (***P* <
1033 0.001, * *P* < 0.05, ns: not significant). **(D)** Volcano plot with DESeq2⁹¹ differential
1034 expression analysis between WT MCF-7 cells and breast epithelium⁸⁰. Blue: 7,937
1035 genes that significantly lost expression (log₂ FC < -1; FDR-adjusted Q < 0.01) in WT
1036 MCF-7 cells. Pink: 5,335 genes that significantly gained expression (log₂ FC > 1; Q <
1037 0.01) in WT MCF-7 cells. Grey: 25,342 genes that maintained similar (-1 ≤ log₂ FC ≤ 1)
1038 expression between WT MCF-7 and breast epithelium cells. **(E)** RT-qPCR analysis of
1039 *SOX2* transcript levels in the MCF-7, T47D, PC-9 and H520 cell lines. Error bars:
1040 standard deviation. Significance analysis by Tukey's test (n = 3; *** *P* < 0.001).

1041
1042 **Supplementary Figure S2: (A)** SOX2 protein levels in mouse embryonic stem cells
1043 (mESC, positive control), WT, $\Delta\text{ENH}^{+/-}$, and $\Delta\text{ENH}^{-/-}$ MCF-7 clones. Cyclophilin A
1044 (CypA) was used as a loading control across all samples. **(B)** RT-qPCR analysis of
1045 SOX2 transcript levels in SRR124–134 heterozygous- ($\Delta\text{ENH}^{+/-}$) and homozygous-
1046 ($\Delta\text{ENH}^{-/-}$) deleted H520 (LUSC) clones compared to WT cells. Error bars: standard
1047 deviation. Significance analysis by Dunnett's test ($n = 3$, ns: not significant). **(C)**
1048 Euclidean distance pairwise comparison between WT and $\Delta\text{ENH}^{-/-}$ MCF-7 replicates (n
1049 = 3) using variance stabilizing transformed RNA-seq reads from DESeq2⁹¹. Darker
1050 colors indicate a higher correlation. **(D)** Euclidean hierarchical clustering of 529
1051 differentially expressed genes ($|\log_2 \text{FC}| > 1$; FDR-adjusted $Q < 0.01$) based on RNA-
1052 seq analysis between WT and $\Delta\text{ENH}^{-/-}$ MCF-7 replicates ($n = 3$). Reads were
1053 normalized for each gene across treatments (Z-score). Blue color indicates
1054 downregulated genes (Z-score < 0). Red color indicates upregulated genes (Z-score $>$
1055 0).

1056
1057 **Supplementary Figure S3: (A)** ATAC-seq metagene enrichment plot ± 2 kb around the
1058 transcription start site (TSS) across all genes from WT and $\Delta\text{ENH}^{-/-}$ MCF-7 cells ($n =$
1059 3). Reads were normalized by library size (RPM). Grey: TSS. Shaded area: standard
1060 deviation. **(B)** ATAC-seq metagene enrichment plot ± 2 kb around the transcription start
1061 site (TSS) across 12,167 high (average \log_2 counts = 9.12), 12,167 medium (average
1062 \log_2 counts = 2.94), and 12,167 low (average \log_2 counts = 0.43) expressed genes in
1063 WT MCF-7 cells. Genes were split into each group according to RNA-seq data. RNA-

1064 seq reads were normalized to library size using DESeq2⁹¹. Grey: TSS. **(C)** ATAC-seq
1065 metagene enrichment plot \pm 2 kb around the transcription start site (TSS) across 12,167
1066 high (average \log_2 counts = 9.10), 12,167 medium (average \log_2 counts = 2.97), and
1067 12,167 low (average \log_2 counts = 0.47) expressed genes in Δ ENH^{-/-} MCF-7 cells.
1068 Genes were split into each group according to RNA-seq data. RNA-seq reads were
1069 normalized to library size using DESeq2⁹¹. Grey: TSS. **(D)** Pairwise Pearson correlation
1070 comparison between WT and Δ ENH^{-/-} MCF-7 replicates (n = 3) using ATAC-seq
1071 normalized signal from diffBind¹⁰⁶. Darker colors indicate a higher correlation. **(E)**
1072 Overlap between GRHL2 ChIP-seq peaks⁸⁹ and ATAC-seq peaks that significantly
1073 (\log_2 FC < -1; P < 0.01) lost chromatin accessibility in Δ ENH^{-/-} MCF-7 cells.
1074 Significance analysis by the hypergeometric test¹⁶⁴. **(F)** Overlap between RUNX2 ChIP-
1075 seq peaks¹⁰⁴ and ATAC-seq peaks that significantly (\log_2 FC < -1; P < 0.01) lost
1076 chromatin accessibility in Δ ENH^{-/-} MCF-7 cells. Significance analysis by the
1077 hypergeometric test¹⁶⁴.

1078

1079 **Supplementary Figure S4: (A)** ATAC-seq signal at the *RAB7A* promoter (pRAB7A),
1080 *SOX2* promoter (pSOX2), *SRR1*, *SRR2*, *SRR124*, *SRR134*, human SCR (hSCR) and
1081 desert region versus the background signal at the repressed *OR5K1* promoter
1082 (pOR5K1) in BLCA (n = 10), BRCA (n = 74), COAD (n = 38), ESCA (n = 18), GBM (n =
1083 9), HNSC (n = 9), LGG (n = 13), LIHC (n = 16), LUAD (n = 22), LUSC (n = 16), PRAD (n
1084 = 26), STAD (n = 21), TGCT (n = 9), and UCEC (n = 13) patient tumors. Dashed line:
1085 regions with \log_2 RPM > 0 were considered “accessible”. Error bars: standard deviation.
1086 Significance analysis by Dunn’s test with Holm correction (* P < 0.05, ** P < 0.01, *** P

1087 < 0.001, ns: not significant). **(B)** ATAC-seq signal at the SOX2 desert region (desert)
1088 against ATAC-seq signal for the SOX2 promoter (pSOX2) from 74 BRCA, 22 LUAD,
1089 and 16 LUSC patient tumors. Dashed line: regions with \log_2 RPM > 0 were considered
1090 “accessible”. Significance analysis by Pearson correlation. Bolded line: fitted linear
1091 regression model. Shaded area: 95% confidence region for the regression fit. **(C)**
1092 ATAC-seq signal at SRR124 and SRR134 regions against ATAC-seq signal for the
1093 SOX2 promoter (pSOX2) from BRCA patient tumors separated into luminal A (n = 31),
1094 luminal B (n = 16), HER2⁺ (n = 10), and basal-like (n = 14) subtypes. Correlation is
1095 shown for accessible chromatin (\log_2 RPM > 0). Grey: tumors with closed chromatin
1096 (\log_2 RPM < 0) at either region, not included in the correlation analysis. Significance
1097 analysis by Pearson correlation. Bolded line: fitted linear regression model. Shaded
1098 area: 95% confidence region for the regression fit.

1099
1100 **Supplementary Figure S5: (A)** RT-qPCR analysis of *FOXA1* and *NFIB* expression in
1101 the H520, MCF-7, PC-9, and T47D cell lines. Error bars: standard deviation.
1102 Significance analysis by Tukey’s test (n = 3; ** $P < 0.01$, *** $P < 0.001$). **(B)** FACS plot of
1103 tagBFP signal (450 nm) over side scatter (SSC) in WT and SOX2-P2A-tagBFP MCF-7
1104 cells. Cell populations within the top 10% tagBFP signal were considered “tagBFP
1105 positive” (BFP^{+ve}), whereas populations within the bottom 10% BFP signal were
1106 considered “tagBFP negative” (BFP^{-ve}).

1107
1108 **Supplementary Figure S6: (A)** Dot-plot alignment of human (GRCh38/hg38, y-axis)
1109 SRR124 and SRR134, and mouse (GRCm38/mm10, x-axis) mSRR96 and mSRR102

1110 homologous sequences (1,500 bp). Lines indicate high conservation scores (> 80%)
1111 across both species. Sequence alignment using Clustal Omega¹⁶⁸. **(B)** ATAC-seq
1112 quantification (\log_2 RPM) at the promoter of the housekeeping gene *Egf* (pEgf, positive
1113 control), *Sox2* promoter (pSox2), mSRR1, mSRR2, mSRR96, mSRR102, mSCR, and a
1114 mouse desert (mdesert) region compared to the background signal at the repressed
1115 *Olf266* promoter (pOlf266) in lung and stomach embryonic tissues from the mouse⁸¹.
1116 mSRR96: homologous to SRR124. mSRR102: homologous to SRR134. Dashed line:
1117 regions with a sum of reads above our threshold (\log_2 RPM > 0) were considered
1118 “accessible”. Error bars: standard deviation. Significance analysis by Dunn’s test with
1119 Holm correction (* $P < 0.05$, ** $P < 0.01$, *** $P < 0.001$, ns: not significant).
1120

- 1121 **Supplementary Table S1:** List of TCGA tumor type abbreviations.
- 1122 **Supplementary Table S2:** SOX2 differential expression analysis between primary
1123 tumor vs. normal tissue across TCGA cancer types.
- 1124 **Supplementary Table S3:** *PUM1* differential expression analysis between primary
1125 tumor vs. normal tissue across TCGA cancer types.
- 1126 **Supplementary Table S4:** TCGA cancer patient overall survival analysis relative to
1127 SOX2 expression levels.
- 1128 **Supplementary Table S5:** TCGA copy number variation (CNV) and SOX2 expression
1129 analysis.
- 1130 **Supplementary Table S6:** RNA-seq differential expression analysis between WT MCF-
1131 7 vs. breast epithelium (ENCODE).
- 1132 **Supplementary Table S7:** Differential ATAC-seq, H3K4me1, and H3K27ac analysis
1133 within ± 1 Mb of the SOX2 gene in WT MCF-7 vs. Breast epithelium (ENCODE).
- 1134 **Supplementary Table S8:** RNA-seq differential expression analysis comparing ΔENH^{-}
1135 $^{-}$ versus WT MCF-7 cells.
- 1136 **Supplementary Table S9:** Gene set enrichment analysis (GSEA) in WT versus ΔENH^{-}
1137 $^{-}$ MCF-7 cells.
- 1138 **Supplementary Table S10:** Significantly changing ATAC-seq peaks in ΔENH^{-} versus
1139 WT MCF-7 cells.
- 1140 **Supplementary Table S11:** ATAC-seq peaks that commonly gained signal in WT MCF-
1141 7 vs. breast epithelium and lost signal in ΔENH^{-} MCF-7 cells.
- 1142 **Supplementary Table S12:** ATAC-seq footprint analysis in ΔENH^{-} vs. WT MCF-7
1143 cells.

1144 **Supplementary Table S13:** ChIP-seq motif analysis of GRHL2 peaks in WT MCF-7
1145 cells.

1146 **Supplementary Table S14:** ChIP-seq motif analysis of RUNX2 peaks in WT MCF-7
1147 cells.

1148 **Supplementary Table S15:** Coordinates of regions used in genome-wide analysis in
1149 humans (GRCh38/hg38).

1150 **Supplementary Table S16:** Chromatin accessibility analysis across TCGA cancer
1151 types.

1152 **Supplementary Table S17:** ATAC-seq quantification used to separate patient tumors
1153 into expression groups and their *SOX2* expression levels.

1154 **Supplementary Table S18:** Significantly correlated transcription factors to accessible
1155 chromatin at the SRR124–134 cluster in BRCA, LUAD, and LUSC tumors.

1156 **Supplementary Table S19:** *FOXA1* transcript levels and chromatin accessibility at the
1157 SRR124–134 cluster in BRCA, LUAD, and LUSC patient tumors.

1158 **Supplementary Table S20:** *NFIB* transcript levels and chromatin accessibility at the
1159 SRR124–134 cluster in BRCA, LUAD, and LUSC patient tumors.

1160 **Supplementary Table S21:** WT, *FOXA1*, and *NFIB* mutated SRR134 sequences.

1161 **Supplementary Table S22:** Chromatin accessibility analysis in human (GRCh38/hg38)
1162 lung and stomach embryonic tissues.

1163 **Supplementary Table S23:** Coordinates of regions used in genome-wide analysis in
1164 the mouse (GRCm38/mm10).

1165 **Supplementary Table S24:** Chromatin accessibility analysis in mouse
1166 (GRCm38/mm10) embryonic lung and stomach tissues

- 1167 **Supplementary Table S25:** List of gRNA sequences used for CRISPR/Cas9.
- 1168 **Supplementary Table S26:** List of primers used for enhancer cloning.
- 1169 **Supplementary Table S27:** List of primers used in RT-qPCR experiments.
- 1170 **Supplementary Table S28:** ENCODE datasets used in this paper.
- 1171 **Supplementary Table S29:** GEO datasets used in this paper.
- 1172

1173 **DATA AVAILABILITY**

1174 Sequencing and processed data files were submitted to the Gene Expression Omnibus
1175 (GEO; <https://www.ncbi.nlm.nih.gov/geo/>) repository (GSE132344).

1176

1177 **ACKNOWLEDGMENTS**

1178 We thank all the members of the Mitchell laboratory for helpful discussions and support.

1179 We also thank the ENCODE Consortium and the TCGA project for generating and

1180 releasing data to the scientific community. Finally, we thank the contribution of the staff

1181 at TCP (The Centre for Phenogenomics), including Kyle Robertson, who handled the

1182 embedding, cutting, and H&E staining of E18.5 mouse embryos, and Vivian Bradaschia,

1183 responsible for the IHC staining. This work was supported by the Canadian Institutes of

1184 Health Research (FRN PJT153186 and PJT180312), the Canada Foundation for

1185 Innovation, and the Ontario Ministry of Research and Innovation (operating and

1186 infrastructure grants held by J.A.M.). BioRender.com was used to create parts of

1187 Figures 6E and 6G.

1188

1189 **AUTHOR CONTRIBUTIONS**

1190 L.E.A. designed and performed bioinformatic analyses, cell culture work, CRISPR

1191 deletions, data curation, gene expression quantification and led the conceptualization

1192 and writing of the manuscript; P.L.F. assessed cellular phenotypes, including the colony

1193 formation assay; L.H. acquired and processed TCGA ATAC-seq data and assisted in

1194 writing review & editing; M.C. assisted in the writing review & editing; M.M.H. provided

1195 TCGA data access and assisted in writing review & editing; J.A.M. was involved in

1196 supervision, funding acquisition, data interpretation, experimental design and writing
1197 review & editing. All authors have participated in the editing and approval of the
1198 manuscript.

1199 **REFERENCES**

- 1200 1. Spitz, F. & Furlong, E. E. M. Transcription factors: from enhancer binding to
1201 developmental control. *Nat Rev Genet* **13**, 613–626 (2012).
- 1202 2. Zhu, J. *et al.* Genome-wide chromatin state transitions associated with
1203 developmental and environmental cues. *Cell* **152**, 642–654 (2013).
- 1204 3. Hawkins, R. D. *et al.* Distinct epigenomic landscapes of pluripotent and lineage-
1205 committed human cells. *Cell Stem Cell* **6**, 479–491 (2010).
- 1206 4. Rada-Iglesias, A. *et al.* A unique chromatin signature uncovers early developmental
1207 enhancers in humans. *Nature* **470**, 279–283 (2011).
- 1208 5. Creighton, M. P. *et al.* Histone H3K27ac separates active from poised enhancers
1209 and predicts developmental state. *Proc. Natl. Acad. Sci. U.S.A.* **107**, 21931–21936
1210 (2010).
- 1211 6. Rubin, A. J. *et al.* Lineage-specific dynamic and pre-established enhancer-
1212 promoter contacts cooperate in terminal differentiation. *Nat Genet* **49**, 1522–1528
1213 (2017).
- 1214 7. Stergachis, A. B. *et al.* Developmental Fate and Cellular Maturity Encoded in
1215 Human Regulatory DNA Landscapes. *Cell* **154**, 888–903 (2013).
- 1216 8. Fu, X. *et al.* FOXA1 upregulation promotes enhancer and transcriptional
1217 reprogramming in endocrine-resistant breast cancer. *Proc Natl Acad Sci USA* **116**,
1218 26823–26834 (2019).
- 1219 9. Roe, J.-S. *et al.* Enhancer Reprogramming Promotes Pancreatic Cancer
1220 Metastasis. *Cell* **170**, 875-888.e20 (2017).

- 1221 10. Bi, M. *et al.* Enhancer reprogramming driven by high-order assemblies of
1222 transcription factors promotes phenotypic plasticity and breast cancer endocrine
1223 resistance. *Nat Cell Biol* **22**, 701–715 (2020).
- 1224 11. Lupien, M. *et al.* FoxA1 Translates Epigenetic Signatures into Enhancer-Driven
1225 Lineage-Specific Transcription. *Cell* **132**, 958–970 (2008).
- 1226 12. Avilion, A. A. Multipotent cell lineages in early mouse development depend on
1227 SOX2 function. *Genes & Development* **17**, 126–140 (2003).
- 1228 13. Chew, J.-L. *et al.* Reciprocal Transcriptional Regulation of Pou5f1 and Sox2 via the
1229 Oct4/Sox2 Complex in Embryonic Stem Cells. *Mol Cell Biol* **25**, 6031–6046 (2005).
- 1230 14. Takahashi, K. & Yamanaka, S. Induction of Pluripotent Stem Cells from Mouse
1231 Embryonic and Adult Fibroblast Cultures by Defined Factors. *Cell* **126**, 663–676
1232 (2006).
- 1233 15. Takahashi, K. *et al.* Induction of Pluripotent Stem Cells from Adult Human
1234 Fibroblasts by Defined Factors. *Cell* **131**, 861–872 (2007).
- 1235 16. Yu, J. *et al.* Induced Pluripotent Stem Cell Lines Derived from Human Somatic
1236 Cells. *Science* **318**, 1917–1920 (2007).
- 1237 17. Ferri, A. L. M. *et al.* Sox2 deficiency causes neurodegeneration and impaired
1238 neurogenesis in the adult mouse brain. *Development* **131**, 3805–3819 (2004).
- 1239 18. Tomioka, M. *et al.* Identification of Sox-2 regulatory region which is under the
1240 control of Oct-3/4–Sox-2 complex. *Nucleic Acids Res* **30**, 3202–3213 (2002).
- 1241 19. Zappone, M. V. *et al.* Sox2 regulatory sequences direct expression of a (beta)-geo
1242 transgene to telencephalic neural stem cells and precursors of the mouse embryo,

- 1243 revealing regionalization of gene expression in CNS stem cells. *Development* **127**,
1244 2367–2382 (2000).
- 1245 20. Zhou, H. Y. *et al.* A Sox2 distal enhancer cluster regulates embryonic stem cell
1246 differentiation potential. *Genes & Development* **28**, 2699–2711 (2014).
- 1247 21. Li, Y. *et al.* CRISPR Reveals a Distal Super-Enhancer Required for Sox2
1248 Expression in Mouse Embryonic Stem Cells. *PLoS ONE* **9**, e114485 (2014).
- 1249 22. Chakraborty, S. *et al.* Enhancer–promoter interactions can bypass CTCF-mediated
1250 boundaries and contribute to phenotypic robustness. *Nat Genet* (2023)
1251 doi:10.1038/s41588-022-01295-6.
- 1252 23. Graham, V., Khudyakov, J., Ellis, P. & Pevny, L. SOX2 Functions to Maintain
1253 Neural Progenitor Identity. *Neuron* **39**, 749–765 (2003).
- 1254 24. Taranova, O. V. *et al.* SOX2 is a dose-dependent regulator of retinal neural
1255 progenitor competence. *Genes Dev* **20**, 1187–1202 (2006).
- 1256 25. Que, J. *et al.* Multiple dose-dependent roles for Sox2 in the patterning and
1257 differentiation of anterior foregut endoderm. *Development* **134**, 2521–2531 (2007).
- 1258 26. Kiernan, A. E. *et al.* Sox2 is required for sensory organ development in the
1259 mammalian inner ear. *Nature* **434**, 1031–1035 (2005).
- 1260 27. Gontan, C. *et al.* Sox2 is important for two crucial processes in lung development:
1261 Branching morphogenesis and epithelial cell differentiation. *Developmental Biology*
1262 **317**, 296–309 (2008).
- 1263 28. Driskell, R. R., Giangreco, A., Jensen, K. B., Mulder, K. W. & Watt, F. M. Sox2-
1264 positive dermal papilla cells specify hair follicle type in mammalian epidermis.
1265 *Development* **136**, 2815–2823 (2009).

- 1266 29. Francis, R. *et al.* Gastrointestinal transcription factors drive lineage-specific
1267 developmental programs in organ specification and cancer. *Sci. Adv.* **5**, eaax8898
1268 (2019).
- 1269 30. Okubo, T., Pevny, L. H. & Hogan, B. L. M. Sox2 is required for development of
1270 taste bud sensory cells. *Genes & Development* **20**, 2654–2659 (2006).
- 1271 31. Que, J., Luo, X., Schwartz, R. J. & Hogan, B. L. M. Multiple roles for Sox2 in the
1272 developing and adult mouse trachea. *Development* **136**, 1899–1907 (2009).
- 1273 32. Teramoto, M. *et al.* The absence of SOX2 in the anterior foregut alters the
1274 esophagus into trachea and bronchi in both epithelial and mesenchymal
1275 components. *Biology Open* **9**, bio048728 (2020).
- 1276 33. Zenteno, J. C., Perez-Cano, H. J. & Aguinaga, M. Anophthalmia-esophageal
1277 atresia syndrome caused by an SOX2 gene deletion in monozygotic twin brothers
1278 with markedly discordant phenotypes. *Am J Med Genet A* **140**, 1899–1903 (2006).
- 1279 34. Williamson, K. A. *et al.* Mutations in SOX2 cause anophthalmia-esophageal-genital
1280 (AEG) syndrome. *Hum Mol Genet* **15**, 1413–1422 (2006).
- 1281 35. Brunner, H. G. & van Bokhoven, H. Genetic players in esophageal atresia and
1282 tracheoesophageal fistula. *Curr Opin Genet Dev* **15**, 341–347 (2005).
- 1283 36. Que, J., Choi, M., Ziel, J. W., Klingensmith, J. & Hogan, B. L. M. Morphogenesis of
1284 the trachea and esophagus: current players and new roles for noggin and Bmps.
1285 *Differentiation* **74**, 422–437 (2006).
- 1286 37. Arnold, K. *et al.* Sox2(+) adult stem and progenitor cells are important for tissue
1287 regeneration and survival of mice. *Cell Stem Cell* **9**, 317–329 (2011).

- 1288 38. Tompkins, D. H. *et al.* Sox2 Is Required for Maintenance and Differentiation of
1289 Bronchiolar Clara, Ciliated, and Goblet Cells. *PLoS ONE* **4**, e8248 (2009).
- 1290 39. Wuebben, E. L. & Rizzino, A. The dark side of SOX2: cancer-a comprehensive
1291 overview. *Oncotarget* **8**, 44917–44943 (2017).
- 1292 40. Chen, Y. *et al.* The Molecular Mechanism Governing the Oncogenic Potential of
1293 SOX2 in Breast Cancer. *Journal of Biological Chemistry* **283**, 17969–17978 (2008).
- 1294 41. Leis, O. *et al.* Sox2 expression in breast tumours and activation in breast cancer
1295 stem cells. *Oncogene* **31**, 1354–1365 (2012).
- 1296 42. Liu, P. *et al.* SOX2 Promotes Cell Proliferation and Metastasis in Triple Negative
1297 Breast Cancer. *Front Pharmacol* **9**, (2018).
- 1298 43. Meng, Y., Xu, Q., Chen, L., Wang, L. & Hu, X. The function of SOX2 in breast
1299 cancer and relevant signaling pathway. *Pathology - Research and Practice* **216**,
1300 153023 (2020).
- 1301 44. Piva, M. *et al.* Sox2 promotes tamoxifen resistance in breast cancer cells. *EMBO*
1302 *Molecular Medicine* **6**, 66–79 (2014).
- 1303 45. Takeda, K. *et al.* Sox2 is associated with cancer stem-like properties in colorectal
1304 cancer. *Scientific Reports* **8**, 17639 (2018).
- 1305 46. Talebi, A., Kianersi, K. & Beiraghdar, M. Comparison of gene expression of SOX2
1306 and OCT4 in normal tissue, polyps, and colon adenocarcinoma using
1307 immunohistochemical staining. *Adv Biomed Res* **4**, 234 (2015).
- 1308 47. Zhang, X.-H., Wang, W., Wang, Y.-Q., Zhu, L. & Ma, L. The association of SOX2
1309 with clinical features and prognosis in colorectal cancer: A meta-analysis. *Pathol*
1310 *Res Pract* **216**, (2020).

- 1311 48. Zhu, Y. *et al.* SOX2 promotes chemoresistance, cancer stem cells properties, and
1312 epithelial-mesenchymal transition by β -catenin and Beclin1/autophagy signaling in
1313 colorectal cancer. *Cell Death Dis* **12**, 449 (2021).
- 1314 49. Alonso, M. M. *et al.* Genetic and Epigenetic Modifications of Sox2 Contribute to the
1315 Invasive Phenotype of Malignant Gliomas. *PLoS ONE* **6**, e26740 (2011).
- 1316 50. Cox, J. L., Wilder, P. J., Desler, M. & Rizzino, A. Elevating SOX2 levels
1317 deleteriously affects the growth of medulloblastoma and glioblastoma cells. *PLoS*
1318 *One* **7**, e44087 (2012).
- 1319 51. Gangemi, R. M. R. *et al.* SOX2 silencing in glioblastoma tumor-initiating cells
1320 causes stop of proliferation and loss of tumorigenicity. *Stem Cells* **27**, 40–48
1321 (2009).
- 1322 52. Hägerstrand, D. *et al.* Identification of a SOX2-dependent subset of tumor- and
1323 sphere-forming glioblastoma cells with a distinct tyrosine kinase inhibitor sensitivity
1324 profile. *Neuro Oncol* **13**, 1178–1191 (2011).
- 1325 53. Sun, C. *et al.* Sox2 expression predicts poor survival of hepatocellular carcinoma
1326 patients and it promotes liver cancer cell invasion by activating Slug. *Med Oncol*
1327 **30**, 503 (2013).
- 1328 54. Chou, Y.-T. *et al.* The Emerging Role of SOX2 in Cell Proliferation and Survival and
1329 Its Crosstalk with Oncogenic Signaling in Lung Cancer. *STEM CELLS* **31**, 2607–
1330 2619 (2013).
- 1331 55. Sholl, L. M., Barletta, J. A., Yeap, B. Y., Chirieac, L. R. & Hornick, J. L. Sox2
1332 protein expression is an independent poor prognostic indicator in stage I lung
1333 adenocarcinoma. *Am J Surg Pathol* **34**, 1193–1198 (2010).

- 1334 56. Nakatsugawa, M. *et al.* SOX2 is overexpressed in stem-like cells of human lung
1335 adenocarcinoma and augments the tumorigenicity. *Lab Invest* **91**, 1796–1804
1336 (2011).
- 1337 57. Bass, A. J. *et al.* SOX2 is an amplified lineage-survival oncogene in lung and
1338 esophageal squamous cell carcinomas. *Nature Genetics* **41**, 1238–1242 (2009).
- 1339 58. Hussenet, T. *et al.* SOX2 Is an Oncogene Activated by Recurrent 3q26.3
1340 Amplifications in Human Lung Squamous Cell Carcinomas. *PLOS ONE* **5**, e8960
1341 (2010).
- 1342 59. Domenici, G. *et al.* A Sox2–Sox9 signalling axis maintains human breast luminal
1343 progenitor and breast cancer stem cells. *Oncogene* (2019) doi:10.1038/s41388-
1344 018-0656-7.
- 1345 60. Simões, B. M. *et al.* Effects of estrogen on the proportion of stem cells in the
1346 breast. *Breast Cancer Research and Treatment* **129**, 23–35 (2011).
- 1347 61. Bulstrode, H. *et al.* Elevated FOXG1 and SOX2 in glioblastoma enforces neural
1348 stem cell identity through transcriptional control of cell cycle and epigenetic
1349 regulators. *Genes Dev* **31**, 757–773 (2017).
- 1350 62. Jeon, H.-M. *et al.* ID4 imparts chemoresistance and cancer stemness to glioma
1351 cells by derepressing miR-9*-mediated suppression of SOX2. *Cancer Res* **71**,
1352 3410–3421 (2011).
- 1353 63. Zhang, L.-H. *et al.* TRIM24 promotes stemness and invasiveness of glioblastoma
1354 cells via activating Sox2 expression. *Neuro Oncol* **22**, (2020).

- 1355 64. Singh, S. *et al.* EGFR/Src/Akt signaling modulates Sox2 expression and self-
1356 renewal of stem-like side-population cells in non-small cell lung cancer. *Molecular*
1357 *Cancer* **11**, 73 (2012).
- 1358 65. Boumahdi, S. *et al.* SOX2 controls tumour initiation and cancer stem-cell functions
1359 in squamous-cell carcinoma. *Nature* **511**, 246–250 (2014).
- 1360 66. Berezovsky, A. D. *et al.* Sox2 promotes malignancy in glioblastoma by regulating
1361 plasticity and astrocytic differentiation. *Neoplasia* **16**, 193–206, 206.e19–25 (2014).
- 1362 67. Fang, X. *et al.* The SOX2 response program in glioblastoma multiforme: an
1363 integrated ChIP-seq, expression microarray, and microRNA analysis. *BMC*
1364 *Genomics* **12**, 11 (2011).
- 1365 68. Stolzenburg, S. *et al.* Targeted silencing of the oncogenic transcription factor SOX2
1366 in breast cancer. *Nucleic Acids Research* **40**, 6725–6740 (2012).
- 1367 69. Hoadley, K. A. *et al.* Multiplatform analysis of 12 cancer types reveals molecular
1368 classification within and across tissues of origin. *Cell* **158**, 929–944 (2014).
- 1369 70. Krasnov, G. S. *et al.* Pan-Cancer Analysis of TCGA Data Revealed Promising
1370 Reference Genes for qPCR Normalization. *Frontiers in Genetics* **10**, (2019).
- 1371 71. Maier, S. *et al.* SOX2 amplification is a common event in squamous cell
1372 carcinomas of different organ sites. *Human Pathology* **42**, 1078–1088 (2011).
- 1373 72. Liu, Y. *et al.* A predominant enhancer co-amplified with the SOX2 oncogene is
1374 necessary and sufficient for its expression in squamous cancer. *Nat Commun* **12**,
1375 7139 (2021).

- 1376 73. Buenrostro, J. D., Giresi, P. G., Zaba, L. C., Chang, H. Y. & Greenleaf, W. J.
1377 Transposition of native chromatin for multimodal regulatory analysis and personal
1378 epigenomics. *Nat Methods* **10**, 1213–1218 (2013).
- 1379 74. Boyle, A. P. *et al.* High-Resolution Mapping and Characterization of Open
1380 Chromatin across the Genome. *Cell* **132**, 311–322 (2008).
- 1381 75. Heintzman, N. D. *et al.* Distinct and predictive chromatin signatures of
1382 transcriptional promoters and enhancers in the human genome. *Nat. Genet.* **39**,
1383 311–318 (2007).
- 1384 76. Heintzman, N. D. *et al.* Histone modifications at human enhancers reflect global
1385 cell-type-specific gene expression. *Nature* **459**, 108–112 (2009).
- 1386 77. Soule, H. D., Vazquez, J., Long, A., Albert, S. & Brennan, M. A human cell line
1387 from a pleural effusion derived from a breast carcinoma. *Journal of the National*
1388 *Cancer Institute* **51**, 1409–1416 (1973).
- 1389 78. Liang, S. *et al.* Isolation and characterization of human breast cancer cells with
1390 SOX2 promoter activity. *Biochemical and Biophysical Research Communications*
1391 **437**, 205–211 (2013).
- 1392 79. Ling, G.-Q., Chen, Dong-bo, Wang, Bao-Qing & Zhang, Lan-Sheng. Expression of
1393 the pluripotency markers Oct3/4, Nanog and Sox2 in human breast cancer cell
1394 lines. *Oncology Letters* (2012) doi:10.3892/ol.2012.916.
- 1395 80. Roadmap Epigenomics Consortium *et al.* Integrative analysis of 111 reference
1396 human epigenomes. *Nature* **518**, 317–330 (2015).
- 1397 81. ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the
1398 human genome. *Nature* **489**, 57–74 (2012).

- 1399 82. Chen, C., Morris, Q. & Mitchell, J. A. Enhancer identification in mouse embryonic
1400 stem cells using integrative modeling of chromatin and genomic features. *BMC*
1401 *Genomics* **13**, 152 (2012).
- 1402 83. Mitchell, J. A. *et al.* Nuclear RNA sequencing of the mouse erythroid cell
1403 transcriptome. *PLoS ONE* **7**, e49274 (2012).
- 1404 84. Singh, G. *et al.* A flexible repertoire of transcription factor binding sites and a
1405 diversity threshold determines enhancer activity in embryonic stem cells. *Genome*
1406 *Res.* **31**, 564–575 (2021).
- 1407 85. Tolhuis, B., Palstra, R.-J., Splinter, E., Grosveld, F. & de Laat, W. Looping and
1408 Interaction between Hypersensitive Sites in the Active β -globin Locus. *Molecular*
1409 *Cell* **10**, 1453–1465 (2002).
- 1410 86. Carter, D., Chakalova, L., Osborne, C. S., Dai, Y. & Fraser, P. Long-range
1411 chromatin regulatory interactions in vivo. *Nature Genetics* **32**, 623 (2002).
- 1412 87. Fullwood, M. J. *et al.* An oestrogen-receptor-alpha-bound human chromatin
1413 interactome. *Nature* **462**, 58–64 (2009).
- 1414 88. Chan, H. L. *et al.* Polycomb complexes associate with enhancers and promote
1415 oncogenic transcriptional programs in cancer through multiple mechanisms. *Nat*
1416 *Commun* **9**, 3377 (2018).
- 1417 89. Cocce, K. J. *et al.* The Lineage Determining Factor GRHL2 Collaborates with
1418 FOXA1 to Establish a Targetable Pathway in Endocrine Therapy-Resistant Breast
1419 Cancer. *Cell Rep* **29**, 889-903.e10 (2019).

- 1420 90. Sato, T. *et al.* Epigenomic Profiling Discovers Trans-lineage SOX2 Partnerships
1421 Driving Tumor Heterogeneity in Lung Squamous Cell Carcinoma. *Cancer Research*
1422 **79**, 6084–6100 (2019).
- 1423 91. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and
1424 dispersion for RNA-seq data with DESeq2. *Genome Biol* **15**, (2014).
- 1425 92. Dunn, O. J. Multiple Comparisons Using Rank Sums. *Technometrics* **6**, 241–252
1426 (1964).
- 1427 93. Holm, S. A Simple Sequentially Rejective Multiple Test Procedure. *Scandinavian*
1428 *Journal of Statistics* **6**, 65–70 (1979).
- 1429 94. Kent, W. J. *et al.* The Human Genome Browser at UCSC. *Genome Res.* **12**, 996–
1430 1006 (2002).
- 1431 95. Moorthy, S. D. *et al.* Enhancers and super-enhancers have an equivalent
1432 regulatory role in embryonic stem cells through regulation of single or multiple
1433 genes. *Genome Res.* **27**, 246–258 (2017).
- 1434 96. Tobias, I. C. *et al.* Transcriptional enhancers: from prediction to functional
1435 assessment on a genome-wide scale. *Genome* **64**, 426–448 (2021).
- 1436 97. Dunnett, C. W. A Multiple Comparison Procedure for Comparing Several
1437 Treatments with a Control. *Journal of the American Statistical Association* **50**,
1438 1096–1121 (1955).
- 1439 98. Tukey, J. W. Comparing Individual Means in the Analysis of Variance. *Biometrics* **5**,
1440 99 (1949).
- 1441 99. Tompkins, D. H. *et al.* Sox2 Activates Cell Proliferation and Differentiation in the
1442 Respiratory Epithelium. *Am J Respir Cell Mol Biol* **45**, 101–110 (2011).

- 1443 100. Soufi, A., Donahue, G. & Zaret, K. S. Facilitators and Impediments of the
1444 Pluripotency Reprogramming Factors' Initial Engagement with the Genome. *Cell*
1445 **151**, 994–1004 (2012).
- 1446 101. Soufi, A. *et al.* Pioneer transcription factors target partial DNA motifs on
1447 nucleosomes to initiate reprogramming. *Cell* **161**, 555–568 (2015).
- 1448 102. Bentsen, M. *et al.* ATAC-seq footprinting unravels kinetics of transcription factor
1449 binding during zygotic genome activation. *Nat Commun* **11**, 4267 (2020).
- 1450 103. Castro-Mondragon, J. A. *et al.* JASPAR 2022: the 9th release of the open-access
1451 database of transcription factor binding profiles. *Nucleic Acids Research* **50**, D165–
1452 D173 (2022).
- 1453 104. Jeselsohn, R. *et al.* Embryonic transcription factor SOX9 drives breast cancer
1454 endocrine resistance. *Proceedings of the National Academy of Sciences* **114**,
1455 E4482–E4491 (2017).
- 1456 105. Yu, G., Wang, L.-G., Han, Y. & He, Q.-Y. clusterProfiler: an R Package for
1457 Comparing Biological Themes Among Gene Clusters. *OMICS: A Journal of*
1458 *Integrative Biology* **16**, 284–287 (2012).
- 1459 106. Stark R, Brown G. DiffBind: differential binding analysis of ChIP-Seq peak data.
1460 [http://bioconductor.org/packages/release/bioc/vignettes/DiffBind/inst/doc/DiffBind.p](http://bioconductor.org/packages/release/bioc/vignettes/DiffBind/inst/doc/DiffBind.pdf)
1461 [df](http://bioconductor.org/packages/release/bioc/vignettes/DiffBind/inst/doc/DiffBind.pdf) (2011).
- 1462 107. Corces, M. R. *et al.* The chromatin accessibility landscape of primary human
1463 cancers. *Science* **362**, eaav1898 (2018).

- 1464 108. Miyagi, S. *et al.* The Sox2 Regulatory Region 2 Functions as a Neural Stem Cell-
1465 specific Enhancer in the Telencephalon. *Journal of Biological Chemistry* **281**,
1466 13374–13381 (2006).
- 1467 109. Lambert, S. A. *et al.* The Human Transcription Factors. *Cell* **172**, 650–665 (2018).
- 1468 110. Meuleman, W. *et al.* Index and biological spectrum of human DNase I
1469 hypersensitive sites. *Nature* **584**, 244–251 (2020).
- 1470 111. Pennacchio, L. A. *et al.* In vivo enhancer analysis of human conserved non-coding
1471 sequences. *Nature* **444**, 499–502 (2006).
- 1472 112. Woolfe, A. *et al.* Highly conserved non-coding sequences are associated with
1473 vertebrate development. *PLoS Biol.* **3**, e7 (2005).
- 1474 113. Pollard, K. S., Hubisz, M. J., Rosenbloom, K. R. & Siepel, A. Detection of
1475 nonneutral substitution rates on mammalian phylogenies. *Genome Res* **20**, 110–
1476 121 (2010).
- 1477 114. Liu, C. *et al.* An ATAC-seq atlas of chromatin accessibility in mouse tissues. *Sci*
1478 *Data* **6**, 65 (2019).
- 1479 115. Bonnet, D. & Dick, J. E. Human acute myeloid leukemia is organized as a hierarchy
1480 that originates from a primitive hematopoietic cell. *Nat Med* **3**, 730–737 (1997).
- 1481 116. Chaffer, C. L. *et al.* Normal and neoplastic nonstem cells can spontaneously
1482 convert to a stem-like state. *PNAS* **108**, 7950–7955 (2011).
- 1483 117. Lapidot, T. *et al.* A cell initiating human acute myeloid leukaemia after
1484 transplantation into SCID mice. *Nature* **367**, 645–648 (1994).
- 1485 118. Gupta, P. B. *et al.* Stochastic state transitions give rise to phenotypic equilibrium in
1486 populations of cancer cells. *Cell* **146**, 633–644 (2011).

- 1487 119. Thurman, R. E. *et al.* The accessible chromatin landscape of the human genome.
1488 *Nature* **489**, 75–82 (2012).
- 1489 120. Bernardo, G. M. *et al.* FOXA1 is an essential determinant of ERalpha expression
1490 and mammary ductal morphogenesis. *Development* **137**, 2045–2054 (2010).
- 1491 121. Liu, Y. *et al.* Foxa1 is essential for mammary duct formation. *Genesis* **54**, 277–285
1492 (2016).
- 1493 122. Besnard, V., Wert, S. E., Kaestner, K. H. & Whitsett, J. A. Stage-specific regulation
1494 of respiratory epithelial cell differentiation by Foxa1. *Am J Physiol Lung Cell Mol*
1495 *Physiol* **289**, L750-759 (2005).
- 1496 123. Paranjapye, A., Mutolo, M. J., Ebron, J. S., Leir, S.-H. & Harris, A. The FOXA1
1497 transcriptional network coordinates key functions of primary human airway
1498 epithelial cells. *Am J Physiol Lung Cell Mol Physiol* **319**, L126–L136 (2020).
- 1499 124. Camolotto, S. A. *et al.* FoxA1 and FoxA2 drive gastric differentiation and suppress
1500 squamous identity in NKX2-1-negative lung cancer. *Elife* **7**, e38579 (2018).
- 1501 125. Fu, X. *et al.* FOXA1 overexpression mediates endocrine resistance by altering the
1502 ER transcriptome and IL-8 expression in ER-positive breast cancer. *PNAS* **113**,
1503 E6600–E6609 (2016).
- 1504 126. Orstad, G. *et al.* FoxA1 and FoxA2 control growth and cellular identity in NKX2-1-
1505 positive lung adenocarcinoma. *Dev Cell* **57**, 1866-1882.e10 (2022).
- 1506 127. Liu, K.-C. *et al.* The multiple roles for Sox2 in stem cell maintenance and
1507 tumorigenesis. *Cell Signal* **25**, (2013).

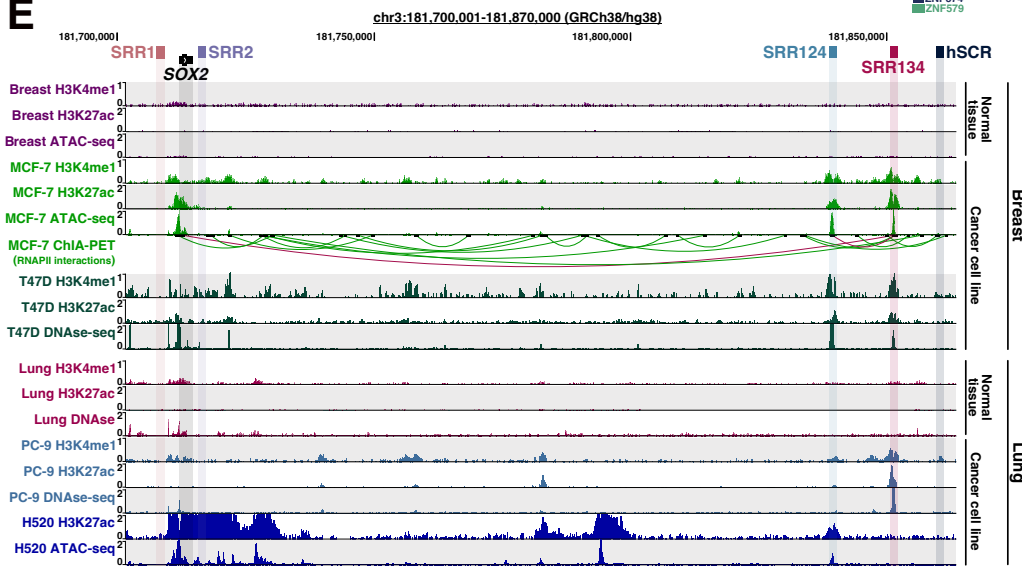
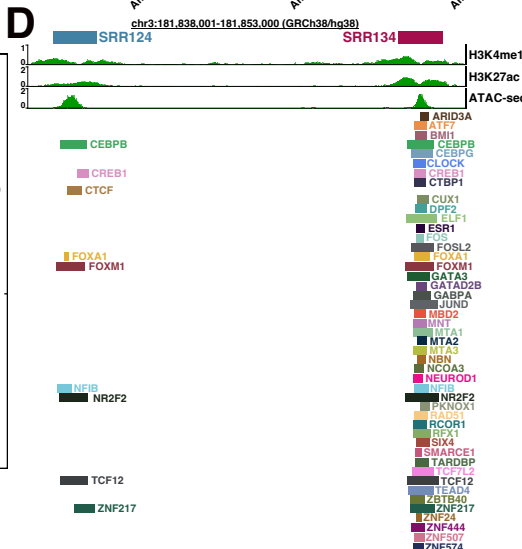
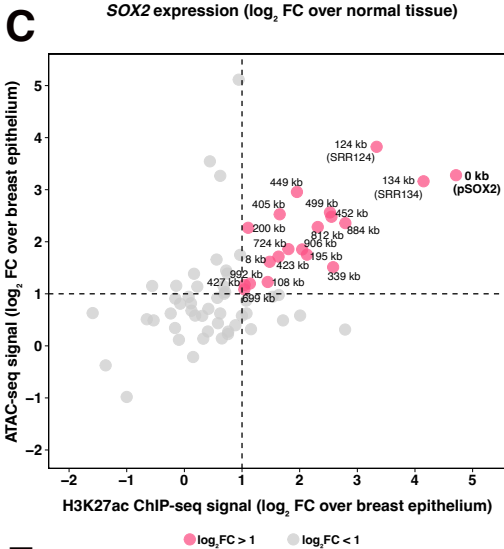
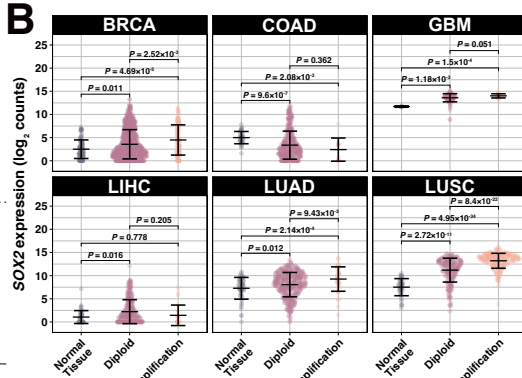
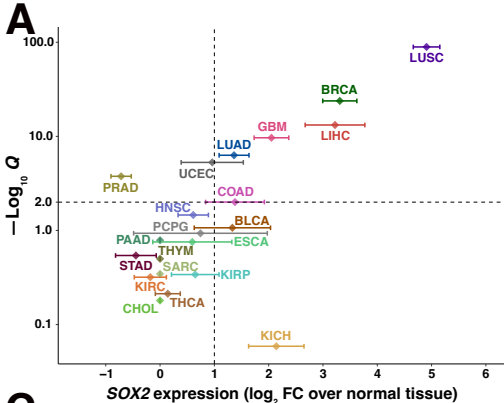
- 1508 128. Pierrou, S., Hellqvist, M., Samuelsson, L., Enerbäck, S. & Carlsson, P. Cloning and
1509 characterization of seven human forkhead proteins: binding site specificity and
1510 DNA bending. *The EMBO Journal* **13**, 5002–5012 (1994).
- 1511 129. Li, L., Wu, D., Yu, Q., Li, L. & Wu, P. Prognostic value of FOXM1 in solid tumors: a
1512 systematic review and meta-analysis. *Oncotarget* **8**, 32298–32308 (2017).
- 1513 130. Harris, L., Genovesi, L. A., Gronostajski, R. M., Wainwright, B. J. & Piper, M.
1514 Nuclear factor one transcription factors: Divergent functions in developmental
1515 versus adult stem cell populations. *Dev Dyn* **244**, 227–238 (2015).
- 1516 131. Steele-Perkins, G. *et al.* The transcription factor gene Nfib is essential for both lung
1517 maturation and brain development. *Mol Cell Biol* **25**, 685–698 (2005).
- 1518 132. Gründer, A. *et al.* Nuclear factor I-B (Nfib) deficient mice have severe lung
1519 hypoplasia. *Mech Dev* **112**, 69–77 (2002).
- 1520 133. Hsu, Y.-C. *et al.* Mesenchymal nuclear factor I B regulates cell proliferation and
1521 epithelial differentiation during lung maturation. *Dev Biol* **354**, 242–252 (2011).
- 1522 134. Favaro, R. *et al.* Hippocampal development and neural stem cell maintenance
1523 require Sox2-dependent regulation of Shh. *Nature Neuroscience* **12**, 1248–1256
1524 (2009).
- 1525 135. Becker-Santos, D. D., Lonergan, K. M., Gronostajski, R. M. & Lam, W. L. Nuclear
1526 Factor I/B: A Master Regulator of Cell Differentiation with Paradoxical Roles in
1527 Cancer. *EBioMedicine* **22**, 2–9 (2017).
- 1528 136. Zhou, M. *et al.* miR-365 promotes cutaneous squamous cell carcinoma (CSCC)
1529 through targeting nuclear factor I/B (NFIB). *PLoS One* **9**, e100620 (2014).

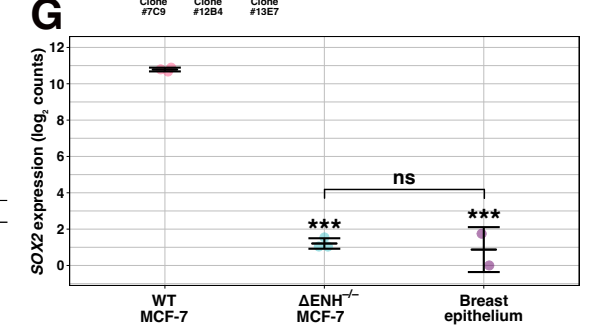
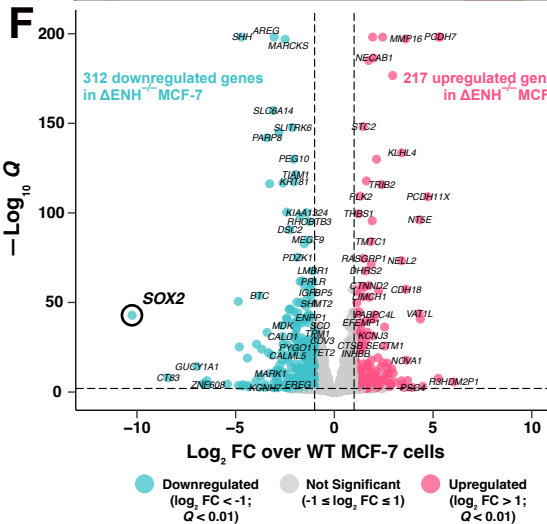
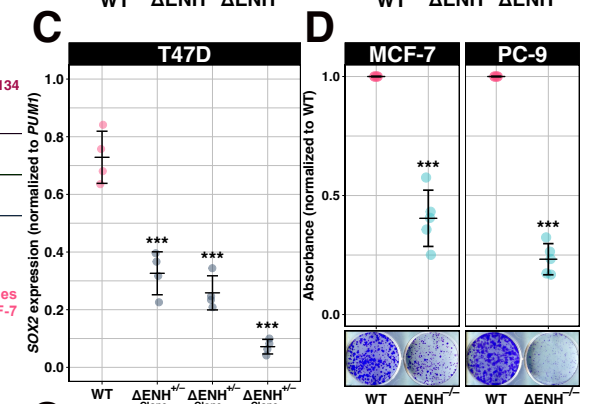
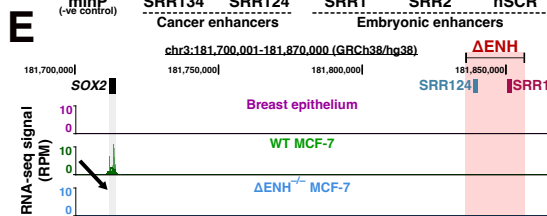
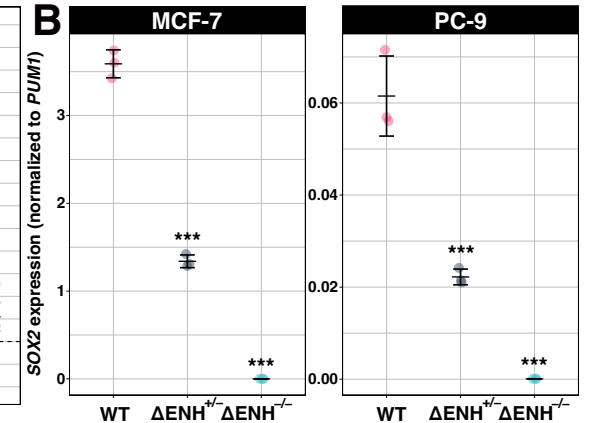
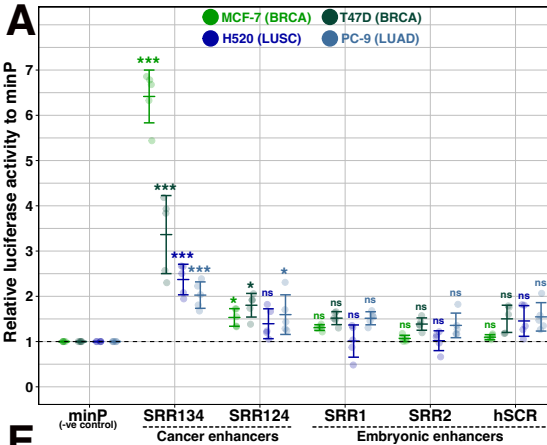
- 1530 137. Becker-Santos, D. D. *et al.* Developmental transcription factor NFIB is a putative
1531 target of oncofetal miRNAs and is associated with tumour aggressiveness in lung
1532 adenocarcinoma. *J Pathol* **240**, 161–172 (2016).
- 1533 138. Ferone, G. *et al.* SOX2 Is the Determining Oncogenic Switch in Promoting Lung
1534 Squamous Cell Carcinoma from Different Cells of Origin. *Cancer Cell* **30**, 519–532
1535 (2016).
- 1536 139. Wu, C. *et al.* NFIB promotes cell growth, aggressiveness, metastasis and EMT of
1537 gastric cancer through the Akt/Stat3 signaling pathway. *Oncol Rep* **40**, 1565–1573
1538 (2018).
- 1539 140. Otsubo, T., Akiyama, Y., Yanagihara, K. & Yuasa, Y. SOX2 is frequently
1540 downregulated in gastric cancers and inhibits cell growth through cell-cycle arrest
1541 and apoptosis. *Br J Cancer* **98**, 824–831 (2008).
- 1542 141. Wang, S. *et al.* SOX2, a predictor of survival in gastric cancer, inhibits cell
1543 proliferation and metastasis by regulating PTEN. *Cancer Lett* **358**, 210–219 (2015).
- 1544 142. Zhang, X. *et al.* SOX2 in gastric carcinoma, but not Hath1, is related to patients'
1545 clinicopathological features and prognosis. *J Gastrointest Surg* **14**, 1220–1226
1546 (2010).
- 1547 143. Miyagi, S. *et al.* The Sox-2 regulatory regions display their activities in two distinct
1548 types of multipotent stem cells. *Mol. Cell. Biol.* **24**, 4207–4220 (2004).
- 1549 144. Aran, D. *et al.* Embryonic Stem Cell (ES)-Specific Enhancers Specify the
1550 Expression Potential of ES Genes in Cancer. *PLoS Genet* **12**, e1005840 (2016).

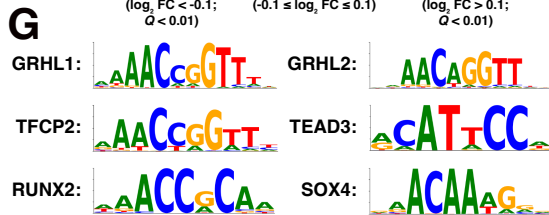
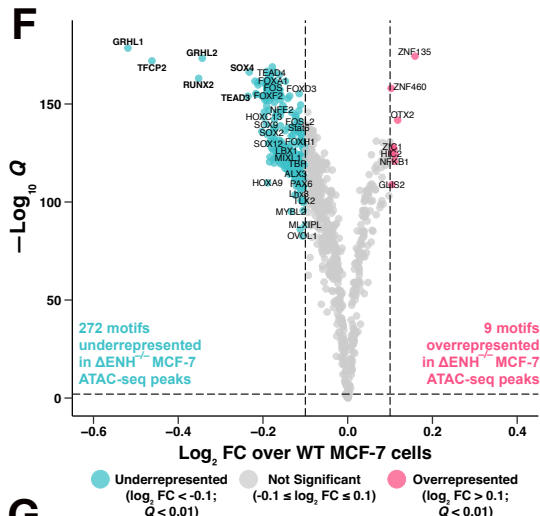
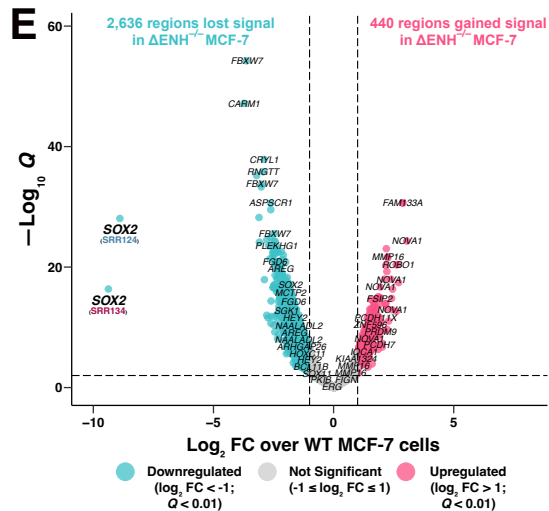
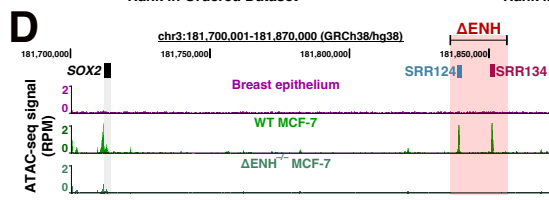
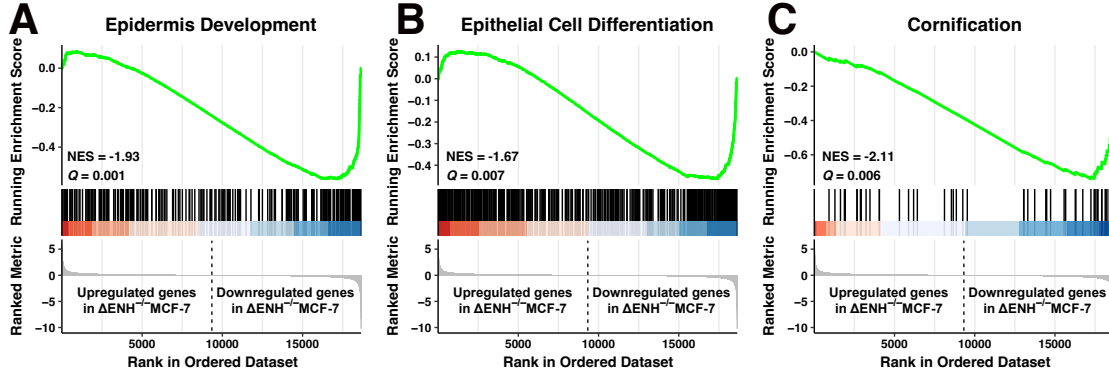
- 1551 145. Iglesias, J. M. *et al.* The Activation of the Sox2 RR2 Pluripotency Transcriptional
1552 Reporter in Human Breast Cancer Cell Lines is Dynamic and Labels Cells with
1553 Higher Tumorigenic Potential. *Front Oncol* **4**, 308 (2014).
- 1554 146. Jung, K. *et al.* Triple negative breast cancers comprise a highly tumorigenic cell
1555 subpopulation detectable by its high responsiveness to a Sox2 regulatory region 2
1556 (SRR2) reporter. *Oncotarget* **6**, 10366–10373 (2015).
- 1557 147. Saenz-Antoñanzas, A. *et al.* CRISPR/Cas9 Deletion of SOX2 Regulatory Region 2
1558 (SRR2) Decreases SOX2 Malignant Activity in Glioblastoma. *Cancers (Basel)* **13**,
1559 1574 (2021).
- 1560 148. Björkqvist, A. M. *et al.* DNA gains in 3q occur frequently in squamous cell
1561 carcinoma of the lung, but not in adenocarcinoma. *Genes Chromosomes Cancer*
1562 **22**, 79–82 (1998).
- 1563 149. Mali, P. *et al.* RNA-Guided Human Genome Engineering via Cas9. *Science* **339**,
1564 823–826 (2013).
- 1565 150. Ding, Q. *et al.* Enhanced efficiency of human pluripotent stem cell genome editing
1566 through replacing TALENs with CRISPRs. *Cell Stem Cell* **12**, 393–394 (2013).
- 1567 151. Ran, F. A. *et al.* Genome engineering using the CRISPR-Cas9 system. *Nat Protoc*
1568 **8**, 2281–2308 (2013).
- 1569 152. Ahier, A. & Jarriault, S. Simultaneous Expression of Multiple Proteins Under a
1570 Single Promoter in *Caenorhabditis elegans* via a Versatile 2A-Based Toolkit.
1571 *Genetics* **196**, 605–613 (2014).

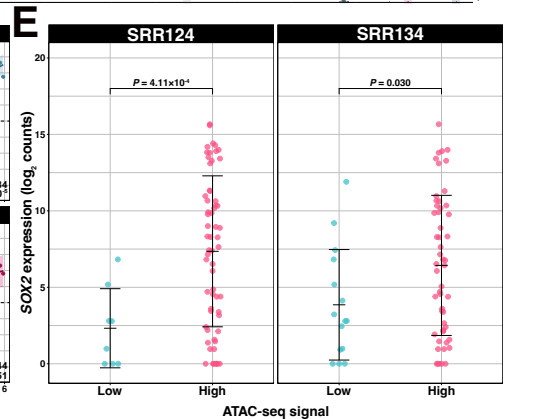
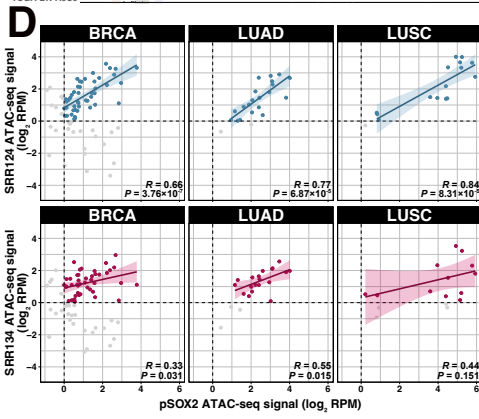
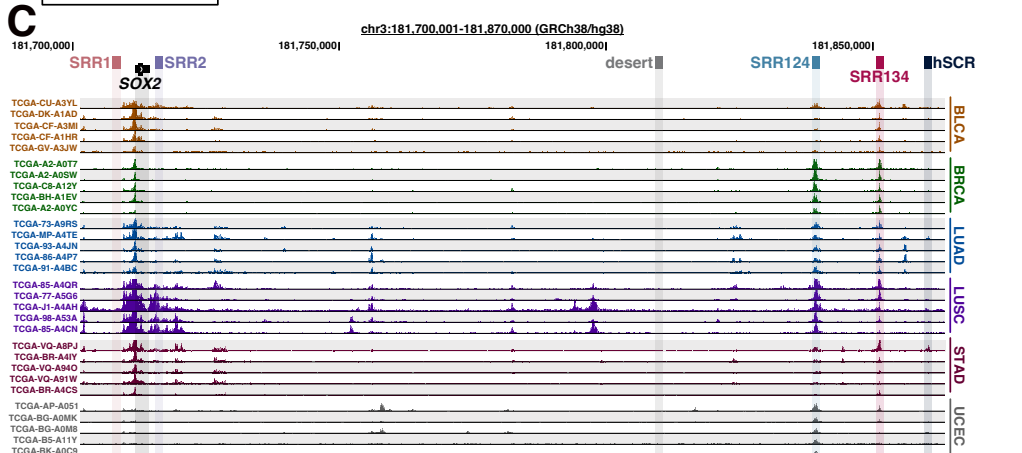
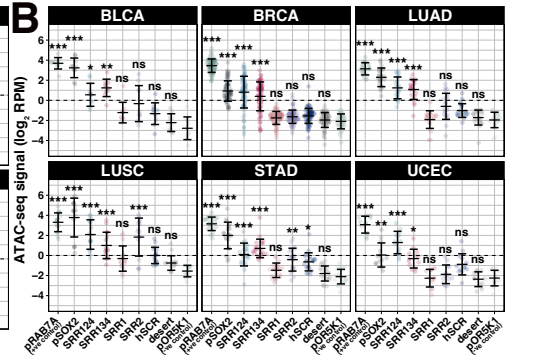
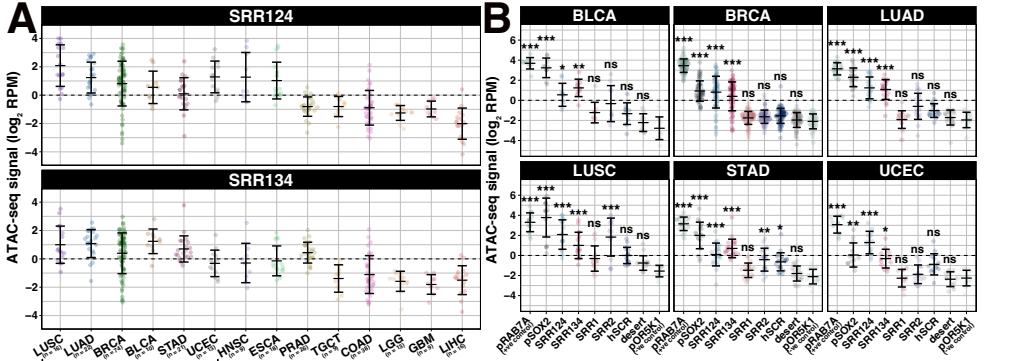
- 1572 153. Zhang, J.-P. *et al.* Efficient precise knockin with a double cut HDR donor after
1573 CRISPR/Cas9-mediated double-stranded DNA cleavage. *Genome Biology* **18**, 35
1574 (2017).
- 1575 154. Kılıç, Y., Çelebiler, A. Ç. & Sakızlı, M. Selecting housekeeping genes as references
1576 for the normalization of quantitative PCR data in breast cancer. *Clinical and*
1577 *Translational Oncology* **16**, 184–190 (2014).
- 1578 155. Lyng, M. B., Lænkholm, A.-V., Pallisgaard, N. & Ditzel, H. J. Identification of genes
1579 for normalization of real-time RT-PCR data in breast carcinomas. *BMC Cancer* **8**,
1580 (2008).
- 1581 156. Chen, S., Zhou, Y., Chen, Y. & Gu, J. fastp: an ultra-fast all-in-one FASTQ
1582 preprocessor. *Bioinformatics* **34**, i884–i890 (2018).
- 1583 157. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15–
1584 21 (2013).
- 1585 158. Liao, Y., Smyth, G. K. & Shi, W. featureCounts: an efficient general purpose
1586 program for assigning sequence reads to genomic features. *Bioinformatics* **30**,
1587 923–930 (2014).
- 1588 159. Shen, L., Shao, N., Liu, X. & Nestler, E. ngs.plot: Quick mining and visualization of
1589 next-generation sequencing data by integrating genomic databases. *BMC*
1590 *Genomics* **15**, 284 (2014).
- 1591 160. Colaprico, A. *et al.* TCGAAbiolinks: an R/Bioconductor package for integrative
1592 analysis of TCGA data. *Nucleic Acids Research* **44**, e71–e71 (2016).
- 1593 161. Liu, J. *et al.* An Integrated TCGA Pan-Cancer Clinical Data Resource to Drive High-
1594 Quality Survival Outcome Analytics. *Cell* **173**, 400-416.e11 (2018).

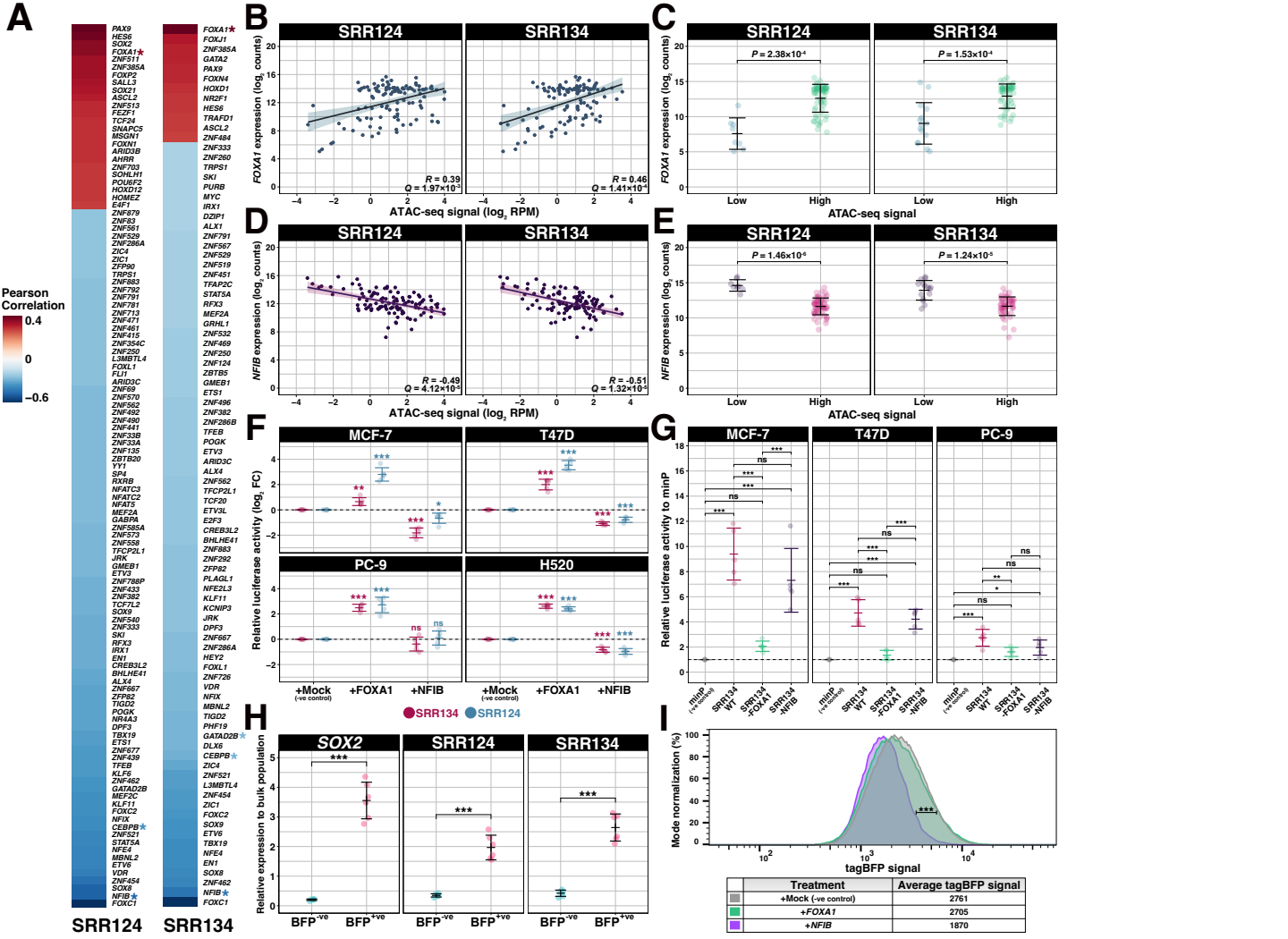
- 1595 162. Corces, M. R. *et al.* An improved ATAC-seq protocol reduces background and
1596 enables interrogation of frozen tissues. *Nature Methods* **14**, 959–962 (2017).
- 1597 163. Heinz, S. *et al.* Simple combinations of lineage-determining transcription factors
1598 prime cis-regulatory elements required for macrophage and B cell identities. *Mol.*
1599 *Cell* **38**, 576–589 (2010).
- 1600 164. Zhu, L. J. *et al.* ChIPpeakAnno: a Bioconductor package to annotate ChIP-seq and
1601 ChIP-chip data. *BMC Bioinformatics* **11**, (2010).
- 1602 165. Seibt, K. M., Schmidt, T. & Heitkam, T. FlexiDot: highly customizable, ambiguity-
1603 aware dotplots for visual sequence analyses. *Bioinformatics* **34**, 3575–3577 (2018).
- 1604 166. Gertsenstein, M. & Nutter, L. M. J. Production of knockout mouse lines with Cas9.
1605 *Methods* **191**, 32–43 (2021).
- 1606 167. Kaplan, E. L. & Meier, P. Nonparametric Estimation from Incomplete Observations.
1607 *Journal of the American Statistical Association* **53**, 457–481 (1958).
- 1608 168. Sievers, F. *et al.* Fast, scalable generation of high-quality protein multiple sequence
1609 alignments using Clustal Omega. *Mol Syst Biol* **7**, 539 (2011).
- 1610





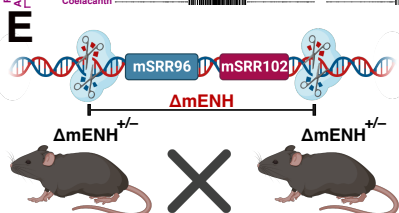
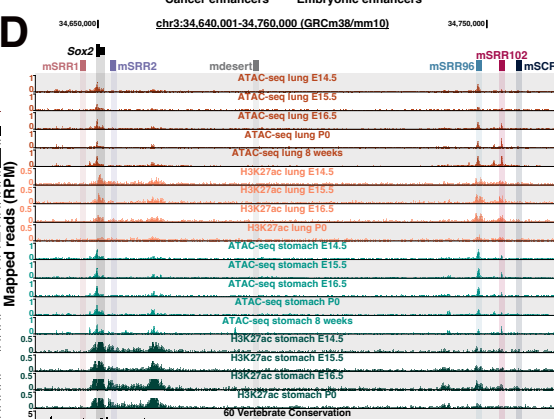
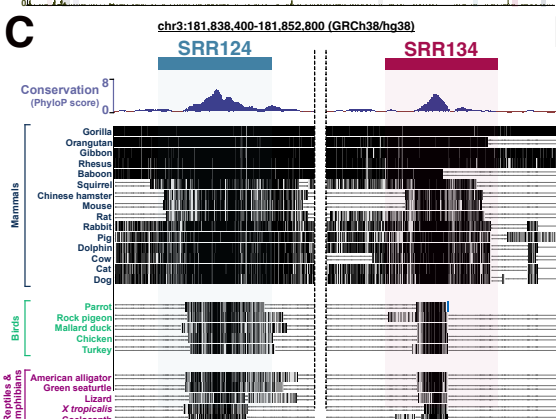
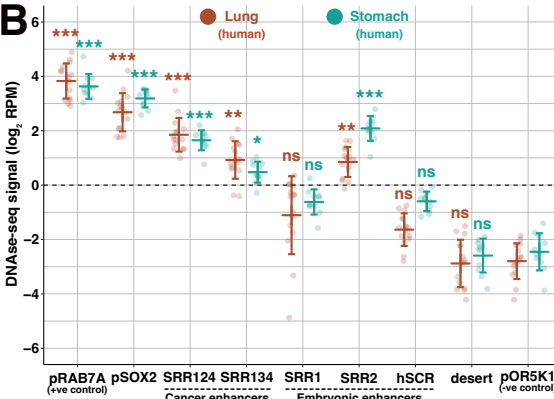
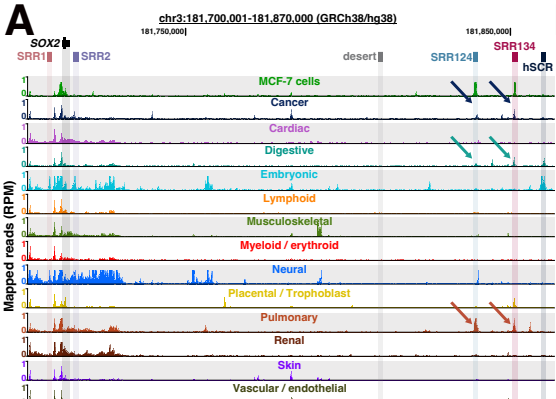






Pearson Correlation
0.4
0
-0.6

SRR124 SRR134



F

		Genotype			Total	P
		WT	ΔmENH ^{+/-}	ΔmENH ^{-/-}		
Expected	n. of live animals	20	40	20	80	
	% of total	25.00	50.00	25.00	100.00	
Observed	n. of live animals	20	26	0	46	3.92 × 10 ⁻⁶
	% of total	43.48	56.52	0.00	100.00	

