

1 **Sox9 regulates cell state and activity of embryonic mouse mammary**  
2 **progenitor cells**

3

4 **Authors:**

5 Naoko Kogata<sup>1,2</sup>, Philip Bland<sup>1</sup>, Mandy Tsang<sup>1</sup>, Erik Oliemuller<sup>1</sup>, Anne Lowe<sup>1</sup>,  
6 and Beatrice A. Howard<sup>1\*</sup>

7

8 **Affiliations:**

9 <sup>1</sup>The Breast Cancer Now Toby Robins Research Centre, Division of Breast  
10 Cancer Research, The Institute of Cancer Research, London, UK

11 <sup>2</sup>current address: Cellular Signalling and Cytoskeletal Function Lab, The  
12 Francis Crick Institute, London, UK

13

14

15 \*Author for correspondence:

16 Beatrice A. Howard

17 e-mail: [beatrice.howard@icr.ac.uk](mailto:beatrice.howard@icr.ac.uk)

18 phone: +44-203-437-7359

19

20 **Keywords:**

21 embryonic mammary gland, embryonic mammary progenitor cell, mammary  
22 lineage, multipotency, Sox9

23

24

25

26

27 **Abstract**

28 Embryonic mammary cells are a unique population comprised of  
29 undifferentiated, highly plastic progenitor cells that create normal mammary  
30 tissues. The mammary gland continues to develop after birth from  
31 descendants of embryonic mammary cells. Here, we establish novel cell lines  
32 from mouse mammary organs, immediately after they formed during prenatal  
33 development, to facilitate studies of primitive mammary cells, which are  
34 difficult to isolate in sufficient quantities for use in functional experiments. We  
35 show some lines can be induced to secrete milk, a distinguishing feature of  
36 mammary epithelial cells. Targeted deletion of *Sox9*, from one clone, leads to  
37 decreased ability to respond to lactogenic stimuli, consistent with a previously  
38 identified role for *Sox9* in regulating luminal progenitor function. *Sox9* ablation  
39 also leads to alterations in 3D morphology and downregulation of *Zeb1*, a key  
40 epithelial-mesenchymal transition regulator. Prenatal mammary cell lines are  
41 an invaluable resource to study regulation of mammary progenitor cell biology  
42 and development.

43

44 **Introduction**

45 Embryonic breast epithelial cells are a unique cell population comprised of  
46 undifferentiated and highly plastic progenitor cells that ultimately give rise to  
47 all other postnatal breast epithelial cells. Lineage tracing studies have  
48 indicated that embryonic mammary cells are multipotent *in vivo*<sup>1-3</sup>. An  
49 important area of research in mammary gland biology is to determine the  
50 roles of genes and signalling pathways that regulate embryonic stages of  
51 mammary gland development, as many of these are also relevant to  
52 processes that are deregulated in cancer<sup>4,5</sup>. Despite their relevance to breast  
53 cancer research, the routine use of primary mid-gestation embryonic  
54 mammary cells for functional study is not currently feasible, due to the small  
55 size of the nascent organ.

56 In mice, mammary gland development commences at embryonic day  
57 11 (E11) with the sequential appearance of five pairs of mammary primordia<sup>6</sup>.  
58 Local epithelial thickenings invaginate to the underlying tissue to form buds,  
59 which from E12.5 onwards are surrounded by a specialized condensed  
60 mammary mesenchyme. Mammary buds grow relatively slowly in size until  
61 E14 when the epithelial cells within the bud start to proliferate extensively and  
62 then invade into the underlying mesenchymal tissues<sup>6</sup>. These early stages of  
63 development are of particular interest as the cells have a number of unique  
64 properties. The epithelial cells within the E11-E13 stage mammary organ are  
65 largely quiescent<sup>7,8</sup>. At these stages of development, epithelial cells are  
66 thought to accrue within the mammary organ via localised cell movements<sup>9,10</sup>.  
67 Dissociated embryonic mammary cells from E12-E13-stage organs have  
68 minimal ability to efficiently repopulate cleared mammary fat pads, whilst cells  
69 from E16-18-stage organs have a much higher ability to engraft<sup>11,12</sup>. Intact  
70 mammary bud epithelium from E13-stage embryos can repopulate cleared fat  
71 pads suggesting a stem cell population has been delimited by mid-gestation<sup>13</sup>.  
72 Recent results from lineage tracing experiments indicate that embryonic  
73 mammary cells at E12-E13 stages are multipotent at the cellular level and  
74 become lineage-restricted prior to birth<sup>1-3</sup>. Although embryonic mammary  
75 progenitor cells from E12-13-stages harbour very distinct biological properties  
76 from embryonic mammary progenitor cells isolated from E16-17-stages and

77 from the postnatal mammary gland, a lack of appropriate *in vitro* models has  
78 limited their accessibility for many researchers.

79 Most studies of the embryonic mammary gland have relied on analyses  
80 of embryos from genetically modified mice or embryonic mammary organ  
81 explant cultures<sup>14-17</sup>. These methods require considerable training, expertise  
82 and the use of animals. However, using these novel embryonic mammary cell  
83 lines and standard two- and three-dimensional culture techniques, we model  
84 several key aspects of embryonic mammary gland development *in vitro*. Using  
85 CRISPR-Cas9 genome editing, we investigate the role of Sox9, an embryonic  
86 Sox transcription factor, which has been implicated in conferring stem cell  
87 state to differentiated postnatal mammary epithelial cells<sup>18</sup>, in the regulation of  
88 stem cell activity and the differentiation potential of cells formed during early  
89 stages of embryonic mammary gland development. Our findings highlight the  
90 distinct biological features and context dependent regulation of embryonic  
91 mammary progenitor cells and are a novel resource for studying this unique  
92 cell population.

93

94

95

96

97

98

99 **Results**

100 **Novel embryonic mammary progenitor cell lines established**

101 To date, the establishment and maintenance of primary mammary embryonic  
102 epithelial cell cultures from mouse embryos at stages between E12-13 has  
103 not been possible. To overcome this issue, we have taken advantage of the  
104 Immortomouse, a genetically modified mouse, to introduce a temperature-  
105 sensitive Simian virus (SV) 40 antigen under control of an interferon regulated  
106 promoter that enables immortalisation of certain types of cells, including  
107 epithelial and mesenchymal cells<sup>19</sup>. Immortomice were bred with s-SHIP-GFP  
108 mice, which GFP expression marks epithelial progenitor cells, including those  
109 in the embryonic mammary primordia to facilitate and confirm dissection of the  
110 mammary organ when it is first morphologically distinct<sup>20</sup>. The majority of s-  
111 SHIP-GFP embryonic mammary epithelial cells are GFP+ and adjacent  
112 embryonic mammary mesenchymal cells are GFP- (Figure 1A and  
113 Supplementary Fig. 1A). Mammary primordia number three were  
114 microdissected from E12.0-stage embryos so that the mammary epithelium,  
115 mammary mesenchyme, and fat pad precursors were isolated and such that  
116 other GFP+ cell types marked in s-SHIP-GFP embryos<sup>21</sup> were excluded  
117 (Figure 1A). Excised mammary organs were embedded in thick Basement  
118 Membrane Extract (BME) for a month, so that embryonic mammary progenitor  
119 cells (eMPC) from the mammary organoid were able to proliferate (Figure 1B).  
120 This adaptation period appeared to be important for successful cell isolation  
121 since enzymatic digestion of Immorto:s-SHIP-GFP mammary organs  
122 immediately after microdissection did not give rise to cultivable eMPCs  
123 beyond one month. After proliferating for one month, eMPCs were harvested  
124 by enzymatically digesting the organoid culture grown in thick BME and  
125 expanding the cells in two-dimension culture. The bulk cells of this expanded  
126 eMPC are referred to subsequently as ePool (Figure 1B). To obtain clonal cell  
127 lines, single ePool cells were separated using fluorescence-activated cell  
128 sorter (FACS) according to s-SHIP-GFP expression status (GFP+ or GFP-) at  
129 the time of sorting (Figure 1B and Supplementary Fig. 1B). ~9% of GFP+ and  
130 ~2% of GFP- cells gave rise to viable clones. Sixteen clones were derived  
131 from the GFP+ fraction (designated eG1, eG2, eG2E9, etc). e1 and e2 were  
132 the only clones that survived from GFP- fraction (Supplementary Fig. 1C).

133 However, GFP expression is not necessarily indicative of tissue of origin since  
134 the activity of s-SHIP promoter is variable in cell populations expanded in 2D  
135 culture. eMPC clones expanded from single s-SHIP-GFP<sup>+</sup> or single s-SHIP-  
136 GFP<sup>-</sup> cells contain both GFP<sup>+</sup> and GFP<sup>-</sup> cells upon passage in 2D culture  
137 (Supplementary Fig. 1C).

138 Hierarchical cluster analysis of gene expression profiles obtained by  
139 RNA sequencing shows that eMPCs are distinct from other types of  
140 progenitor and stem cells, including embryonic stem cells (ESC), mouse  
141 embryonic fibroblasts (MEF), mesenchymal stem cells (MSC) and bone  
142 marrow mononuclear cells (BMMC) (Figure 2A). From 17204 significantly  
143 modulated genes, the 2059 most highly variable were selected when DESeq2  
144 software was employed with the Likelihood ratio test (LRT), and using the  
145 intensity difference test (detailed selection strategy and gene lists are in  
146 Supplementary Dataset 1). Principal component analysis (PCA), based on 92  
147 genes, produces a similar clustering of the eMPC clones with other stem cell  
148 types (Figure 2B and Supplementary Dataset 1).

149 Mammary organs from C57BL/6 E12.5-stage mouse embryos were  
150 microdissected so that mammary epithelium, mammary stroma (mammary  
151 mesenchyme, and adjacent fat pad precursor), and surface epithelium were  
152 separated and the tissues were used for gene expressing profiling using RNA-  
153 sequencing. Hierarchical clustering based on transcriptome profiles of the  
154 separated embryonic mammary primordia tissues and eMPCs and other stem  
155 cell types was performed (Supplementary Fig. 2 and Supplementary Dataset  
156 2). To assess the global landscape of expression in this dataset, we  
157 performed PCA and 26 genes revealed clustering of eMPC lines into three  
158 distinct groups based on their expression of genes encoding Keratin and cell  
159 adhesion regulators (Figure 2C). Our previous results have shown that  
160 embryonic mammary cells are largely devoid of cells expressing differentiation  
161 markers, although a few cells present within the E12-stage embryonic  
162 mammary bud epithelium express stem cell/progenitor markers associated  
163 with mature postnatal mammary epithelial cell (MEC) lineages<sup>22</sup>. PCA  
164 identifies clusters based on expression of several genes, some of which  
165 control cell adhesion, including the protocadherin, *Pcdh7*, integrin, *Itga3*, and  
166 extracellular matrix (ECM) ligand, *Lama4*, that are highly expressed by a

167 subset of the cell lines (eG1, eG2C7, eG1C10), which also express Keratins  
168 (*Krt7*, *Krt18*, *Krt19*). These clones belonging to the cluster marked in green in  
169 Figure 2 may represent mammary epithelial progenitor cells. Clones  
170 belonging to another cluster, marked in yellow in Figure 2, including e1, which  
171 was selected further study, express high levels of a number of markers  
172 including, *Cd34*, *Cd74*, *Ebf2*, *Fabp4*, *Runx2*, and *Sox9* that are associated  
173 with a variety of stem cell types<sup>23-28</sup> and may represent less-differentiated  
174 progenitor/stem cells types. Two cells lines with less distinct features, (marked  
175 in orange in Figure 2), e2 and eG2 that cluster with the pool, were also  
176 selected for further study. eG1 was selected for further study since it  
177 expresses high levels of markers associated with an epithelial-state, including  
178 Keratins, the epithelial basement membrane protein *Lamc2*, and *Lgr4*, a  
179 regulator of mammary gland development and stem cell activity<sup>29</sup>.

180 In total, two clones expanded from GFP+ cells (eG1 and eG2) and the  
181 two clones expanded from GFP- cells (e1 and e2) lines, along with the ePool,  
182 were used for further characterization based on their distinct morphologies,  
183 transcriptome profiles and cluster analysis results (Figure 3). The two GFP-  
184 cell lines (e1, e2) do not cluster together. However, single GFP+ and GFP-  
185 cells were expanded prior to other characterization including, RNA  
186 sequencing, and these cell populations contained a mixture of both GFP+ and  
187 GFP- cells (Supplementary Fig. 1C).

188 We selected a subset of markers that have direct links to stem cell  
189 biology and confirmed the RNA sequencing profiles using qRT-PCR (Figure  
190 3B). We also used immunohistochemistry (IHC) to confirm expression of Sox9  
191 and Twist1 protein in specific eMPC clones (Figure 3C). Sox9 and Twist1  
192 protein expression were discordant with the RNA expression results, but  
193 mismatch between transcriptional and translational levels is a frequently  
194 occurring phenomena<sup>30</sup>.

195

### 196 **eMPCs harbour varying capacities for differentiation into the** 197 **mesodermal lineage**

198 We examined the capacity of eG1, eG2, e1, e2 and the ePool to differentiate  
199 into mesodermal lineages. Cells were exposed to differentiation media,

200 previously reported to induce MSCs toward adipocyte or endothelial cell  
201 lineages. Under adipogenic conditions, only e1 cells produced lipid-filled  
202 adipocytes in a substantial fraction of cells, indicating notable adipogenic  
203 ability, whereas eG1 and eG2 do not respond to adipogenic stimuli. (Figure  
204 4A). Under endotheliogenic conditions, e1, eG2 and the ePool formed tubular  
205 network comparable to control Human umbilical vein endothelial cells  
206 (HUVECs) cultured in endothelium growth medium (Figure 4B). Both e1 and  
207 eG2 express a vascular progenitor cell marker *Cd34* prior to induction (Figure  
208 3A). Overall, these functional studies suggest that clone e1 is the most  
209 responsive of the four cell lines to mesodermal lineage differentiation cues.

210

### 211 **eMPCs are clonogenic and form distinct sphere morphologies**

212 eG1, eG2, e1, e2, and the ePool were evaluated for mammary sphere-  
213 forming ability using the standard sphere formation technique. Each clone  
214 was cultured under sphere-forming conditions and formed spheres with highly  
215 distinct morphologies (Figure 5A). All eMPCs formed spheres, although e1  
216 had the highest sphere-forming rate (3.93%) and eG1 the lowest (0.69%)  
217 (Figure 5B). Two clones (e1, eG1) were selected for more detailed functional  
218 assessments based on their distinct responses in functional assays and  
219 marker profiles. eG1 spheres show higher form factor, a measure of  
220 sphericity, when compared to those of e1, suggesting that eG1 spheres are  
221 relatively compact compared to the spheres formed from the other clones  
222 (Figure 5C-D). e1 cells had limited engraftment ability when injected into  
223 cleared mammary fat pads (Supplementary Figure 3). No teratoma formation  
224 was observed, as would be expected to occur after xenografting of ESC cells.

225

### 226 **eMPCs form differentiated mammary acini *in vitro***

227 eMPCs were evaluated for their ability to undergo functional alveolar  
228 differentiation by assessing cellular production of milk proteins by  
229 immunofluorescence staining with an antibody that detects milk specific  
230 proteins including,  $\beta$ -casein, free secretory component, and lactoferrin.



231 eMPCs were grown using conditions that induce lactogenic differentiation of  
232 mammary cells with prolactin. Under these conditions, e1, eG1, eG2 and  
233 ePool formed domes that produce milk protein, indicative of lactogenic  
234 differentiation (Figure 6A). In these assays, e1 and ePool produced the  
235 greatest number of alveoli-like structures, with an alveologenic capacity to  
236 produce milk similar to postnatal MECs. eG1, and eG2 also showed  
237 lactogenic ability, but it was not possible to directly compare with the others  
238 since it was not possible to maintain the confluent cell culture needed for  
239 assessment at identical timepoints with this assay (Figure 6B).

240

#### 241 **Sox9 ablation in e1 leads to reduced expression of Zeb1**

242 Mammary tissues originate from multipotent embryonic progenitors,  
243 which give rise to unipotent basal and luminal stem cells found in postnatal  
244 mammary tissues<sup>31</sup>. *Sox9* was selected for further study since it is amongst  
245 the earliest transcription factors expressed by cells within the embryonic  
246 mammary epithelium and its expression increases so that most cells express  
247 high levels of *Sox9* by E14.5 (Supplementary Fig. 4)<sup>22</sup>. We assessed the role  
248 of *Sox9* in regulating embryonic mammary stem cell function using e1, a line  
249 that expresses a number of markers typically associated with stem/progenitor  
250 cells, including high levels of *Sox9*. CRISPR-Cas9 was used to create 3  
251 independent subclones with gene deletions (or knockouts (KO)) of *Sox9*  
252 derived from e1, e1/*Sox9*-KO#1, e1/*Sox9*-KO#2, and e1-*Sox9*-KO#3 (Figure  
253 7A).

254 We confirmed a reduction in *Sox9* levels in the 3 subclones with *Sox9*  
255 deletions compared to non-targeting guide control cells by qRT-PCR (Figure  
256 7B), immunohistochemistry (Figure 7C), and western blotting (Supplementary  
257 Fig. 5B). We also detected a marked reduction of expression of *Zeb1*, a key  
258 epithelial-to-mesenchymal transition (EMT) regulator, and concomitant with  
259 the level of reduction of *Sox9* (Figure 7B). Using IHC, we confirmed total loss  
260 of *Sox9* expression in e1/*Sox9*-KO#1 subclone, but detected *Sox9*<sup>+</sup> cells in 1-  
261 4% of e1/*Sox9*-KO#2 and e1/*Sox9*-KO#3 cells (Figure 7C). Reduced

262 detection of Sox9 in the 3 Sox9-KO subclones is correlated with reduction of  
263 Zeb1-expressing cells (Figure 7C).

264

### 265 **Sox9 ablation increases e1 clonogenicity**

266 In non-targeting guide control cells derived from e1, (e1/Co#1 and  
267 e1/Co#2), sphere-forming efficiency is similar in both the parental e1 clone  
268 and non-targeted e1 subclones (Figure 5B and 7D). In all three  
269 independently-derived e1 subclones with Sox9 deletions, sphere-forming  
270 efficiencies increased by approximately two-fold, indicating that deletion of  
271 Sox9 in e1 enhances clonogenic ability, a measure of stem cell activity  
272 (Figure 7D). When plated as multicellular aggregates on low-attachment  
273 plates, no change in spheroid morphology is observed (Supplementary Fig.  
274 6). However, when spheroids are embedded with BME or matrigel, the  
275 spheroids from Sox9-KO cells are larger and exhibit a greater number of  
276 cellular protrusions when compared to controls indicative of altered  
277 morphogenetic and enhanced migratory capacity (Figure 7E-F and  
278 Supplementary Figure 6).

279

### 280 **Sox9 ablation decreases response to lactogenic stimuli**

281 We assessed expression of several stem/lineage markers in e1/Sox9-  
282 KO#1, which lacks detectable Sox9 when cells are stained by IHC. *Acta2* ( $\alpha$ -  
283 SMA), a marker normally expressed by basal stem cells and differentiated  
284 myoepithelial cells, was increased when Sox9 was ablated. Levels of *Procr*, a  
285 marker for a rare basal mammary stem cell population and *Trp63*, a basal cell  
286 marker, were both decreased, suggesting profound alterations of the basal  
287 features of these cells (Figure 8A). We found other EMT regulators including,  
288 *Snail2/Slug* and *Twist* were expressed at lower levels in cells lacking Sox9,  
289 whilst *E-cadherin* levels, an epithelial marker, were slightly increased,  
290 consistent with a reduced EMT signature that would be expected from  
291 reduced Zeb1 levels.

292 RNA sequencing analysis confirmed that Sox9 transcripts were not  
293 completely deleted from the three e1/Sox9-KO subclones. This was

294 anticipated since the targeting strategy does not remove the 5' transcript,  
295 including ATG start codon). e1/Sox9-KO#1 expresses the lowest Sox9 levels  
296 compared to e1/Sox9-KO#2 and -KO#3 which show reduced Sox9 levels  
297 compared to control cells (Supplementary Fig. 7). A small number of genes  
298 were consistently up- and down-regulated in all three e1/Sox9-KO lines when  
299 compared to e1/control subclones in transcriptomic analysis (Figure 8B). One  
300 is *Csf1*, which promotes postnatal mammary stem cell activity, and is  
301 expressed at higher levels in all Sox9-deficient e1 subclones. Two clusters of  
302 interest were identified after clustering of all hits identified based on either  
303 DESeq or Intensity Difference tests. One group includes Sox9, whilst the  
304 other group displays the opposite expression pattern. e1/Sox9-KO#1 and  
305 e1/Sox9-KO#3 clustered more closely together than with e1/Sox9-KO#2, and  
306 both lines displayed pronounced changes in morphologies when cultured in  
307 3D, so these two lines were analysed together to identify genes consistently  
308 modulated in these e1/Sox9-KO cells compared to the control cells. Genes  
309 associated with innate immunity, cholesterol biosynthesis, cell fate are down-  
310 regulated, while genes regulating cell adhesion, Rap1 signalling and kinase  
311 activity are upregulated in e1/Sox9-KO#1 and e1/Sox9-KO#3 cells, identifying  
312 potential regulators of the morphological changes observed when cultured in  
313 3D (Supplementary Dataset 3).

314 Results from lineage tracing studies have shown Sox9<sup>+</sup> cells contribute  
315 to alveologenesis<sup>32</sup>. We therefore investigated the ability of e1/Sox9-KO cells  
316 to form fully differentiated mammary acini and secrete milk using an *in vitro*  
317 assay. Fewer alveoli-like structures derived from e1/Sox9-KO cells formed  
318 from all three e1/Sox9-KO lines and these were smaller and ill-defined  
319 structures compared to those formed from e1/control cells (Figure 8C-D). Our  
320 results demonstrate a diminished ability of e1/Sox9-KO embryonic mammary  
321 cells to respond to lactogenic stimuli and undergo functional alveologenesis.  
322 These results indicate Sox9 expression is required for embryonic mammary  
323 progenitor cells to attain a mature luminal progenitor cell phenotype capable  
324 of efficiently differentiating to form fully functional alveoli.

325

## 326 Discussion

327 In this study, we have established eighteen novel, cultivable, embryonic  
328 mammary progenitor cell lines (eMPCs). eMPCs will be a valuable resource  
329 for studying molecular regulators of normal breast development and  
330 understanding mammary cell fate and its plasticity. Based on gene expression  
331 and principal component analysis of the embryonic mammary progenitor cell  
332 lines, we identified three major clusters. One cluster (marked in yellow in  
333 Figure 2 and eleven out of the sixteen clones) is likely to represent mammary  
334 epithelial progenitor cells in an undifferentiated, plastic state since these cells  
335 express high levels of several stem cell makers; another cluster of 3 clones  
336 (marked in green in Figure 2) is likely to be representing more differentiated  
337 epithelial mammary progenitor cells since cells from this cluster express  
338 higher levels of Keratins. The third cluster of two clones and the pool (marked  
339 in magenta in Figure 2) appears to be composed of less plastic progenitor  
340 cells with restricted developmental potential. Four of the eighteen cell lines  
341 were selected for use in plasticity assays to represent the three major clusters  
342 and included both GFP- clones we had derived from single GFP- cells, as well  
343 as two of the sixteen GFP+ clones derived from single GFP+ cells. Although  
344 GFP+ expression was retained within the epithelial-appearing component  
345 after sustained organoid culture of the microdissected mammary organ, GFP  
346 expression is variable after expansion in 2D culture such that both GFP+ and  
347 GFP- cells are present in all eMPC lines, and we cannot assume GFP+  
348 clones are derived from embryonic mammary epithelial cells or that GFP-  
349 cells represent mesenchymal cells. We cannot exclude the possibility that  
350 GFP+ clones could have originated from other GFP+ cells present within the  
351 embryo although our dissection protocol should exclude this possibility<sup>21</sup>. In  
352 addition, we observe a small percentage of cells that lack GFP expression  
353 within the mammary primordial epithelium.

354 Embryonic mammary lineage commitment is not completely  
355 understood<sup>33</sup>. It is still not clear how plastic or committed embryonic  
356 mammary cells are after removal from their native microenvironment. A  
357 subset of clones appeared to give rise to cell lines with mammary epithelial  
358 progenitor features and other clones appear less differentiated and more

359 plastic. Another small group clustered together with the ePool suggesting it  
360 consisted of a mixture of cell types. Heterogeneity exists between the clone  
361 profiles and their properties. This may reflect different levels of commitment of  
362 individual cells to the mammary lineage or loss of mammary epithelial  
363 phenotype. In addition, stem cells exist in a variety or spectrum of cell states  
364 that can interconvert<sup>34,35</sup>, so the observed heterogeneity could represent lines  
365 comprised of various components of a mammary progenitor spectrum.

366 During embryonic development, tissues reshape during a combination  
367 of morphogenetic processes. Tissue compaction is a morphogenetic process  
368 by which a tissue adopts a tighter structure. Both cell adhesion and cell  
369 contraction play roles in tissue compaction<sup>36</sup>. Of the four clones assayed, only  
370 eG1, a line that expresses Keratins, as well as other adhesion regulators  
371 (including protocadherins, integrins, ECM ligands) undergoes compaction.  
372 eMPC are clonogenic and form mammary acini *in vitro*. All four of the eMPC  
373 lines tested can be induced to secrete milk with the appropriate hormonal  
374 stimulation (although e2 and eG2 with low efficiency), which demonstrates  
375 these cells have been committed to the mammary lineage, an event that  
376 occurs during embryonic development<sup>6</sup>. Clone e1 displayed limited ability to  
377 engraft in cleared mammary fat pads; this is consistent with other reports that  
378 found dissociated mid-gestation stage mouse mammary cells have very low  
379 stem cell activity when assessed using the cleared fat pad assay<sup>12</sup>. Baseline  
380 sequencing permits monitoring for changes after passage and we have been  
381 able to recluster eMPC cells together after early and later passages when re-  
382 sequenced, suggesting these cells can be maintained without substantial  
383 changes.

384 Sox9 is amongst the earliest transcription factors expressed by cells  
385 within the prenatal mammary epithelium<sup>22</sup>. Sox9 is also expressed by basal  
386 and luminal Estrogen Receptor (ER)- postnatal mammary epithelial cells<sup>32,37</sup>.  
387 Sox9 is thought to act together with Slug to confer a stem cell state to  
388 postnatal mammary epithelial cells<sup>18</sup>. Regulation of prenatal and postnatal  
389 mammary stem cell function by Sox9 produces distinct effects on clonogenic  
390 ability. In contrast to our finding of increased sphere-forming ability in e1 after  
391 Sox9 deletion, knockdown of Sox9 in primary postnatal mouse MECs causes

392 a marked reduction in organoid-forming ability, reduced *Slug* levels, and  
393 diminishes gland-reconstituting activity<sup>18</sup>. Impaired ability to undergo EMT is  
394 likely to be a common feature in all *Sox9*-deficient mammary cells. Mammary  
395 stem cell function is highly impacted by *Sox9*, but cell context also appears to  
396 be a significant factor, including developmental stage and state. A recent  
397 study shows a role for SOX9 in the regulation of hormone resistance in breast  
398 cancer<sup>38</sup>, highlighting the relevance of embryonic mammary factor expression  
399 in breast cancers and the need for further investigations into their ability to  
400 modulate cell states.

401 Cell state is thought to influence epithelial cell phenotype and  
402 plasticity<sup>39</sup>. Mature postnatal MECs express Keratins, indicative of being in an  
403 epithelial state (E). Loss of *Sox9* and ability to undergo EMT in postnatal  
404 MECs should push them towards a more highly epithelial state, which, in  
405 theory, would have less stem cell activity according to several recent  
406 models<sup>40</sup>. Stem cell activity is thought to be highest when cells are in an  
407 intermediate epithelial/mesenchymal state<sup>35</sup>. Deleting *Sox9* from e1 leads to  
408 reduced *Zeb1* levels and reduced propensity to undergo EMT. As a result, e1/  
409 *Sox9*-KO cells would be expected to shift from the embryonic mesenchymal  
410 (M) state towards epithelial to an intermediate hybrid M/E state, which, in  
411 theory, should have more stem cell activity and mixed epithelial (adhesive)  
412 and mesenchymal (migratory) properties that confer ability to move  
413 collectively as clusters (Supplementary Fig. 8)<sup>41,42</sup>. Our results clearly show  
414 that there are profound differences between the effects from loss of *Sox9*  
415 function on clonogenic ability in embryonic mammary progenitor cells  
416 compared to postnatal MECs, which exist in relatively high epithelial cell  
417 states. Despite these differences, these findings are consistent with current  
418 models of EMT/MET mediation of cell state in regulating stem cell function.

419 Analysis from lineage tracing data indicated that *Sox9* marks luminal  
420 progenitors that give rise to alveolar progenitors<sup>32</sup>. Fewer alveoli-like  
421 structures formed from e1/*Sox9*-KO cells and were smaller than control  
422 alveoli, which is consistent with a requirement for active *Sox9* signalling in  
423 luminal progenitors to produce functional alveoli. *Sox9* is required for luminal  
424 progenitor cell survival, but basal/myoepithelial cells can withstand loss of

425 Sox9 in studies using a conditional deletion of Sox9 from the mammary gland  
426 using Mouse mammary tumour virus (MMTV)-Cre<sup>37</sup>. Mice with Sox9-deficient  
427 mammary glands could lactate normally, but mosaic MMTV-Cre mediated  
428 gene deletion was reported, so it is likely that sufficient numbers of Sox9+  
429 luminal progenitors were retained to produce functional alveoli that had  
430 undergone terminal differentiation<sup>37</sup>. Our results support a model in which  
431 Sox9 is required for the specification of mammary progenitor cells from early  
432 embryonic stages of through postnatal development as a key regulator of  
433 mammary gland development and luminal progenitor homeostasis and  
434 maintenance.

435 Our study reveals complex interactions that underlie mammary lineage  
436 regulation, cell state and progenitor cell activity. Plasticity, defined as the  
437 ability of cells to dynamically change cell state, is crucial for many processes  
438 during mammary gland development and tissue homeostasis. Reversible  
439 EMT is central to tissue development, epithelial stemness, and cancer  
440 metastasis. SOX9 has been linked to all three of these processes<sup>18,37,43</sup>.

441 Cell context is key in studies of mammary cell function and embryonic  
442 cells are highly plastic and multipotent, giving rise to both basal and luminal  
443 postnatal mammary cells. Cancer progression involves the loss of  
444 differentiated phenotype and acquisition of progenitor/stem cell features,  
445 which have recently been put forward as a hallmark of cancer<sup>44</sup>. These cell  
446 lines will be an invaluable resource to explore unresolved questions related to  
447 biology of mammary progenitor and stem cells, including those relevant to the  
448 origin of breast cancer.

449  
450

451 **Methods**

452 **Animal experiments**

453 All animal work was carried out under UK Home Office project and personal  
454 licenses following local ethical approval from The Institute of Cancer  
455 Research Ethics Committee and in accordance with local and national  
456 guidelines. Female mice (*Mus musculus*) were used for all experiments  
457 except breeding. Immortomouse and SCID/Beige mice were purchased from  
458 Charles River (Harlow, UK). *s-SHIP-GFP* mice were a gift from the late  
459 Professor Larry Rohrschneider (Fred Hutchinson Cancer Research Center,  
460 Seattle, WA, USA). Mice were housed in IVC cages on a 12h light/dark cycle  
461 and received food and water *ad libitum*. Mouse genotyping for detecting  
462 *Immorto* and *s-SHIP-GFP* transgenes was performed at Transnetyx (Cordova,  
463 TN, USA) and the results were used for colony maintenance.

464

465 **Isolation and culture of eMPC**

466 Immortomouse were backcrossed with C57BL6/J (B6) for four generations.  
467 The Immortomouse and *s-SHIP-GFP* mice were bred to obtain *Immorto:s-*  
468 *SHIP-GFP* embryo. E12.0-stage *Immorto:s-SHIP-GFP* embryo was dissected  
469 to isolate only the organ of mammary primordium number three as a micro-  
470 tissue, using Dumont #5 forceps and Tübingen Spring Scissors (Fine Science  
471 Tools GmbH, Heidelberg, Germany) using a Leica M205FA fluorescent  
472 stereomicroscope (Leica, Wetzlar, Germany). The remaining embryonic tissue  
473 was used for genotyping the *Immorto* transgene. Micro-tissue was  
474 immediately embedded into 200 µl of ice-cold Cultrex BME type 1, stem cell  
475 qualified, growth factor reduced (3434-005-02, AMS Biotechnology Ltd,  
476 Abingdon, Oxford, UK) within a well of a 48-well plate. The matrix-embedded  
477 tissue was transferred into CO<sub>2</sub> incubator equilibrated to 33°C, 5% CO<sub>2</sub>  
478 atmosphere for 3 hours. Culture medium was overlaid on the solidified BME  
479 during the pre-incubation. Culture medium contains Mesencult MSC basal  
480 medium for mouse (05501), Mesencult Stimulatory Supplements for Mouse  
481 (05502), 2 mM L-Glutamine (17100; all from STEMCELL Technologies UK  
482 Ltd., Cambridge, UK), 5 U/ml Mouse Interferon gamma (IFN $\gamma$ ) (315-05, Pepro  
483 Tech EC Ltd., London, UK), and 0.1 mg/ml Primocin (ant-pm1, Source



484 Bioscience, Nottingham, UK). eMPCs migrated and proliferated within BME  
485 from the *Immorto:s-SHIP-GFP* micro-tissue for 1 month. Cells were harvested  
486 by trypsinisation and further expanded on two-dimensional dishes pre-coated  
487 thinly with BME, and referred to henceforth, as eMPC pool, abbreviated as  
488 ePool.

489

## 490 **Generation of eMPC clones**

### 491 ***Fluorescence activated cell sorting***

492 ePool cells were harvested from the culture dish and adjusted to the  
493 concentration of  $2 \times 10^6$  cells/ml with PBS. To stain apoptotic and necrotic  
494 cells, cells were stained with DAPI at the final concentration of 1  $\mu\text{g/ml}$  for 5  
495 min at room temperature prior to cell sorting. ePool were passed through a 40  
496  $\mu\text{m}$  filter and sorted using BD Aria fluorescence-activated cell sorter (BD  
497 Biosciences, Oxford, UK). Live single ePool cells were positively gated using  
498 their forward and side scatter profiles (Supplementary Fig. 1B). ePool  
499 population with high s-SHIP-GFP expression and ePool GFP- population were  
500 subsequently gated based on green channel profile.

### 501 ***Clone derivation from single cells***

502 Live single ePool GFP+ and ePool GFP- cells were individually sorted into 96  
503 multiwell plates. The multiwell plates were precoated with either Cultrex BME  
504 or Matrigel, growth factor reduced (734-0268, VWR International, Lutterworth,  
505 Leicestershire, UK) and filled with eMPC culture medium. The single cells on  
506 multiwell plates were incubated at CO<sub>2</sub> incubator equilibrated to 33°C, 5%  
507 CO<sub>2</sub> atmosphere without changing medium for a month. eMPC clones  
508 proliferated for a month after which they were trypsinised and transferred to  
509 dishes precoated with BME for further expansion. Only two GFP- eMPC  
510 clones were grown from a Matrigel-precoated multiwell plate and named as  
511 e1 and e2 (Figure 1B). Fourteen GFP+ eMPC clones were recovered from  
512 two Matrigel-precoated multiwell plates, which are named as eG1, eG2,  
513 eG2E9..., depending on well position. Only one GFP+ eMPC clone, eG1B3,  
514 was recovered from a BME-precoated plate, suggesting that Matrigel is  
515 preferable for initial clonal cell expansion.

516

517 **Cell culture**

518 The following mouse C57BL/6 (B6) cell lines were purchased and cultured for  
519 RNA isolation according to manufacturer's instructions; Bone Marrow  
520 Mononuclear cells, (c57-6271F, Caltag Medsystems Buckingham, UK);  
521 Embryonic Stem Cells, (SCRC-1002, LGC standards, Middlesex, UK); Mouse  
522 Embryonic Fibroblasts, untreated, (GSC-6002, AMS Biotechnology Ltd);  
523 Mesenchymal Stem cells, (MUBMX-01001 Cambridge Bioscience,  
524 Cambridge, UK). B6 postnatal mammary epithelial cells were isolated as  
525 described<sup>45</sup>. The mammary epithelial fraction was further purified using  
526 EasySep Mouse Epithelial Cell Enrichment Kit according to manufacturer's  
527 instructions (19758, STEMCELL Technologies UK). Human Umbilical Vein  
528 Endothelial Cells were cultured according to manufacturer's instruction  
529 (191027, Lonza, Basel, Switzerland).

530

531 **Total RNA isolation and quality control**

532 Embryonic mammary tissues (embryonic mammary primordial epithelium  
533 (MPE), embryonic mammary stroma (includes mammary mesenchyme (MM)  
534 and fat pad precursor (FPP)), and surface epithelium (Epi), were micro-  
535 dissected from E12.5 B6 mice following the protocol described in<sup>46</sup>. Total RNA  
536 from cells and mammary tissue were extracted using RNeasy Plus Micro kit  
537 and RNeasy Mini kit (74034 and 74104, respectively, Qiagen, Manchester,  
538 UK), RNA concentration and gDNA contamination were determined using  
539 Qubit 2.0 Fluorometer and accompanying kits (Thermo Fisher Scientific,  
540 Waltham, MA, USA). RNA samples highly contaminated with gDNA for more  
541 than 10% of RNA were treated with TURBO DNA-free Kit (AM1907, Thermo  
542 Fisher Scientific) and concentrated using RNA Clean and Concentrator-5  
543 (R1015, Zymo Research, Irvine, CA, USA). All samples were assessed for  
544 RNA quality using Agilent 2100 Bioanalyzer system (Agilent Technologies,  
545 Cheshire, UK) and confirmed that RNA Integrity Number was greater than 9.0.

546

547 **RNA sequencing**

548 cDNA library preparation was carried out at Oxford Genomics Centre, The  
549 Wellcome Trust Centre for Human Genetics using PolyA+ RNA enrichment  
550 method for total RNA from cultured cells and SMARTer method for total RNA

551 from embryonic mammary tissue, respectively. mRNA fraction was selected  
552 from the total RNA before conversion to cDNA. Second strand cDNA  
553 synthesis incorporated dUTP. The cDNA was end-repaired, A-tailed and  
554 adapter-ligated. Prior to amplification, samples underwent uridine digestion.  
555 The prepared libraries were size selected, multiplexed and quality checked  
556 before paired end sequencing over three lanes of a flow cell. Amplified cDNA  
557 from embryonic mammary tissues were generated by the SMARTer  
558 amplification kit. The cDNA was end-repaired, A-tailed, adapter-ligated and  
559 amplified. The prepared libraries were size selected, multiplexed and quality  
560 checked before paired end sequencing on four lanes of a flow cell. Data was  
561 aligned to the reference genome, mm10, and quality checked.

562 RNA sequencing files were submitted to ArrayExpress as accession E-  
563 MTAB-6846, MTAB-6856, E-MTAB-6859. FastQ files were truncated to a  
564 consistent length of 75bp using trim galore v0.4.3 and were then aligned  
565 against the mouse GRCm38 genome assembly using hisat2 v2.0.5 using  
566 options --no-mixed and --no-discordant. Mapped positions with MAPQ values  
567 of <20 were discarded. Gene expression was quantitated using the RNA-Seq  
568 quantitation pipeline in SeqMonk software v1.37.0  
569 (<https://www.bioinformatics.babraham.ac.uk/projects/seqmonk/>) in opposing  
570 strand specific library mode. For count based statistics, raw read counts over  
571 exons in each gene were used. For visualisation and other statistics  
572 log2RPM (reads per million reads of library) expression values were used.  
573 Differentially expressed genes were selected based on passing two statistical  
574 filters -the DESeq2 Likelihood Ratio Test with a cutoff of  $p < 0.05$  following  
575 multiple testing correction and the SeqMonk Intensity Difference filter on  
576 log2RPM values with a sample size of 1% of all genes and a cutoff of  $p < 0.05$   
577 after multiple testing correction. Hierarchical clustering was performed on per-  
578 gene median centred log2RPM expression values using Pearson's  
579 correlation. Gene cluster separation was performed by segmenting the tree at  
580 an R value of 0.5. Clusters containing <50 genes were discarded. PCA was  
581 performed on column centred log2RPM values without additional scaling.

582 The intensity difference test used a locally matched subset of 1% of  
583 genes based on average expression. From these a local standard deviation in  
584 log2RPM difference values was calculated and used to calculate the

585 probability of the cumulative distribution function for a normal distribution with  
586 this standard deviation, using the observed difference in the gene being  
587 tested. P-values were corrected for multiple testing using Benjamini and  
588 Hochberg multiple testing correction.

589

#### 590 **cDNA synthesis and real-time quantitative reverse transcription PCR**

591 cDNA library from total RNA was synthesised using QuantiTect Reverse  
592 Transcription Kit (205310, Qiagen). Comparative  $C_T$  ( $\Delta\Delta C_T$ ) real-time  
593 **quantitative** PCR was performed using Taqman Gene Expression assays as  
594 previously described<sup>17</sup> on either ABI Prism 7900HT sequence detection  
595 system (Applied Biosystems, Foster City, CA, USA) or QuantStudio 6 Flex  
596 system (Thermo Fisher Scientific). *Actb* was used as an endogenous control  
597 and fold change normalised to a comparator sample was calculated. The  
598 assay probes are listed in Supplementary Dataset 4.

599

#### 600 **Antibodies**

601 Antibodies used for whole-mount immunofluorescence,  
602 immunohistochemistry, immunofluorescence, and western blotting in this  
603 study are listed in Supplementary Dataset 4.

604

#### 605 **Adipogenesis assay**

606 eMPC and MSCs (positive control) were plated at 100% confluency in  
607 Osteogenic/Adipogenic Base Media (CCM007, Bio-Techne Ltd, Oxford, UK)  
608 containing 1  $\mu\text{g/ml}$  Insulin, 0.5 mM 3-Isobutyl-1-methylxanthine and 0.25 mM  
609 Dexamethasone (19278, I5879 and D2915 from Sigma, Dorset, UK,  
610 respectively). After 6 days, cells were fixed with 4% Paraformaldehyde/PBS  
611 and lipid droplets were detected by HCS LipidTOX™ Red Neutral Lipid Stain  
612 (H34476, Thermo Fisher Scientific) at 1:200 in PBS. Lipid-labelled cells were  
613 photographed using EVOS FL fluorescent microscope equipped with 10x Plan  
614 LWD FL 0.30NA objective (Thermo Fisher Scientific).

615

#### 616 ***In vitro* tube formation assay**

617 Tube formation ability of eMPCs was examined following a published  
618 protocol<sup>47</sup>. Briefly, eMPCs, MSCs (negative control) and HUVECs (positive  
619 control) were plated at ~80% confluency using endothelial growth medium  
620 (EGM) (EBM-2 medium containing EGM-2 SingleQuot Kit, CC-3156 and CC-  
621 3162 from Lonza, respectively). On the next day, 96 well plates filled with 50-  
622 80 µl Cultrex BME were prepared on ice without forming air bubbles and pre-  
623 incubated at CO<sub>2</sub> incubator equilibrated to 37°C, 5% CO<sub>2</sub> atmosphere to allow  
624 matrix gelation. Cells were harvested by trypsinisation and re-suspended at  
625 the concentration of 1.5 x 10<sup>5</sup> cells /ml in either endothelial basal medium 2  
626 (EBM-2) or EGM-2 medium. 100 µl of single cell suspension (containing total  
627 15,000 cells) was gently added on the top of BME gel. After incubation for 24  
628 hours, cells were labelled with Calcein AM solution at the final concentration  
629 of 2 µM in EBM-2 medium and incubated in CO<sub>2</sub> incubator for 30 min. Calcein  
630 AM-labelled cells were photographed using EVOS FL fluorescent microscope  
631 equipped with 2x Plan LWD 0.06NA objective (Thermo Fisher Scientific). The  
632 fluorescent images were analysed to detect mesh structure using a developed  
633 package of ImageJ, Fiji<sup>48</sup> (<http://imagej.net/Fiji/Downloads>), and its Plugin,  
634 Angiogenesis Analyzer (Carpentier G., Angiogenesis Analyzer for ImageJ  
635 (2012);  
636 <https://imagej.nih.gov/ij/macros/toolsets/Angiogenesis%20Analyzer.txt>). Cells  
637 were also tested for tube formation assay with EBM-2 medium containing 10  
638 mM Sulforaphane (angiogenic inhibitor, S4441, Merck KGaA, Darmstadt,  
639 Germany) as a control to observe its inhibitory effect on the assay.

640

#### 641 **Mammary alveologenesis assay**

642 Lactogenic ability of eMPCs was examined following an established  
643 protocol<sup>49</sup>. Briefly, eMPCs and MECs (positive control) were plated at 100%  
644 confluent in RPMI 1640 medium (21875, Thermo Fisher Scientific) containing  
645 10% horse serum (16050-122, Thermo Fisher Scientific), 5 µg/ml Insulin, 10  
646 ng/ml EGF (236-EG-200, Bio-Techne Ltd), 10 mM HEPES (H0887, Sigma),  
647 and Primocin (Bioscience). On the following day, in order to withdraw the EGF  
648 and enhance differentiation, rather than cell proliferation, cell medium was  
649 changed to EGF- medium containing 5% horse serum, 10 mM HEPES, and

650 Primocin in RPMI 1640 medium. After 3 days, mammary alveologenesis was  
651 induced by replacement with priming medium (EGF- medium containing 5  
652  $\mu\text{g/ml}$  mouse Prolactin (757908, Biolegend, London, UK), 5  $\mu\text{g/ml}$  Insulin, and  
653 0.1 mM Dexamethasone). Cells were fixed with 4% Paraformaldehyde/PBS 6  
654 days after induction and stained with anti-Milk antibody and DAPI. Alveoli  
655 were photographed using EVOS FL microscope and the number of alveoli  
656 larger than 100  $\mu\text{m}$  per 24 well was counted manually.

657 For studies of e1 with reduced Sox9 expression, the assay was  
658 modified since e1/Sox9-knockout cells detached upon induction of lactogenic  
659 differentiation using the standard assay. Cells were plated around 80-90%  
660 confluence in MesenCult Expansion medium with MesenCult supplement, 200  
661 mM L-Glutamine, 5 U/ml IFN $\gamma$  (Pepro Tech) and 0.1mg/ml Primocin. Once  
662 cells reached 100% confluence (on day 4), cell medium was changed to RPMI  
663 1640 medium containing 5% horse serum, 10 mM HEPES, and Primocin. On  
664 day 5, mammary alveologenesis was induced with medium containing 5  $\mu\text{g/ml}$   
665 mouse Prolactin, 5  $\mu\text{g/ml}$  Insulin, and 0.1 mM Dexamethasone. Induction  
666 medium was replenished every 3 days. On day 11, cells were fixed with 4%  
667 Paraformaldehyde/PBS and stained with anti-Milk antibody and DAPI. Alveoli  
668 were photographed using EVOS FL microscope and the number of alveoli in  
669 different size ranges were counted manually.

670

### 671 **Sphere-forming assays**

672 A clonogenic mammosphere protocol and media were adapted for eMPCs<sup>50</sup>.  
673 Mouse EpiCult-B complete medium (05610, STEMCELL Technologies) was  
674 supplemented with 0.6% Methylcellulose (HSC011, Bio-Techne Ltd), 20 ng/ml  
675 basic FGF (3139-FB-025/CF, Bio-Techne Ltd), 20 ng/ml EGF, 2%  
676 NeuroCult™ SM1 Without Vitamin A, 10  $\mu\text{g/ml}$  Heparin, 1  $\mu\text{g/ml}$   
677 Hydrocortisone (05731, 07980, and 07926, respectively, from STEMCELL  
678 Technologies UK Ltd), 10  $\mu\text{g/ml}$  Insulin, 5 U/ml Mouse IFN $\gamma$ , and 0.1 mg/ml  
679 Primocin. Trypsinised eMPCs and MSCs were passed through 40  $\mu\text{m}$  filter to  
680 make absolute single cell suspension and the cell number was adjusted to 2-4  
681  $\times 10^5$  cells/ml. 10,000 cells were added into 2 ml of the supplemented Epicult-  
682 B medium. Cells were pipetted thoroughly to be distributed homogenously

683 within the viscous Methylcellulose-containing medium and transferred into  
684 Corning ultra-low attachment 6 well (10154431, Fisher Scientific,  
685 Loughborough, UK). 6 well plates were incubated within a humidified box at  
686 33°C incubator for eMPCs and 37°C for MSCs, respectively. At day 7, sphere  
687 images were photographed and automatically scored using Single Colony  
688 Verification program in Celigo cytometer (Nexcelom Bioscience LLC,  
689 Lawrence, MA, USA).

690

### 691 **Spheroid growth in hydrogel**

692 For analysing 3D morphology, e1/control and e1/Sox9-KO cells were seeded  
693 with 5000 cells/well onto Corning Costar Ultra-Low Attachment 96 well plates  
694 (Corning Kennebunk ME, USA) filled with eMPC culture medium. On day 3,  
695 50 µl of Corning Matrigel Growth Factor Reduced Basement Membrane  
696 Matrix (356231) or Cultrex PathClear 3-D Culture Matrix RGF BME I (3434-  
697 005) was added into each well (except control wells). Medium was  
698 replenished every 3 days. On day 9, 50 µl of 1% agarose was added to the  
699 control wells. All spheroids were subsequently fixed with 4%  
700 Paraformaldehyde (PFA)/PBS and were photographed using EVOS FL  
701 microscope.

702

### 703 **CRISPR/Cas9-mediated gene targeting of Sox9**

704 We took advantage of lentiviral transduction of *Cas9* gene within eMPCs  
705 followed by transfection of pre-made CRISPR RNA (crRNA) and trans-  
706 activating crRNA (tracrRNA) for genetic ablation of *Sox9*, which were  
707 purchased from GE Healthcare (Little Chalfont, UK). Edit-R Lentiviral Blast-  
708 Cas9 Nuclease Particles (VCAS10128, Dharmacon, Lafayette, Colorado,  
709 USA) were transduced into e1 and *Cas9*-expressing e1 population was  
710 expanded with 10 µg/ml Blasticidin (ant-bl-1, Source Bioscience Nottingham,  
711 UK) from day 4 onward.  $2 \times 10^5$  cells of Blasticidin-resistant e1 were seeded  
712 on 6 well and on the next day treated with 5 µl DharmaFECT 1 Transfection  
713 Reagent (T-2001-02, Dharmacon) together with 5 µl of 10 µM Edit-R  
714 CRISPR-Cas9 Synthetic tracrRNA (U-002000-20, Dharmacon) and 5 µl of 10  
715 µM *Sox9*-targeting crRNA or non-targeting crRNA. Cell populations with *Sox9*

716 mutations were screened by Heteroduplex formation assay using  
717 QuickExtract DNA Extraction Solution (QE0905T, Qiagen, Peterborough, UK)  
718 and EnGen™ Mutation Detection Kit (E3321S, New England Biolabs, Herts,  
719 UK) according to manufacturer's instructions. We tested 4 different crRNA  
720 and the primer sets to detect mutation within each crRNA-binding region;  
721 details are provided in Supplementary Dataset 4C. 2 cell populations (Sox9-  
722 01 and Sox9-02 in Supplementary Fig. 5A) that were verified for Sox9  
723 mutation were further sorted using BD Aria to develop single cell-derived  
724 subclones. 3 e1/Sox9-knockout (KO) subclones, KO#1, KO#2, and KO#3,  
725 produced no Sox9 expression detectable by Western blot (Supplementary  
726 Fig. 5B). These three e1/Sox9-KO subclones as well as two e1/control clones,  
727 Co#1 and Co#2, were used in this study.

728

#### 729 **Mammary fat pad injection of e1**

730 To allow bioluminescence imaging of e1 cells, cells were labelled with  
731 carrying red-shifted *Luciola Italica* luciferase transgene using lentivirus  
732 particles, RediFect Red-FLuc-Puromycin (CLS960002, Perkin Elmer,  
733 Buckinghamshire, UK). The transduced e1/Red-FLuc cells were selected with  
734 5 µg/ml Puromycin for 3 weeks. e1/Red-FLuc cells were harvested and  
735 resuspended at a concentration of  $1 \times 10^6$  cells/ml in PBS. 100 µl of cells were  
736 injected into both the left and right inguinal mammary fat pad number 4 of 10-  
737 week-old SCID/Beige female mice. Mice were injected 100 µl of 15 mg/ml D-  
738 Luciferin (119222, Perkin Elmer) in PBS intraperitoneally and imaged after 5  
739 min using IVIS Illumina II (Perkin Elmer). Engrafted mammary glands were  
740 harvested 1-2 weeks after mammary fat pad injections and fixed in 4% PFA in  
741 PBS for immunohistochemistry.

742

#### 743 **Statistical Analysis**

744 The data in the graphs are presented as mean and the standard error of the  
745 mean. The data was analysed by two-tailed ANOVA, Student's t-test, or two-  
746 tailed, paired t-test using GraphPad Prism 7 software. *P*-value  $\leq 0.0001$  is  
747 considered as extremely significant (\*\*\*\*), *P*  $\leq 0.001$  as highly significant (\*\*\*),  
748 *P*  $\leq 0.01$  as very significant (\*\*), *P*  $\leq 0.05$  as significant (\*), and *P*  $> 0.05$  as  
749 not significant (ns), respectively.



750

751 **Data availability statement**

752 The authors declare that the data supporting the findings of this study are  
753 available within the article and its supplementary information files. All source  
754 data underlying the graphs and charts presented in the main figures is  
755 available in Supplementary Dataset 5.

756

757 **References**

- 758 1 Van Keymeulen, A. & Blanpain, C. Tracing epithelial stem cells during  
759 development, homeostasis, and repair. *The Journal of cell biology* **197**,  
760 575-584, doi:10.1083/jcb.201201041 (2012).
- 761 2 Wuidart, A. *et al.* Early lineage segregation of multipotent embryonic  
762 mammary gland progenitors. *Nat Cell Biol* **20**, 666-676,  
763 doi:10.1038/s41556-018-0095-2 (2018).
- 764 3 Lilja, A. M. *et al.* Clonal analysis of Notch1-expressing cells reveals the  
765 existence of unipotent stem cells that retain long-term plasticity in the  
766 embryonic mammary gland. *Nat Cell Biol* **20**, 677-687,  
767 doi:10.1038/s41556-018-0108-1 (2018).
- 768 4 Zvelebil, M. *et al.* Embryonic mammary signature subsets are activated in  
769 Brca1-/- and basal-like breast cancers. *Breast cancer research : BCR* **15**,  
770 R25, doi:10.1186/bcr3403 (2013).
- 771 5 Howard, B. & Ashworth, A. Signalling pathways implicated in early  
772 mammary gland morphogenesis and breast cancer. *PLoS Genet* **2**, e112,  
773 doi:10.1371/journal.pgen.0020112 (2006).
- 774 6 Propper, A. Y., Howard, B. A. & Veltmaat, J. M. Prenatal morphogenesis of  
775 mammary glands in mouse and rabbit. *Journal of mammary gland biology*  
776 *and neoplasia* **18**, 93-104, doi:10.1007/s10911-013-9298-0 (2013).
- 777 7 Balinsky, B. On the developmental processes in mammary glands and  
778 other epidermal structures. *Transactions of the Royal Society Edinburgh*  
779 **62**, 1-31 (1949-1950).
- 780 8 Lee, M. Y. *et al.* Ectodermal influx and cell hypertrophy provide early  
781 growth for all murine mammary rudiments, and are differentially  
782 regulated among them by Gli3. *PloS one* **6**, e26242,  
783 doi:10.1371/journal.pone.0026242 (2011).
- 784 9 Propper, A. Y. Wandering epithelial cells in the rabbit embryo milk line. A  
785 preliminary scanning electron microscope study. *Developmental biology*  
786 **67**, 225-231 (1978).
- 787 10 Propper, A. Y., Howard, B. A. & Veltmaat, J. M. Prenatal Morphogenesis of  
788 Mammary Glands in Mouse and Rabbit. *Journal of mammary gland biology*  
789 *and neoplasia*, doi:10.1007/s10911-013-9298-0 (2013).
- 790 11 Makarem, M. *et al.* Developmental changes in the in vitro activated  
791 regenerative activity of primitive mammary epithelial cells. *PLoS Biol* **11**,  
792 e1001630, doi:10.1371/journal.pbio.1001630 (2013).
- 793 12 Spike, B. T. *et al.* A mammary stem cell population identified and  
794 characterized in late embryogenesis reveals similarities to human breast  
795 cancer. *Cell stem cell* **10**, 183-197, doi:10.1016/j.stem.2011.12.018  
796 (2012).
- 797 13 Sakakura, T., Nishizuka, Y. & Dawe, C. J. Capacity of mammary fat pads of  
798 adult C3H/HeMs mice to interact morphogenetically with fetal mammary  
799 epithelium. *J Natl Cancer Inst* **63**, 733-736 (1979).
- 800 14 Voutilainen, M. *et al.* Ectodysplasin regulates hormone-independent  
801 mammary ductal morphogenesis via NF-kappaB. *Proceedings of the*  
802 *National Academy of Sciences of the United States of America* **109**, 5744-  
803 5749, doi:10.1073/pnas.1110627109 (2012).
- 804 15 Hens, J. R. *et al.* BMP4 and PTHrP interact to stimulate ductal outgrowth  
805 during embryonic mammary development and to inhibit hair follicle

806 induction. *Development* **134**, 1221-1230, doi:10.1242/dev.000182  
807 (2007).

808 16 Veltmaat, J. M. *et al.* Gli3-mediated somitic Fgf10 expression gradients are  
809 required for the induction and patterning of mammary epithelium along  
810 the embryonic axes. *Development* **133**, 2325-2335,  
811 doi:10.1242/dev.02394 (2006).

812 17 Kogata, N., Oliemuller, E., Wansbury, O. & Howard, B. A. Neuregulin-3  
813 regulates epithelial progenitor cell positioning and specifies mammary  
814 phenotype. *Stem Cells Dev* **23**, 2758-2770, doi:10.1089/scd.2014.0082  
815 (2014).

816 18 Guo, W. *et al.* Slug and Sox9 cooperatively determine the mammary stem  
817 cell state. *Cell* **148**, 1015-1028, doi:10.1016/j.cell.2012.02.008 (2012).

818 19 Jat, P. S. *et al.* Direct derivation of conditionally immortal cell lines from an  
819 H-2Kb-tsA58 transgenic mouse. *Proceedings of the National Academy of  
820 Sciences of the United States of America* **88**, 5096-5100 (1991).

821 20 Rohrschneider, L. R., Custodio, J. M., Anderson, T. A., Miller, C. P. & Gu, H.  
822 The intron 5/6 promoter region of the ship1 gene regulates expression in  
823 stem/progenitor cells of the mouse embryo. *Dev Biol* **283**, 503-521,  
824 doi:S0012-1606(05)00270-8 [pii]  
825 10.1016/j.ydbio.2005.04.032 (2005).

826 21 Bai, L. & Rohrschneider, L. R. s-SHIP promoter expression marks activated  
827 stem cells in developing mouse mammary tissue. *Genes Dev* **24**, 1882-  
828 1892, doi:10.1101/gad.1932810 (2010).

829 22 Wansbury, O. *et al.* Transcriptome analysis of embryonic mammary cells  
830 reveals insights into mammary lineage establishment. *Breast Cancer  
831 Research* **13**, doi:Artn R79  
832 Doi 10.1186/Bcr2928 (2011).

833 23 Sidney, L. E., Branch, M. J., Dunphy, S. E., Dua, H. S. & Hopkinson, A. Concise  
834 review: evidence for CD34 as a common marker for diverse progenitors.  
835 *Stem Cells* **32**, 1380-1389, doi:10.1002/stem.1661 (2014).

836 24 dos Santos, C. O. *et al.* Molecular hierarchy of mammary differentiation  
837 yields refined markers of mammary stem cells. *Proc Natl Acad Sci U S A*  
838 **110**, 7123-7130, doi:10.1073/pnas.1303919110 (2013).

839 25 Qian, H. *et al.* Molecular characterization of prospectively isolated  
840 multipotent mesenchymal progenitors provides new insight into the  
841 cellular identity of mesenchymal stem cells in mouse bone marrow. *Mol  
842 Cell Biol* **33**, 661-677, doi:10.1128/MCB.01287-12 (2013).

843 26 Shan, T., Liu, W. & Kuang, S. Fatty acid binding protein 4 expression marks  
844 a population of adipocyte progenitors in white and brown adipose tissues.  
845 *FASEB J* **27**, 277-287, doi:10.1096/fj.12-211516 (2013).

846 27 Wang, C. Q., Jacob, B., Nah, G. S. & Osato, M. Runx family genes, niche, and  
847 stem cell quiescence. *Blood Cells Mol Dis* **44**, 275-286,  
848 doi:10.1016/j.bcmd.2010.01.006 (2010).

849 28 Jo, A. *et al.* The versatile functions of Sox9 in development, stem cells, and  
850 human diseases. *Genes Dis* **1**, 149-161, doi:10.1016/j.gendis.2014.09.004  
851 (2014).

852 29 Wang, Y. *et al.* Lgr4 regulates mammary gland development and stem cell  
853 activity through the pluripotency transcription factor Sox2. *Stem Cells* **31**,  
854 1921-1931, doi:10.1002/stem.1438 (2013).

- 855 30 Liu, Y., Beyer, A. & Aebersold, R. On the Dependency of Cellular Protein  
856 Levels on mRNA Abundance. *Cell* **165**, 535-550,  
857 doi:10.1016/j.cell.2016.03.014 (2016).
- 858 31 Rodilla, V. *et al.* Luminal progenitors restrict their lineage potential during  
859 mammary gland development. *PLoS Biol* **13**, e1002069,  
860 doi:10.1371/journal.pbio.1002069 (2015).
- 861 32 Wang, C., Christin, J. R., Oktay, M. H. & Guo, W. Lineage-Biased Stem Cells  
862 Maintain Estrogen-Receptor-Positive and -Negative Mouse Mammary  
863 Luminal Lineages. *Cell Rep* **18**, 2825-2835,  
864 doi:10.1016/j.celrep.2017.02.071 (2017).
- 865 33 Howard, B. A. In the beginning: the establishment of the mammary lineage  
866 during embryogenesis. *Seminars in cell & developmental biology* **23**, 574-  
867 582, doi:10.1016/j.semcd.2012.03.011 (2012).
- 868 34 Liu, S. *et al.* Breast cancer stem cells transition between epithelial and  
869 mesenchymal states reflective of their normal counterparts. *Stem Cell*  
870 *Reports* **2**, 78-91, doi:10.1016/j.stemcr.2013.11.009 (2014).
- 871 35 Wahl, G. M. & Spike, B. T. Cell state plasticity, stem cells, EMT, and the  
872 generation of intra-tumoral heterogeneity. *NPJ Breast Cancer* **3**, 14,  
873 doi:10.1038/s41523-017-0012-z (2017).
- 874 36 Turlier, H. & Maitre, J. L. Mechanics of tissue compaction. *Semin Cell Dev*  
875 *Biol* **47-48**, 110-117, doi:10.1016/j.semcd.2015.08.001 (2015).
- 876 37 Malhotra, G. K. *et al.* The role of Sox9 in mouse mammary gland  
877 development and maintenance of mammary stem and luminal progenitor  
878 cells. *BMC Dev Biol* **14**, 47, doi:10.1186/s12861-014-0047-4 (2014).
- 879 38 Jeselsohn, R. *et al.* Embryonic transcription factor SOX9 drives breast  
880 cancer endocrine resistance. *Proc Natl Acad Sci U S A* **114**, E4482-E4491,  
881 doi:10.1073/pnas.1620993114 (2017).
- 882 39 Varga, J. & Greten, F. R. Cell plasticity in epithelial homeostasis and  
883 tumorigenesis. *Nat Cell Biol* **19**, 1133-1141, doi:10.1038/ncb3611 (2017).
- 884 40 Hong, T. *et al.* An Ovol2-Zeb1 Mutual Inhibitory Circuit Governs  
885 Bidirectional and Multi-step Transition between Epithelial and  
886 Mesenchymal States. *PLoS Comput Biol* **11**, e1004569,  
887 doi:10.1371/journal.pcbi.1004569 (2015).
- 888 41 Nieto, M. A., Huang, R. Y., Jackson, R. A. & Thiery, J. P. Emt: 2016. *Cell* **166**,  
889 21-45, doi:10.1016/j.cell.2016.06.028 (2016).
- 890 42 Revenu, C. & Gilmour, D. EMT 2.0: shaping epithelia through collective  
891 migration. *Curr Opin Genet Dev* **19**, 338-342,  
892 doi:10.1016/j.gde.2009.04.007 (2009).
- 893 43 Malladi, S. *et al.* Metastatic Latency and Immune Evasion through  
894 Autocrine Inhibition of WNT. *Cell* **165**, 45-60,  
895 doi:10.1016/j.cell.2016.02.025 (2016).
- 896 44 Malta, T. M. *et al.* Machine Learning Identifies Stemness Features  
897 Associated with Oncogenic Dedifferentiation. *Cell* **173**, 338-354 e315,  
898 doi:10.1016/j.cell.2018.03.034 (2018).
- 899 45 Smalley, M. J. Isolation, culture and analysis of mouse mammary epithelial  
900 cells. *Methods Mol Biol* **633**, 139-170, doi:10.1007/978-1-59745-019-  
901 5\_11 (2010).
- 902 46 Sun, L., Lee, M. Y. & Veltmaat, J. M. A non-enzymatic microsurgical  
903 dissection technique of mouse embryonic tissues for gene expression

904 profiling applications. *Int J Dev Biol* **55**, 969-974,  
905 doi:10.1387/ijdb.1134241s (2011).  
906 47 Arnaoutova, I. & Kleinman, H. K. In vitro angiogenesis: endothelial cell  
907 tube formation on gelled basement membrane extract. *Nat Protoc* **5**, 628-  
908 635, doi:10.1038/nprot.2010.6 (2010).  
909 48 Schindelin, J. *et al.* Fiji: an open-source platform for biological-image  
910 analysis. *Nat Methods* **9**, 676-682, doi:10.1038/nmeth.2019 (2012).  
911 49 Morrison, B. & Cutler, M. L. Mouse Mammary Epithelial Cells form  
912 Mammospheres During Lactogenic Differentiation. *J Vis Exp*,  
913 doi:10.3791/1265 (2009).  
914 50 Boras-Granic, K., Dann, P. & Wysolmerski, J. J. Embryonic cells contribute  
915 directly to the quiescent stem cell population in the adult mouse  
916 mammary gland. *Breast Cancer Res* **16**, 487, doi:10.1186/s13058-014-  
917 0487-6 (2014).  
918

919 **Figure legends**

920

921 **Figure 1. Establishment of novel embryonic mammary progenitor cell**  
922 **lines.**

923 (A) Schematic illustration of experimental scheme of eMPC derivation from  
924 E12.0-stage Immorto;s-SHIP-GFP mammary organ number three. After  
925 microdissection of MP3 to include three tissues: GFP+ MPE (GFP staining in  
926 green), GFP-, ER $\alpha$ + MM (ER $\alpha$  staining in magenta) and adjacent ER $\alpha$ - FPP  
927 (DAPI staining in blue), MP3 was cultured on thick BME for 4 weeks.

928 (B) Brightfield and GFP images of cultured MP3 at Day 1, Day 11, Day 25.  
929 After culture, a variety of cell types including lipid-containing adipocyte-like  
930 cells (area shown in black box and yellow arrow), as well as contractile  
931 muscle-like cells (black arrowhead) were observed within the cell population.  
932 After enzymatic dissociation, eMPC were plated and expanded in 2D culture  
933 to obtain a pool of eMPCs (ePool). eMPC cells were subject to single cell  
934 sorting into GFP+ and GFP- cell populations using fluorescence-activated cell  
935 sorting. Clones were expanded as single-cell derived clones to create the  
936 eighteen single cell-derived clones from GFP- cells (e1 and e2) or GFP+ cells  
937 (eG1, eG2E9, etc.) described in this study. scale bar, 200  $\mu$ m.

938 MP3, mammary primordium 3, MPE, mammary primordial epithelium, Epi,  
939 epithelium, MM, mammary mesenchyme, FPP, fat pad precursor. eMPC,  
940 embryonic mammary progenitor cells (eMPC), BF, brightfield.

941

942 **Figure 2. Global gene expression landscape of novel embryonic**  
943 **mammary progenitor cell lines.**

944 (A) Dendrograms from hierarchical clustering of RNA sequencing data using  
945 2059 most significantly-changing genes of eMPC and other types of  
946 progenitor/stem cells. Eighteen eMPC clones are classified into three main  
947 clusters, which are highlighted: e1 cluster (yellow box), e2/eG2/ePool cluster  
948 (magenta box), and eG1 cluster (green box).

949 (B) PCA of eMPC and other types of progenitor/stem cells. Loading plot using  
950 principal component (PC) 2 (18%) and PC3 (14%) of the 17,204 most variable  
951 genes. Three main eMPC clusters (yellow, magenta, and green) recapitulate  
952 the hierarchical clustering results shown in (B). The heatmap shows  
953 hierarchical clustering based on top 0.5% of PC2 (92 genes).

954 (C) PC1 and PC2 loading plot using the most highly significant genes in the  
955 dataset of eMPC and other progenitor/stem cells combined with embryonic  
956 mammary tissues. Three eMPC clusters (yellow, magenta, and green) and an  
957 embryonic mammary tissue cluster (purple) from PCA are highlighted based  
958 PC1. The heatmap shows results of hierarchical clustering based on 26 genes  
959 from PC1 (top high and low rotation genes analysed using the Likelihood ratio  
960 test (LRT) as well as the intensity difference test).

961 eMPC, embryonic mammary progenitor cell; MP3, mammary primordium 3;  
962 BME, basement membrane extract; PC, principal component, MPE,  
963 mammary primordial epithelium, Epi, epithelium, MM, mammary  
964 mesenchyme, FPP, future fat pad precursor. ESC, embryonic stem cells,  
965 MEF, mouse embryonic fibroblasts, MSC, mesenchymal stem cells and  
966 BMMC, bone marrow mononuclear cells. Early and late indicate clones that  
967 were sequenced at early (5th) and later (13th) passage.

968

969 **Figure 3. Progenitor and lineage-associated marker expression in select**  
970 **eMPCs.**

971 (A) Heatmap depicting stem cell and epithelial marker expression in four  
972 eMPC selected for further characterisation.

973 (B) qRT-PCR for *Aldh1a1*, *Krt19*, *Krt18*, *Sox9*, *Snai2*, *Twist1*, *Procr*, *Trp63*,  
974 *Fabp4* in eMPC clones. Relative quantification (RQ) of each sample is  
975 normalised to that of ePool and shown in log2 plot with error bar of triplicates.  
976 ( $n = 3$ , mean  $\pm$  s.e.m.).

977 (C) Immunohistochemistry staining of Sox9 and Twist expression in eMPC.  
978 Scale bar, 200  $\mu$ m.

979 eMPC, embryonic mammary progenitor cell; MEF, mouse embryonic  
980 fibroblast; MSC, mesenchymal stem cell; ESC, embryonic stem cell.

981

982 **Figure 4. eMPCs harbour varying potential for differentiation into**  
983 **mesodermal lineage.**

984 (A) *In vitro* differentiation of eMPC clones and positive control (MSC) with  
985 adipogenic stimuli. Neutral lipid staining merged with bright-field image shows  
986 accumulation of lipid droplets (red) within cells after 6 days of induction. Scale  
987 bar, 100  $\mu\text{m}$ .

988 (B) Vasculogenesis assay results of eMPCs and positive control (HUVEC).  
989 Tubes formed in medium with growth factors (EGM) and without growth  
990 factors were stained with Calcein AM (white) after 24 hours incubation. The  
991 total area of networks (blue, yellow, and green lines) was analysed and  
992 presented in box plots with whiskers denoting minimum and maximum values  
993 ( $n = 4$  mean  $\pm$  s.e.m.). Statistical significance was computed using one-way  
994 ANOVA and Dunnett's multiple comparisons test as \*\*\*\* $P \leq 0.0001$ ,  $P^{***} \leq$   
995  $0.001$ , ns = no significance. Two-tailed  $P$  values are ePool EGM 0.0001, e1  
996 EGM 0.0001, e2 EGM 0.1607, eG1 EGM 0.9994, eG2 EGM 0.0005, MSC  
997 EGM 0.9999, and HUVEC EGM 0.0001, when compared to control HUVEC  
998 EBM. Scale bar, 200  $\mu\text{m}$ .

999 eMPC, embryonic mammary progenitor cell; MSC, mesenchymal stem cell;  
1000 HUVEC, Human umbilical vein endothelial cells; EGM, endothelial cell growth  
1001 media EBM, endothelial basal growth media.

1002

1003 **Figure 5. eMPC lines are clonogenic and display distinct acinar**  
1004 **morphologies when grown as mammospheres in 3-D culture.**

1005 (A) Bright-field images of single cell-derived spheres. Scale, 40  $\mu\text{m}$ .

1006 (B) Colony-formation ability presented in box plots with whiskers denoting  
1007 minimum and maximum values showing number of cells (out of 10,000 cells  
1008 plated) that gave rise to spheres in anchorage-independent cultures.



1009 Statistical significance was computed for each cell line compared to control  
1010 (MSC) using one-way ANOVA and Dunnett's multiple comparisons test,  
1011 where \*\*\*\* $P \leq 0.0001$ , \*\*  $P \leq 0.01$ , \* $P \leq 0.05$ , ns = not significant. Two-tailed  $P$   
1012 values are ePool 0.0152, e1 0.0001, e2 0.3582, eG1 0.0092, and eG2 0.1107  
1013 when compared to control MSC. ( $n = 3$  mean  $\pm$  s.e.m.).

1014 (C) Compactness of individual spheres presented where a circle with Form  
1015 factor = 1 is most condensed. eG1 forms more compact spheres compared to  
1016 those of e1, which is also evident in images in (A). Statistical significance was  
1017 computed using unpaired t-test, where two-tailed \*\*\*\* $P = 0.0001$ . Bar  
1018 indicates mean value  $\pm$  s.e.m.  $n = 3$ .

1019 (D) Images of e1 and eG1 spheres stained by immunofluorescence with  
1020 CD44 (red), an adhesion receptor expressed by basal mammary cells which  
1021 mediates epithelial-stromal and cell-cell interactions and DAPI (blue).

1022 Scale bar, 50  $\mu$ m.

1023 eMPC, embryonic mammary progenitor cell; MSC, mesenchymal stem cell.

1024

## 1025 **Figure 6. eMPCs have lactogenic potential.**

1026 (A) Alveologensis assay results for eMPC clones and positive control (MEC).  
1027 Formation of alveolar-like structures and milk production was assessed by  
1028 morphological changes using bright-field microscopy, as well as IF using anti-  
1029 milk antibody (red). Scale bar, 40  $\mu$ m.

1030 (B) Quantification of number of alveolar-like structures greater than 100  $\mu$ m.  
1031 The plot is shown as count per 24 well. ( $n = 3$ , mean  $\pm$  s.e.m.).

1032 eMPC, embryonic mammary progenitor cell; MEC (mammary epithelial cell).

1033

## 1034 **Figure 7. Loss of Sox9 in e1 enhances colony formation ability and** 1035 **alters sphere morphology.**

1036 (A) Schematic illustration for establishment of Sox9-targeted subclones from  
1037 Cas9-expressing e1.

1038 (B) qRT-PCR for *Sox9* and EMT regulator, *Zeb1*. Each RQ is normalised to  
1039 Co#1. Statistical significance between control group and *Sox9* knockout group  
1040 was computed using unpaired t-test.  $**P \leq 0.01$ ,  $*P \leq 0.05$ . Two-tailed *P*  
1041 values are 0.0306 for *Sox9* and 0.0036 for *Zeb1*, respectively ( $n = 3$ , mean  
1042 with 95% confidence interval).

1043 (C) Immunohistochemistry for *Sox9* and *Zeb1* showing *Zeb1* expression is  
1044 decreased by *Sox9* reduction. Scale, 200  $\mu\text{m}$ .

1045 (D) Representative images of spheres grown in methylcellulose (left) and  
1046 colony-formation ability presented in box plots (right) with minimum and  
1047 maximum values showing number of control or *Sox9*-targeted cells (out of  
1048 10,000 cells plated) that gave rise to spheres in anchorage-independent  
1049 cultures. Statistical significance was calculated using unpaired t-test. Two-  
1050 tailed  $****P = 0.0001$ . ( $n = 3$ , mean  $\pm$  s.e.m.). Scale bar , 100  $\mu\text{m}$ .

1051 (E) Representative images of spheroids grown in BME from control and *Sox9*-  
1052 targeted cells. Scale bar, 400  $\mu\text{m}$ .

1053 (F) Quantification of area of spheroids and (G) increase in area of protrusions  
1054 from spheroids grown in BME from e1/control and e1*Sox9*-KO cells. ( $n = 8$ ,  
1055 mean + S.D.) Statistical significance was calculated using one-way ANOVA  
1056 and multiple comparisons,  $***P = 0.0002$  in F and G.

1057 e1, embryonic mammary progenitor cell 1; Co, control, non-targeted cells; KO,  
1058 knockout; *Sox9*-targeted cells; MEC, mammary epithelial cell.

1059

1060 **Figure 8. Effect of *Sox9* ablation on embryonic mammary progenitor cell**  
1061 **fate and function.**

1062 (A) qRT-PCR for *Sox9*, *Zeb1*, *Snai2*, *Twist*, *Procr*, *Trp63*, *Acta2*, and *Ecad* in  
1063 e1/Co#1 and e1/*Sox9*-KO#1 cells. Relative quantification of each sample is  
1064 normalised to that of e1/Co#1 and shown in Log2 plot with error bar of  
1065 triplicates. ( $n = 3$ , mean  $\pm$  s.e.m.).

1066 (B) Heat map from RNA-seq of e1/control and e1/Sox9-KO cells of genes that  
1067 were significantly up- or down-regulated in all three e1/Sox9-KO subclones  
1068 using both DESeq and Intensity Difference test.

1069 (C) *In vitro* alveologenesis assay results. Formation of alveoli-like structure  
1070 and milk production was assessed by morphological changes using bright-  
1071 field microscopy, as well as IF using anti-milk antibody (red) and DAPI. Scale  
1072 bar, 200  $\mu\text{m}$ .

1073 (D) *In vitro* alveologenesis assay results. Quantification of number of alveoli  
1074 less than 200  $\mu\text{m}$ , 200-400  $\mu\text{m}$ , and greater than 400  $\mu\text{m}$  per well. Statistical  
1075 significance between control group and Sox9 knockout group was computed  
1076 using two-tailed paired t-test.  $P = 0.0044$ . ( $n = 3$ ).

1077 e1, embryonic mammary progenitor cell; Co#1/#2, control #1/#2, non-targeted  
1078 e1 cells; KO#1/#2/#3, knockout #1, #2, #3; Sox9-targeted e1 cells.

1079  
1080  
1081

1082 **Acknowledgements**

1083 This work was funded by Programme Grants from Breast Cancer Now as part  
1084 of Programme Funding to the Breast Cancer Now Toby Robins Research  
1085 Centre. We thank the High-Throughput Genomics Group at the Wellcome  
1086 Trust Centre for Human Genetics (funded by Wellcome Trust grant reference  
1087 090532/Z/09/Z) for the generation of sequencing data and Dr Simon Andrews,  
1088 Babraham Insititute, for performing bioinformatics analysis.

1089

1090 **Conflict of interest**

1091 The authors declare no conflict of interest.

1092

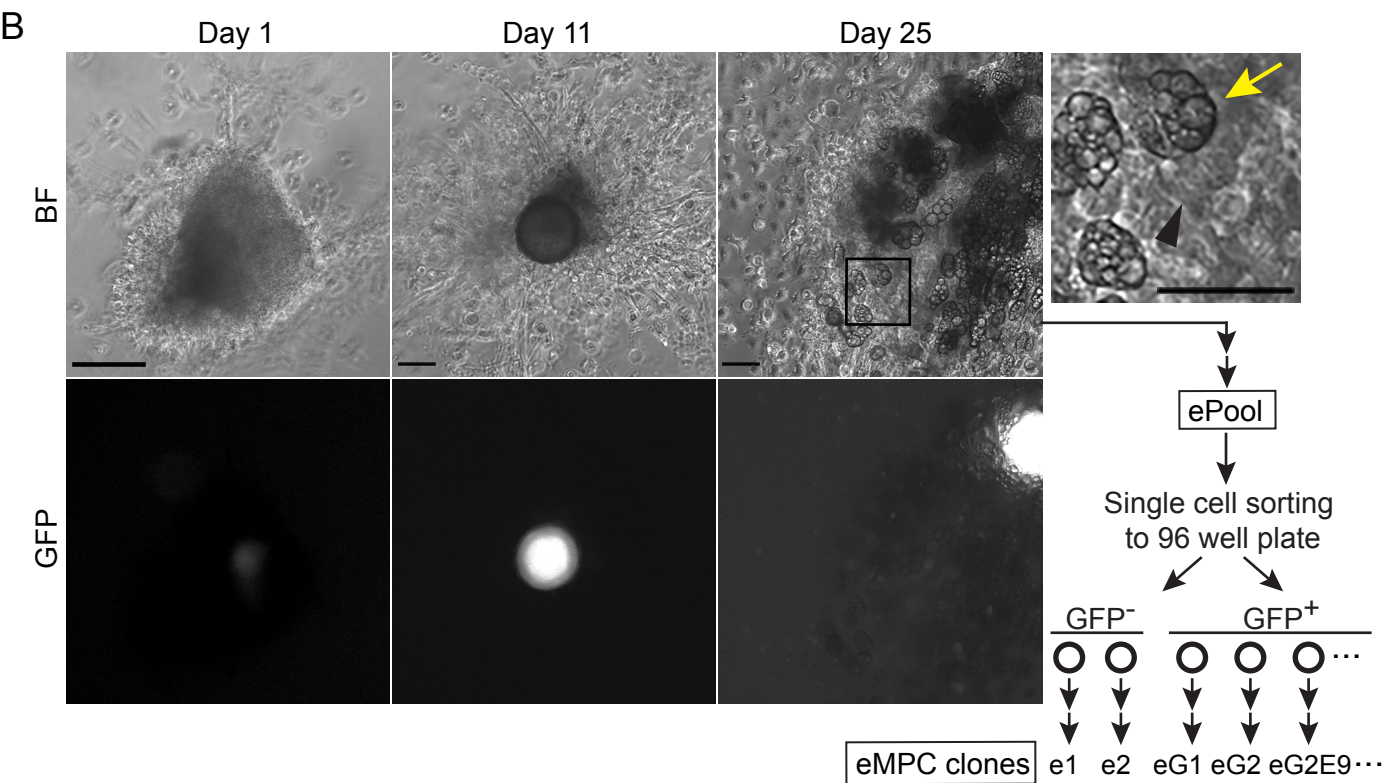
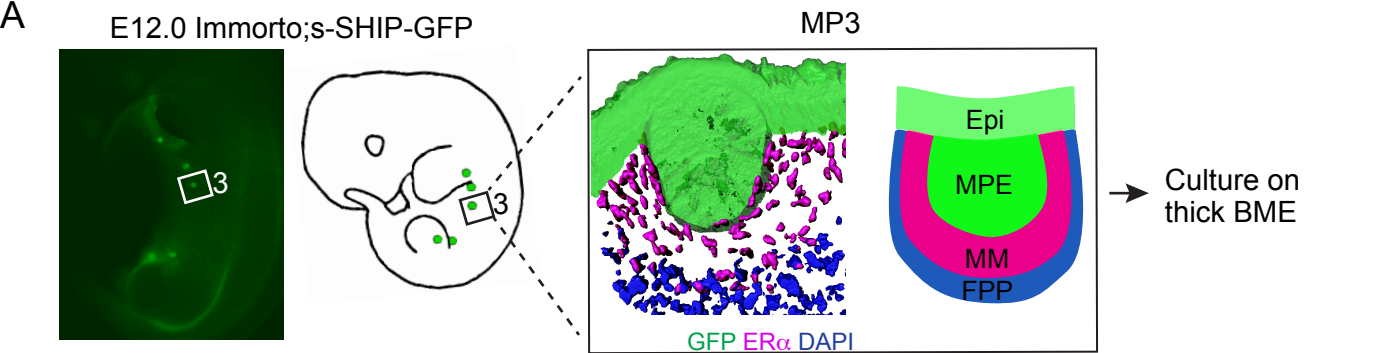
1093 **Author contributions**

1094 NK and BAH were responsible for study design, experimental work, data  
1095 analysis and manuscript preparation; NK, PB, MT, EO, AL performed the  
1096 experiments and analysed the data; NK manuscript editing, experimental  
1097 details; BAH wrote the manuscript; all authors commented on the manuscript.

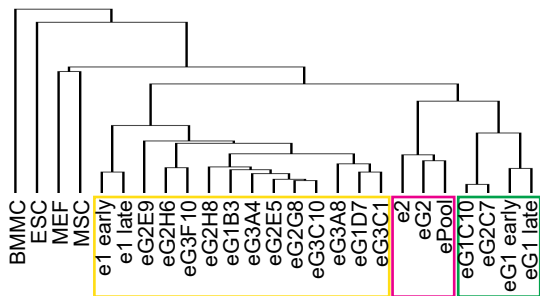
1098

1099

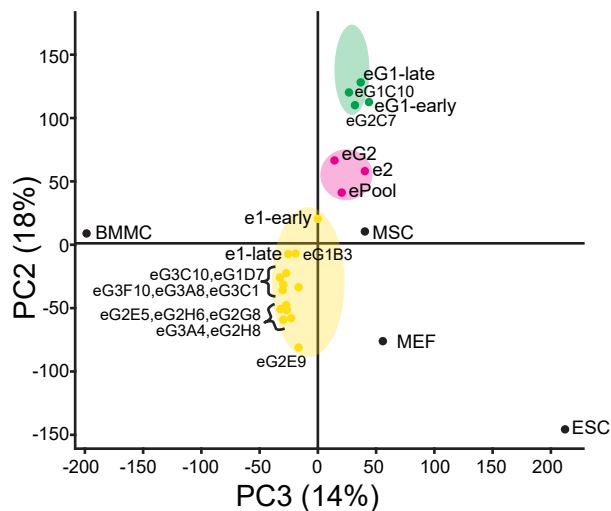
1100



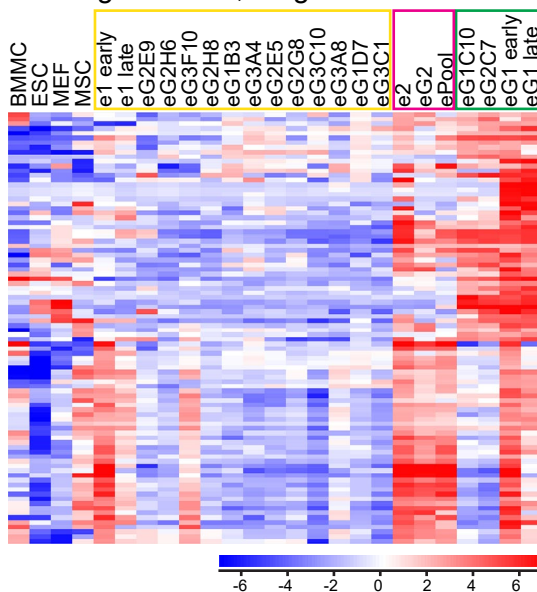
A



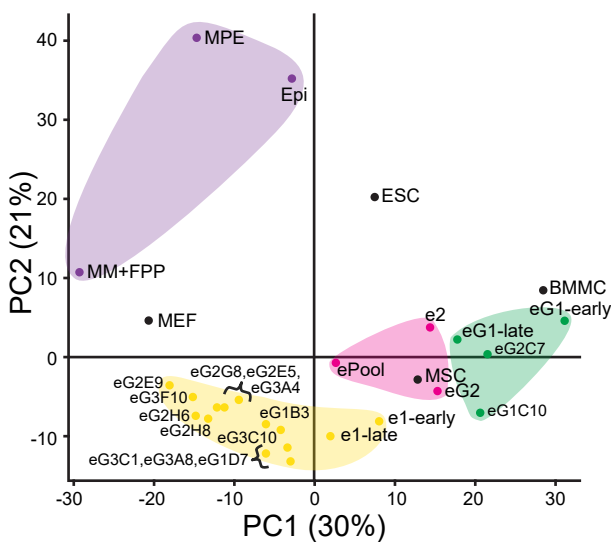
B



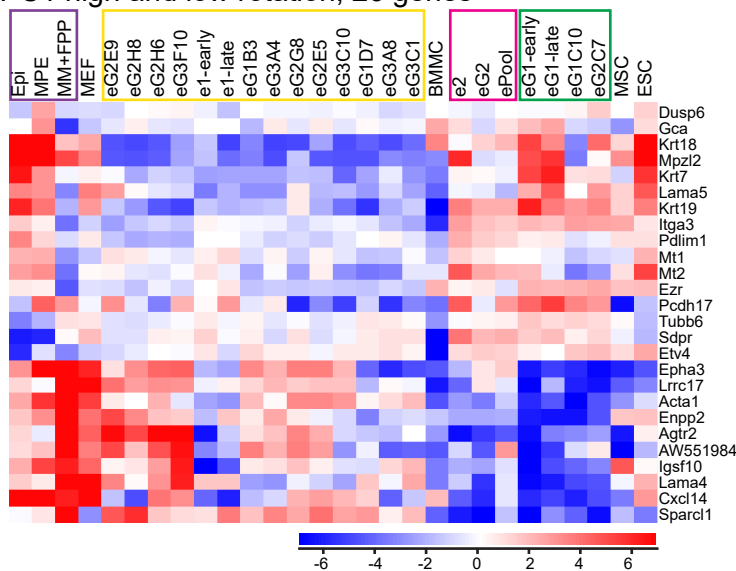
PC2 high rotation, 92 genes

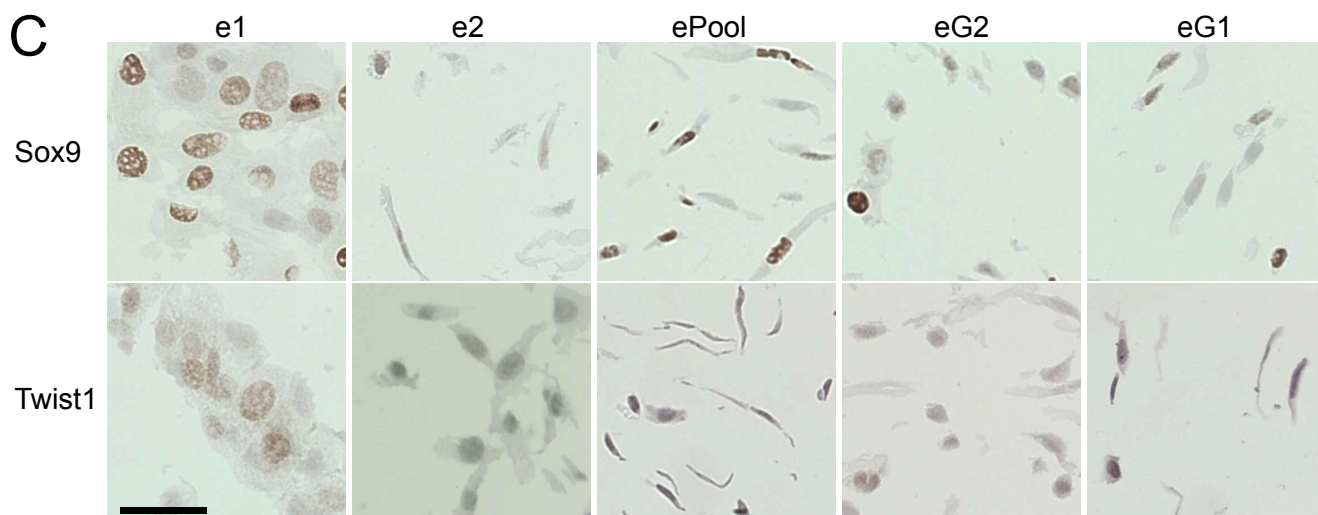
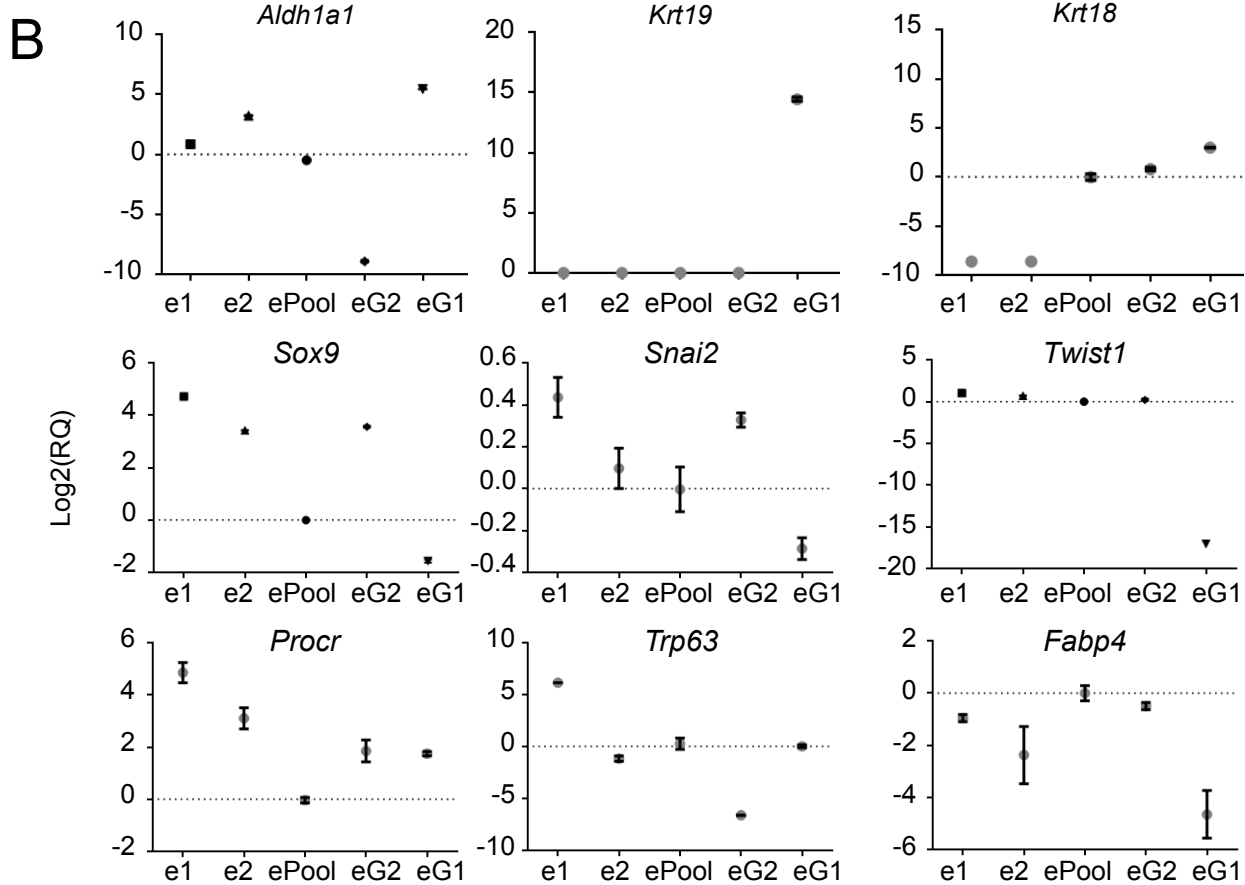
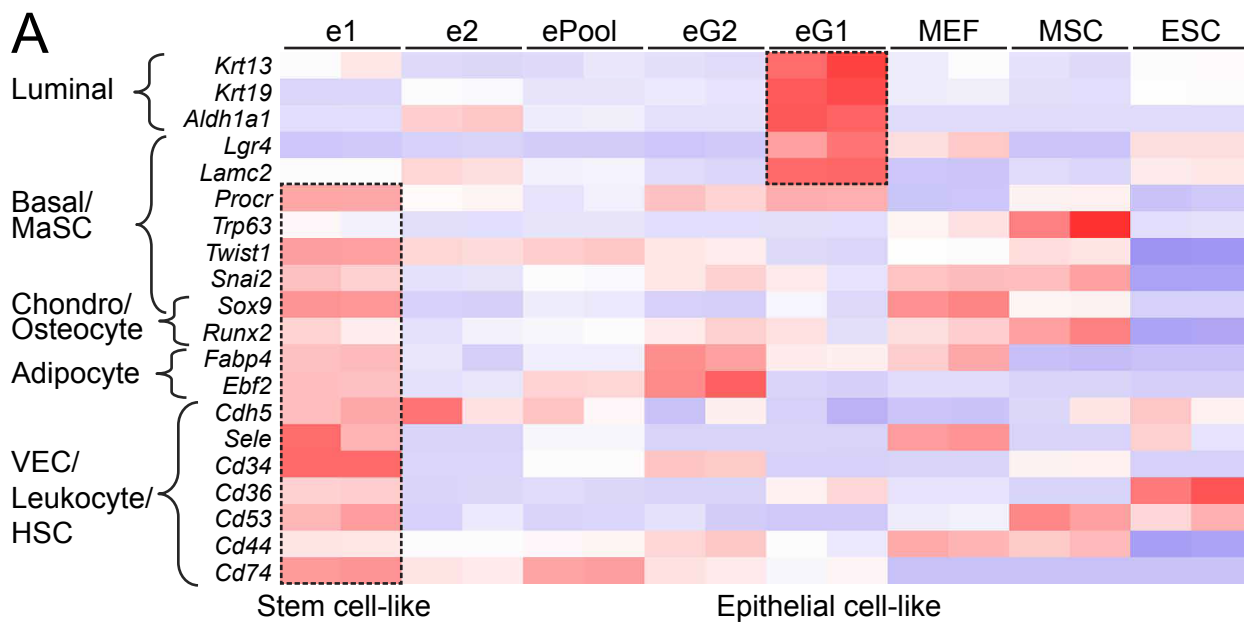


C

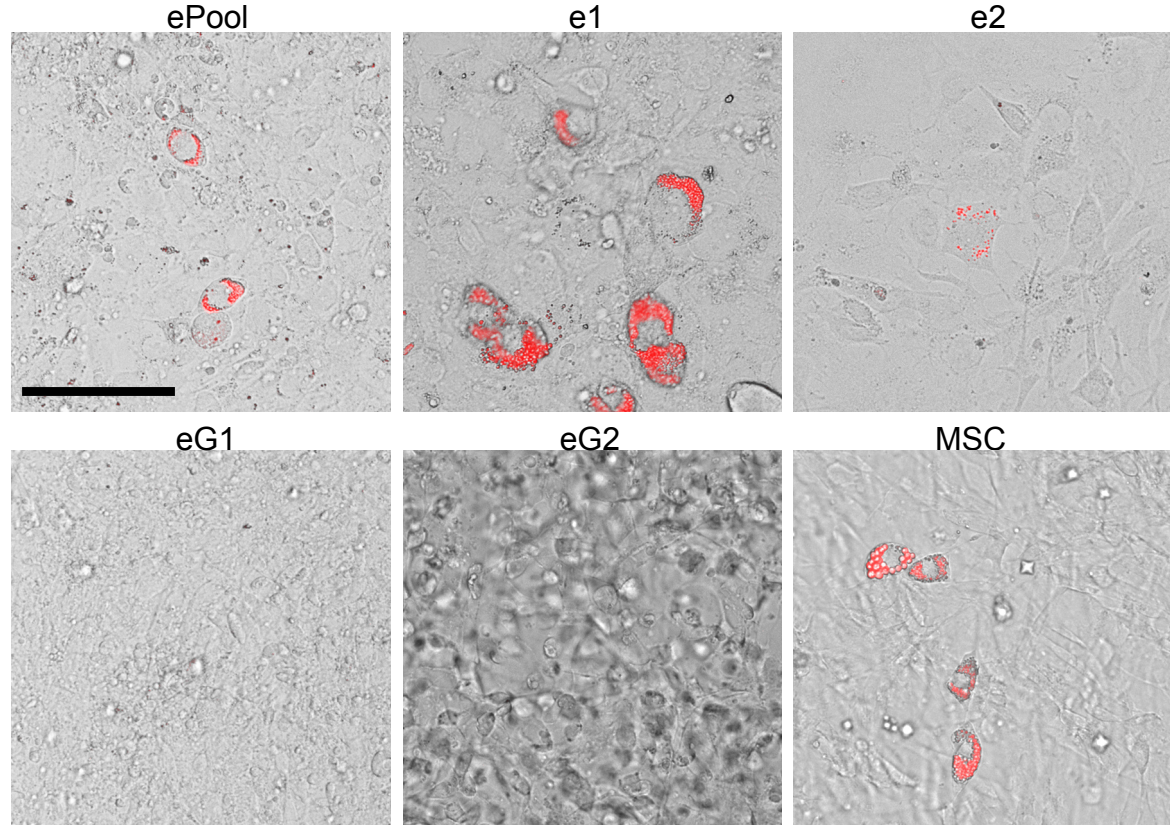


PC1 high and low rotation, 26 genes

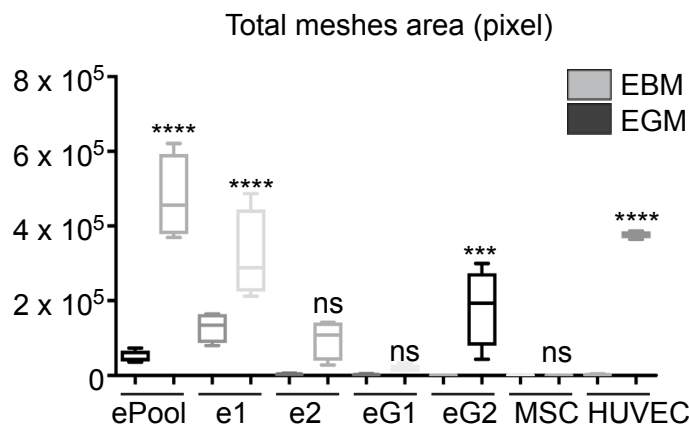
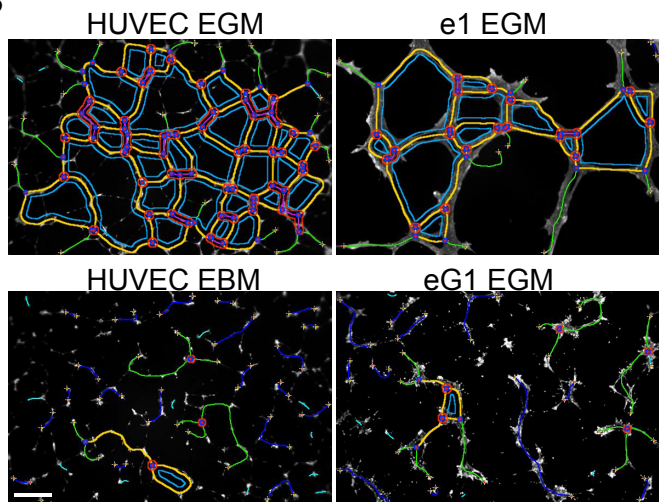




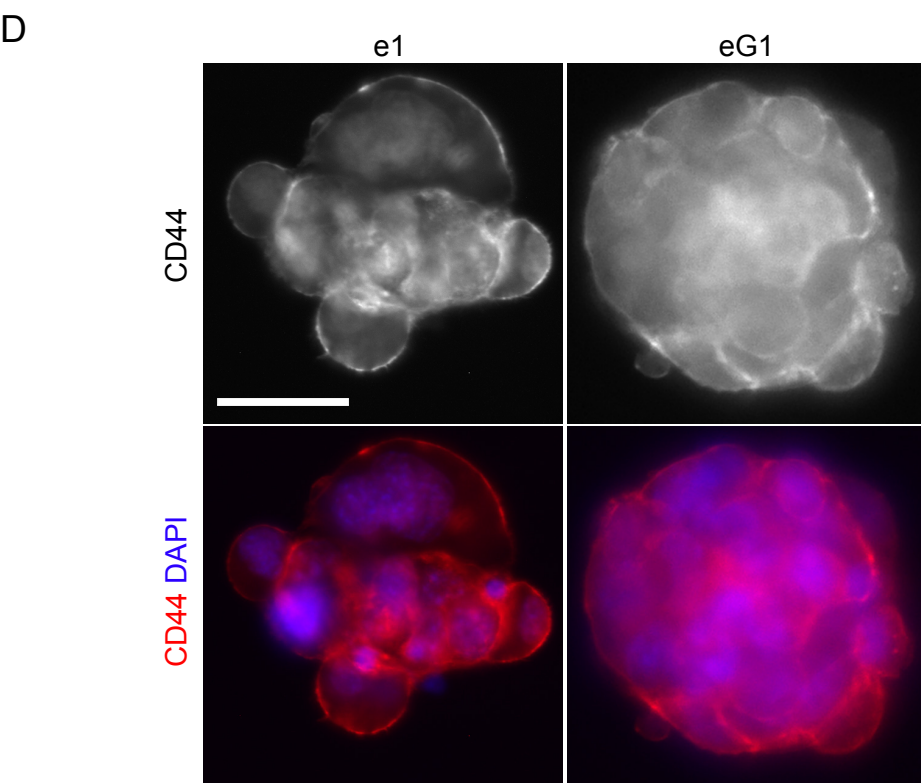
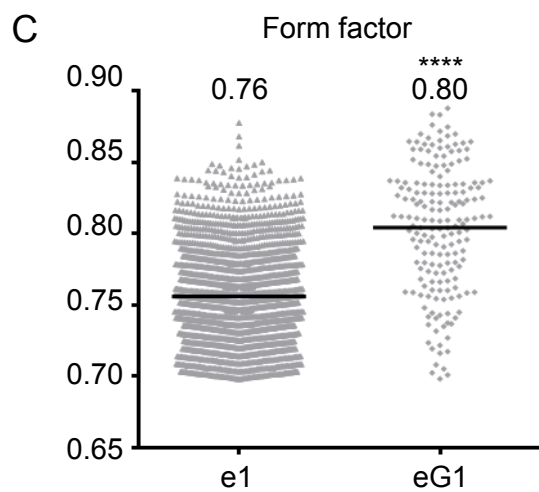
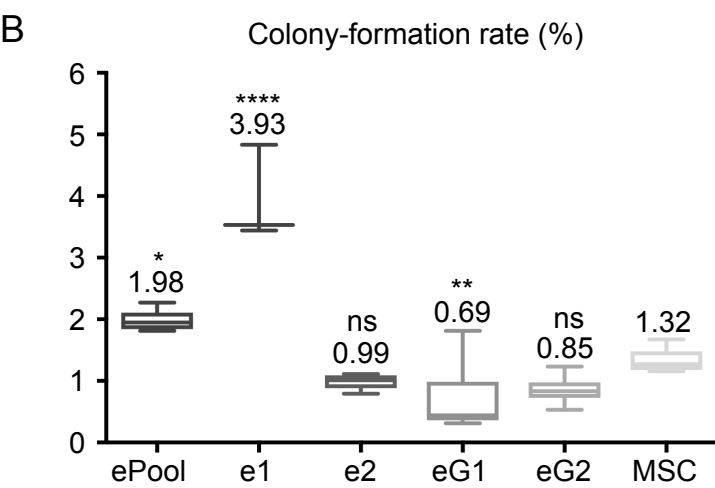
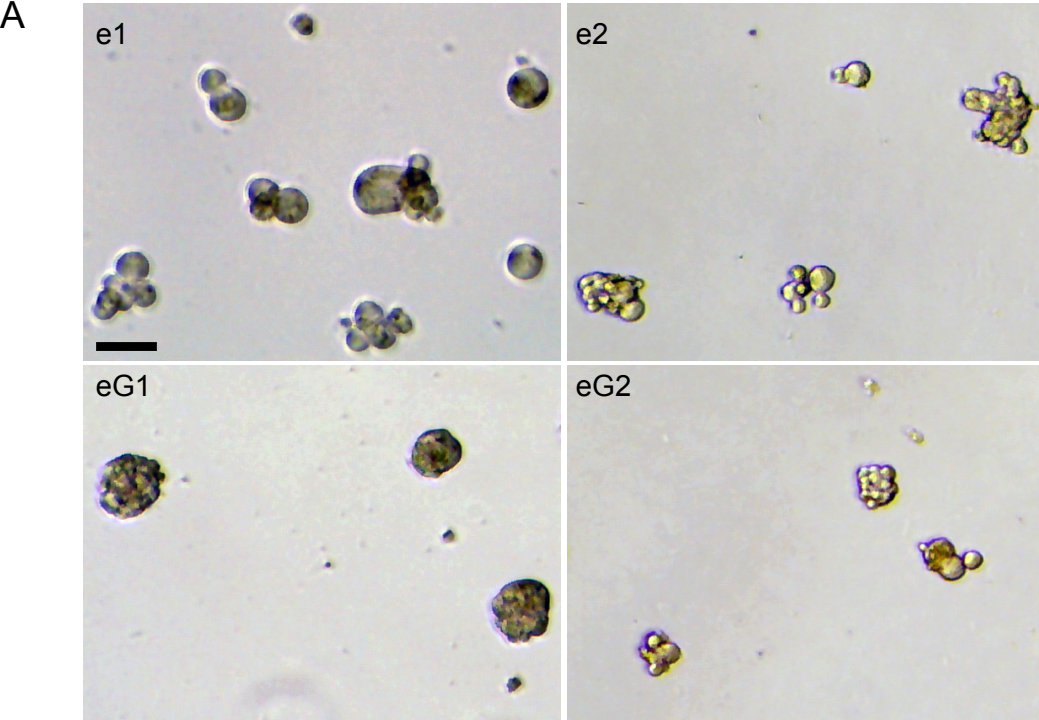
A

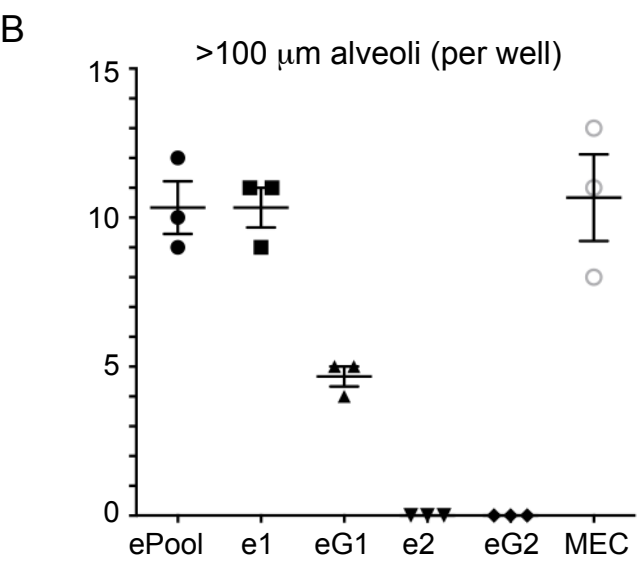
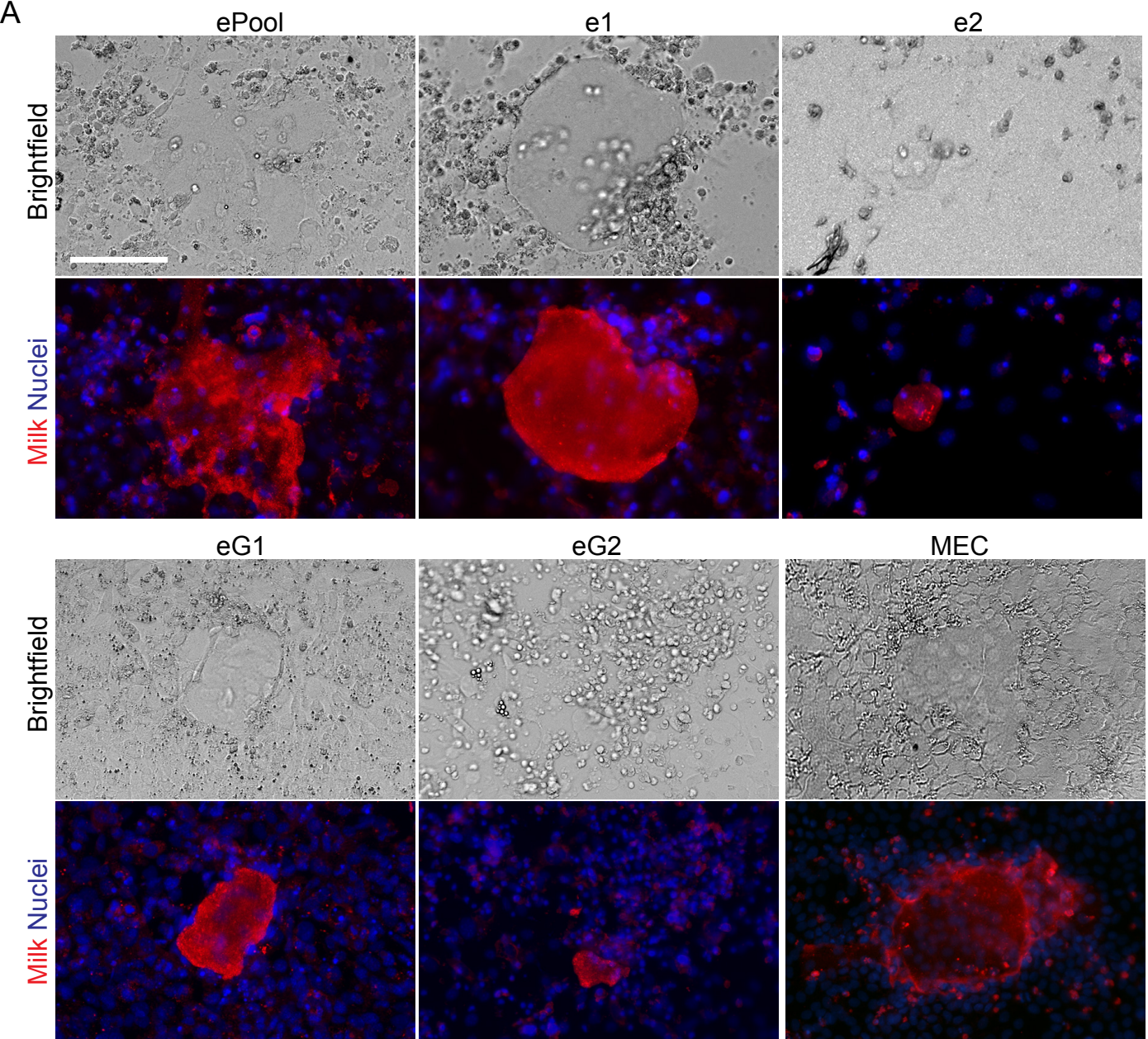


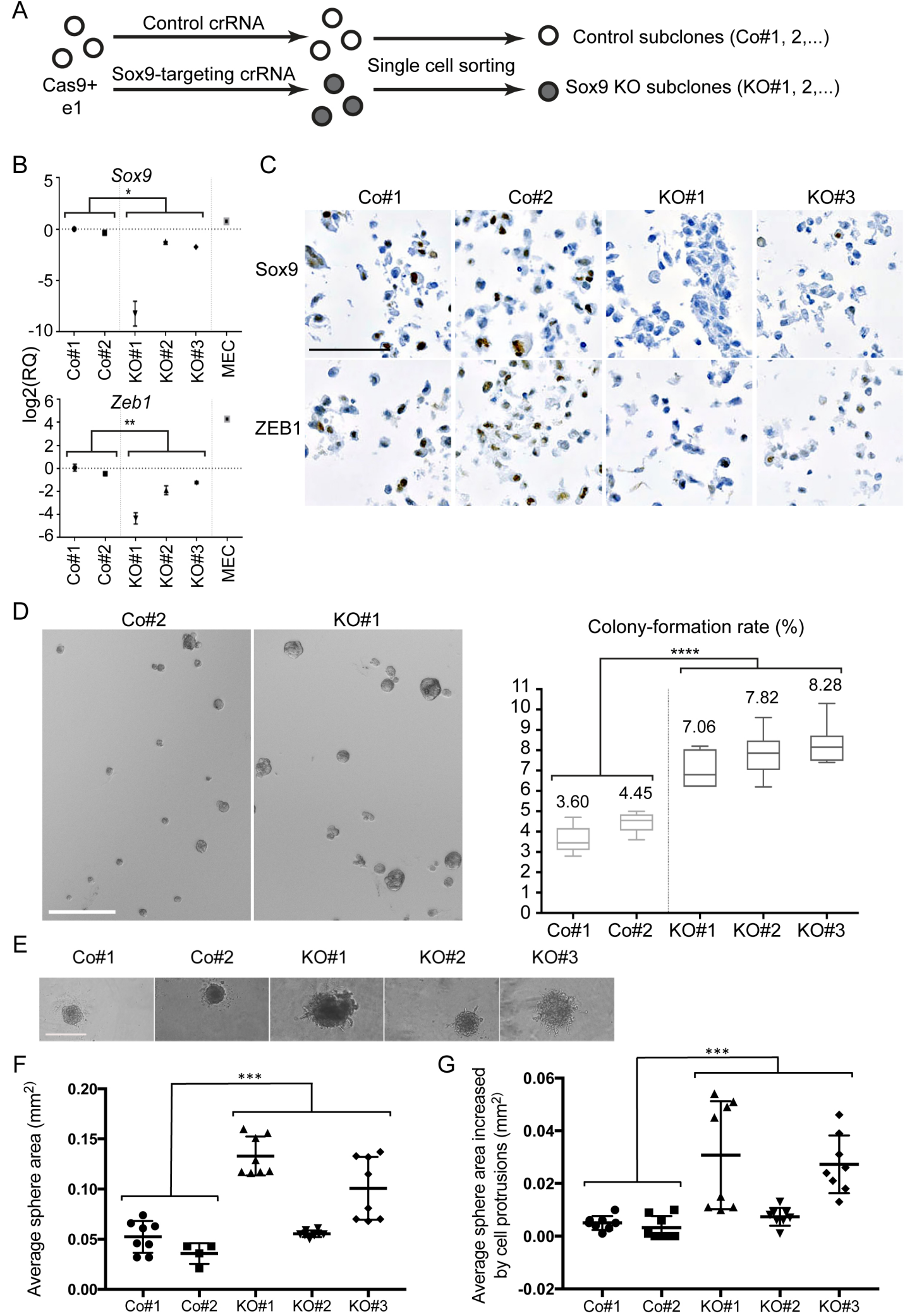
B



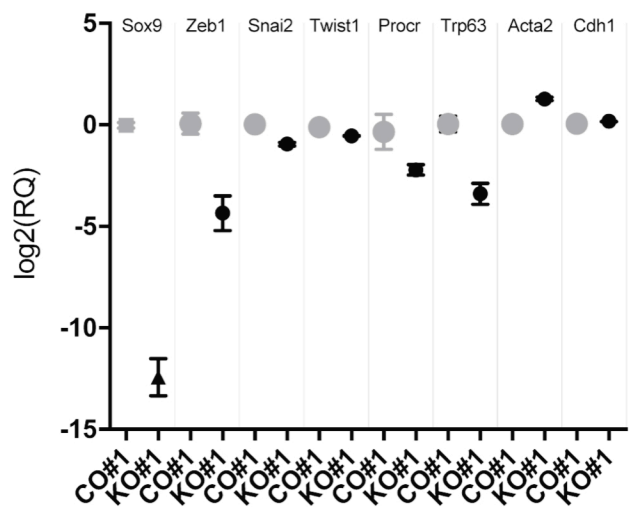




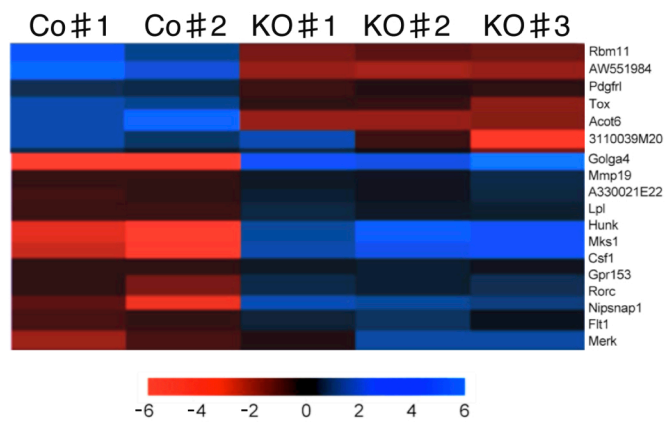




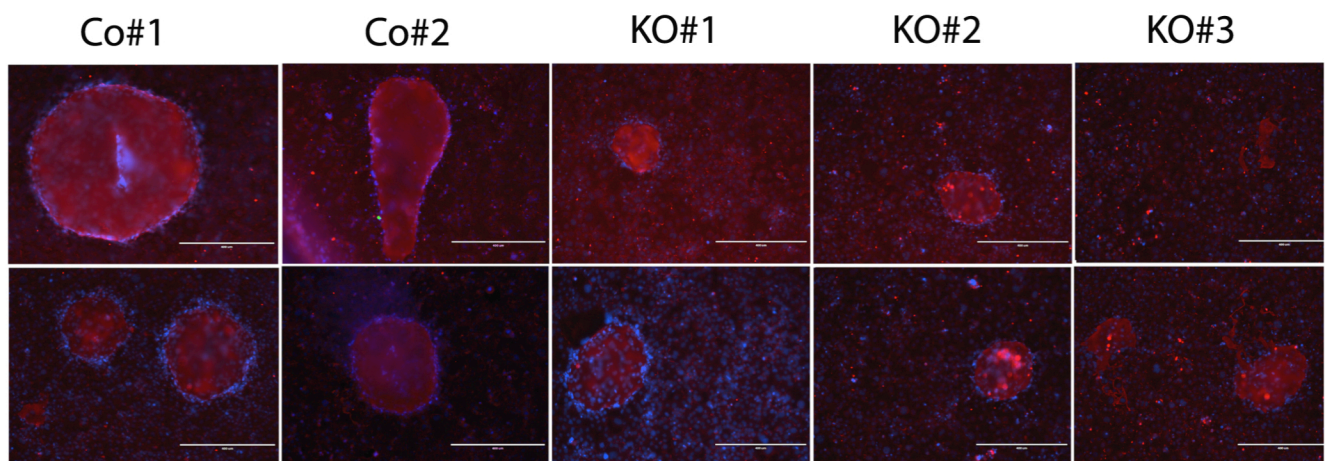
A



B



C



D

