# Integrated single-cell potency and expression landscape in mammary epithelium reveals novel bipotent-like cells associated with breast cancer risk

Andrew E. Teschendorff [1,2,*] , Samuel J Morabito [3] , Kai Kessenbrock [3] and Kerstin Meyer [4]

1. CAS Key Laboratory of Computational Biology, CAS-MPG Partner Institute for Computational Biology, Shanghai Institute of Nutrition and Health, Shanghai Institutes for Biological Sciences, University of Chinese Academy of Sciences, Chinese Academy of Sciences, 320 Yue Yang Road, Shanghai 200031, China.

2. UCL Cancer Institute, Paul O'Gorman Building, University College London, 72 Huntley Street, London WC1E 6BT, United Kingdom.

3. Chao Family Comprehensive Cancer Center, University of California, Irvine 839 Health Science Road, Sprague Hall 114 Irvine, CA 92697-3905, USA.

4. Wellcome Sanger Institute, Cambridge, UK.

*Corresponding author: Andrew E. Teschendorff- a.teschendorff@ucl.ac.uk , andrew@picb.ac.cn

**Running title:** Progenitor single cells and cancer risk

**Keywords:** scRNAseq; differentiation; stem-cell; cancer; Waddington landscape; entropy

## Abstract

22  **The identification of progenitor and stem like cells in epithelial tissues, as well as those**

23  **that may serve as the cell of origin for epithelial cancers, is an outstanding challenge.**

24

25  **Here we present a novel algorithm, called LandSCENT, which constructs a**

26  **3-dimensional integrated landscape of cell-states, encompassing cell-potency and**

27  **expression subtypes, to facilitate the identification of progenitor and stem-like cells.**

28  **Application to thousands of single-cell RNA-Seq profiles from the normal mammary**

29  **epithelium reveals a rare 5% subpopulation of highly potent single-cells. The integrated**

30  **landscape naturally predicts that these cells define a bi-potent-like state, a result not**

31  **obtainable via standard methods or without invoking prior assumptions. The**

32  **bi-potent-like cells are overrepresented within the basal compartment but also overlap**

33  **with an immature luminal phenotype. We characterize the transcriptome of these cells**

34  **and show that is enriched for a mammary stem-cell module. We further identify *YBX1*,**

35  **a regulator of breast cancer risk identified from GWAS, as the key transcription factor**

36  **defining this candidate bi-potent cellular phenotype. We validate the putative bi-potency**

37  **of *YBX1*-marked cells using independent FACS-sorted bulk expression data. In addition,**

38  ***YBX1* is overexpressed in basal breast cancer and correlates with clinical outcome. In**

39  **summary, we here provide a novel computational framework which may serve to**

40  **identify and prioritize candidate normal or cancer progenitor/stem-like single-cell**

41  **phenotypes, for subsequent functional studies.**

42

## Introduction

44  Single-cell RNA-sequencing (scRNA-Seq) studies are revolutionizing our understanding of

45  cellular development, helping us to elucidate the hierarchical organization of cell-types

46  within complex tissues (Patel et al. 2014; Trapnell et al. 2014; Treutlein et al. 2014; Scialdone

47  et al. 2016; Tirosh et al. 2016a; Tirosh et al. 2016b; Treutlein et al. 2016; Haber et al. 2017;

48  Regev et al. 2017; Rozenblatt-Rosen et al. 2017; Hon et al. 2018; Laurenti and Gottgens 2018;

49  Shepherd et al. 2018). In these studies, a common computational task is the clustering of

50  single cells, which may reveal novel cell-types within these tissues (Haber et al. 2017; Han et

51  al. 2018). Relations between known and novel cell-types can be subsequently derived using

52  lineage-trajectory inference type algorithms (Trapnell et al. 2014) (Chen et al. 2016) (Marco

53  et al. 2014) (Haghverdi et al. 2016) (Grun et al. 2016). However, assignment of single cells to

54  cell-types often requires prior knowledge of specific markers, which may inevitably introduce

55  bias (Stingl et al. 2006; Trapnell 2015; Yuan et al. 2017). In certain circumstances this bias

56  can be substantial, specially if knowledge of suitable markers is not available or at best

57  controversial (Costa et al. 2018; Grun 2018). Moreover, lineage-trajectory inference

58  algorithms, including recent state-of-the-art ones such as Monocle-2 (Qiu et al. 2017), often

59  require specification of a "root" cell or node, in order to give the trajectories a "temporal"
60  direction. In the absence of temporal data, the specification of this root node may rely on
61  existing biological knowledge and therefore equally subject to bias. Another related and key
62  problem is that cell-types are typically inferred as clusters of relatively high cell-density in a
63  two dimensional reduced space, a procedure which does not necessarily allow for the
64  identification of cellular states. For instance, how to identify novel progenitor or stem-like
65  states within a cell-type may not be possible using two-dimensional clustering alone since
66  potency/stemness may be defined by additional latent dimensions.

67  To address these outstanding challenges, we here present LandSCENT, a novel computational
68  framework that avoids the aforementioned biases, assigning each cell, not only to a specific
69  cell-type, but also to a specific potency state (Teschendorff and Enver 2017). LandSCENT
70  achieves this without the need for prior knowledge or assumptions. LandSCENT integrates
71  the inferred cell-types and potency states into a multi-layered single-cell landscape, where
72  cell-states are defined by clusters of single-cells within a potency state. This novel approach
73  allows cells to be placed into specific cellular-states, thus allowing novel cellular phenotypes
74  to be identified, for instance novel progenitor or stem-like states within complex epithelial
75  tissues.

76  We illustrate this strategy in the context of the breast epithelium, a tissue for which
77  scRNA-Seq encompassing over 25,000 single epithelial cells from 4 women, has recently
78  been generated using the 10X Genomics Chromium assay (Nguyen et al. 2018). We apply
79  LandSCENT to this data to construct an integrated potency and cell-type landscape at the
80  single-cell level. This landscape reveals a novel putative bi-potent progenitor like cell-state,
81  characterized by overexpression of *YBX1,* a recently discovered regulator of breast cancer
82  risk (Castro et al. 2016), a result which we would not have found had we used standard
83  state-of-the-art clustering methods. We further validate the bi-potent/stem-like nature of the
84  identified single-cells using orthogonal bulk expression data of mammosphere-derived
85  mammary stem cells. Our data support the view that the identified bi-potent cells expressing
86  *YBX1* may give rise to both basal and luminal progenitors, potentially marking the
87  cell-of-origin for basal breast cancer.

88

## Results

90  **Constructing an integrative landscape of cell-states in breast epithelium**

91  We posited that improved clustering of single-cells so as to reveal novel biology, would be
92  possible by integrating cellular state information with the cells' expression profiles when
93  performing the clustering itself. One important feature of a cell that informs on cell-state is its
94  differentiation potency and in previous studies we proposed and validated an *in-silico*
95  measure of single-cell potency, based on the concept of single cell signaling entropy (SCENT)
96  (Banerji et al. 2013; Teschendorff and Enver 2017), which we have further shown is more

97  robust than other proposed single-cell potency models (Grun et al. 2016; Guo et al. 2016; Shi
98  et al. 2018a). We stress that SCENT represents a marker-free systems-biology approach to the
99  quantification of a cell's potency, which has been demonstrated to be very robust, and which
100  is applicable also to bulk samples (Banerji et al. 2013; Teschendorff and Enver 2017; Shi et al.
101  2018b). This is important because the alternative approach, i.e. to use expression of surface
102  markers, is unlikely to capture the full biological complexity underlying cellular potency,
103  while also introducing potential bias. Thus, here we present LandSCENT, a novel extension
104  of SCENT that combines inference of cell potency with single-cell clustering to construct a
105  landscape of single-cell states: these single-cell states integrate the single-cell potency
106  estimates with the inferred cell-type clusters, providing a 3-dimensional landscape
107  representation (**Fig.1, Methods**). Here we applied LandSCENT to a 10X Genomics
108  Chromium assay profiling thousands of single-cells in the breast epithelium (Nguyen et al.
109  2018), in order to define the landscape of cellular states in this tissue (**Fig.1**).
110  First, we phenotypically characterized the single cells, by performing t-SNE (van der Maaten
111  2008) followed by density-based spatial clustering (Ester et al. 1996) on 3473 single
112  epithelial cells (after QC) from one individual and using a reduced subset of 4261 genes that
113  exhibited a significant average and variance in expression across all cells (**Methods**). This
114  revealed three main single-cell clusters (**Fig.2A**), in line with previous observations (Nguyen
115  et al. 2018). One of these clusters expressed high levels of *KRT14*, a well-known basal
116  marker, which was not expressed in the other two main clusters (**Fig.2B**). Instead, the other
117  two clusters expressed *KRT18*, a well-known luminal marker. Consistent with the report of
118  Nguyen et al (Nguyen et al. 2018), the two luminal clusters were distinguished by expression
119  of lactotransferin (*LTF*) and luminal differentiation markers (*GATA3/FOXA1*), as well as
120  hormone receptors (*ESR1/PGR*) (**Fig.2B**), suggesting that the higher *LTF*-expressing cluster
121  represents a more immature (alveolar) luminal phenotype.
122  Next, we applied our Signaling Entropy Rate (SR) measure from SCENT to estimate the
123  differentiation potency of each single cell. To broadly categorize different levels of inferred
124  potency, we applied a Gaussian mixture model to the logit-transformed potency estimates of
125  the 3473 single cells, revealing the existence of three main potency states (**Fig.2C-D,**
126  **Methods**). We observed that the highest potency state represented a minority population,
127  with approximately only 5% of single-cells falling into this putative progenitor or stem-like
128  state (**Fig.2D**).
129
**Validation of potency assignments**
131  Although signaling entropy has been extensively validated as a cell-potency measure
132  (Teschendorff and Enver 2017; Shi et al. 2018a), we sought additional validation of the
133  specific potency assignments in the current dataset. It is well known that *GATA3, FOXA1* and
134  *ESR1* are associated with a more differentiated luminal phenotype and therefore the
135  expectation would be that their expression levels should be higher in the luminal cells of

136  lowest potency. We were able to confirm this with high statistical significance (**Fig.3A**). We
137  also validated the potency assignments within the basal compartment. For instance, we
138  observed that expression of *KRT5* and *EGFR,* two well-known basal differentiation markers,
139  decreased in the basal cells of higher potency (**Fig.3B**). We note that all these negative
140  correlations were apparent only when we restricted to cells where the genes were expressed.
141  If all cells were included, including technical and biological dropouts, we did not observe
142  these genes to exhibit the expected negative correlation: in fact, they showed an opposite
143  trend due to a larger number of dropouts among low potency cells (**SI fig.S1**). To investigate
144  this further and to validate our method to call differential expression (DE), we used bulk
145  mRNA expression data from FACS sorted differentiated luminal and basal cells (Shehata et al.
146  2012) to define a gold-standard list of 5,773 differentially expressed genes between basal and
147  luminal cells. For each of these gold-standard genes, and using only cells expressing the
148  corresponding gene, we derived a t-statistic of differential expression between the single cell
149  basal and luminal clusters, which revealed that for the great majority of gold-standard genes
150  with sufficient single-cell data, these exhibited the expected pattern of differential expression
151  (OR=8.31, Fisher=test P=2e-26, **SI fig.S2**). Based on this, we conclude that performing DE
152  using only cells expressing the gene is a valid procedure, thus also validating our potency
153  assignments.
154
155  **Integrative landscape reveals a putative bi-potent cell state**
156  Having identified and validated the main single-cell clusters and potency states, we next
157  considered the distribution of potency states across these 3 clusters, as well as those cells not
158  assigned to any cluster ("peripheral cells"). Interestingly, cells in the high potency state were
159  found primarily within the basal compartment, but also mapped preferentially to the common
160  peripheral area of the three main clusters, and were therefore also relatively over represented
161  among peripheral cells (**Fig.2C**, **Fig.4A**). To assess this in more detail, we used LandSCENT
162  to create cell-density elevation maps of all cells, and separately also for all highly potent cells,
163  within the two-dimensional t-SNE landscape, which confirmed that the maximum density of
164  the highly potent cells defined a peak within the basal cluster, but with a ridge connecting it
165  to another peak within the immature luminal (L1) cluster (**Fig.4B**), suggestive of a bi-potent
166  cell population. In line with this, we observed that among all cells categorized into the high
167  potency (PS3) state, those falling within this peak also exhibited the highest levels of
168  signaling entropy (**SI fig.S3**). To exclude the possibility that these higher or bi-potent cells
169  may be doublets, we estimated doublet scores for all cells using a novel simulation approach
170  (Dahlin et al. 2018). In line with the expected doublet rate for 10X technology, this analysis
171  revealed that approximately 2% of the assayed cells are potential doublets (**SI fig.S4A**). As
172  expected, most of these mapped to the peripheral area between the major luminal and basal
173  clusters, yet they clearly also did not overlap with the most highly potent cells within the
174  basal and luminal clusters, confirming that our candidate bi-potent cells are generally not

175    doublets (**SI fig.4B**). Supporting this, we observed that the relation between signaling entropy
176    and doublet scores is a non-linear one, with many highly potent cells not necessarily having
177    high doublet scores (**SI fig.4C**). Finally, we verified that similar results were obtained had we
178    used another method for estimating doublet scores (**SI fig.S5, Methods**).
179

180    **Bipotent cells are marked by *YBX1* and *ENO1* overexpression**
181    In order to characterize the highly potent cells we performed DE analysis between high and
182    low potent cells, irrespective of their epithelial subtype. The great majority of genes were
183    downregulated in the more potent cells, with only 72 exhibiting overexpression (Bonferroni
184    adjusted P<0.05, **Fig.4C**). Correspondingly, among the 1369 TFs, 582 exhibited differential
185    expression (Bonferroni adjusted P < 0.05) with only 3 TFs (*ENO1, YBX1* and *BTF3*)
186    exhibiting higher expression in the more potent cells (**Fig.4C-D**). Remarkably, *YBX1* and
187    *ENO1* are two transcription factors whose targets are highly enriched for breast cancer
188    GWAS eQTLs (Castro et al. 2016), thus implicating them in breast cancer risk. In addition,
189    siRNA against *YBX1* in a normal ER- cell-line (MCF10A) resulted in significantly reduced
190    cell-confluence and growth, even when compared to other breast cancer risk TFs (Castro et al.
191    2016).We confirmed that the associations of *YBX1* and *ENO1* expression with potency
192    remained after adjustment for cell-cycle phase (**SI fig.S6, Methods**), and that their expression
193    also correlated with cell potency in the scRNA-Seq data from the other 3 women (**SI fig.S7**).
194

195    **Upregulated bipotent single-cell signature correlates with mammary stemness**
196    If the highly potent cells are bipotent, the expectation would be that they are transcriptionally
197    similar to previously characterized mammary stem cells. We performed rank-based GSEA
198    (Subramanian et al. 2005) on the 72 genes upregulated in the highly potent single cells to
199    further characterize the putative bipotent cells. This revealed strong enrichment for ribosomal
200    genes, but importantly also for genes upregulated in mammary stem-cells (**SI fig.S8**). In
201    particular, we observed a relatively strong enrichment (12 gene overlap, OR=39, BH-adjusted
202    Fisher-test P<1e-10) with a previously characterized mammary stem-cell signature (Pece et al.
203    2010). Of note, among the 12 overlapping genes, 9 (*RPS2, RPS7, RPS10, RPL8, RPS18,*
204    *RPS3, RPL10A*) were ribosomal proteins or ubiquitin ribosomal fusion proteins (*UBA2 &*
205    *FAU)*, consistent with recent findings that expression of ribosomal proteins is a universal
206    marker of stemness and potency (Athanasiadis et al. 2017; Teschendorff and Enver 2017).
207    Among the other 3 genes, we observed *NACA*, a protein that associates with the upregulated
208    transcription factor *BTF3,* and *TXN* (thioredoxin), a protein involved in the response to
209    intracellular nitric oxide.
210    To confirm the results of the GSEA, we obtained and normalized mRNA expression data
211    from mammosphere-derived FACS sorted pools of quiescent mammary stem-cells and
212    transit-amplifying progenitors (Pece et al. 2010) (**Methods**). Validating the association with
213    stemness, the 12-genes exhibited increased expression in three separate pools of quiescent

214    mammary stem-cells compared to their derived transit-amplifying progenitors (**Fig.5A-B,**
215    Wilcox test P=0.001, **Methods**), a result which remained significant compared to randomly
216    selected genes (**Fig.5C**, Monte Carlo P=0.0001). Results remained significant had we used all
217    72 genes (63 genes had representation on the Affymetrix platform used in Pece et al (Pece et
218    al. 2010)) from the upregulated scRNA-Seq bipotent signature (**SI fig.S9**). However,
219    interestingly, *YBX1* and *ENO1* were not upregulated in the quiescent mammary stem cells
220    compared to the transit-amplifying progenitor cells (**SI fig.S9**), suggesting that while
221    potency/stemness is marked by the expression level of ribosomal proteins, the progenitor
222    non-quiescent state is associated with higher expression of *YBX1* and *ENO1*.

223

224

225    *YBX1* **expression correlates with luminal subtype and is increased in luminal**
226    **progenitors**
227    The correlation of *YBX1* expression with potency was particularly evident in the luminal
228    compartment (**Fig.6A, SI fig.S10**), pointing towards *YBX1* as playing not only a key role in
229    defining a basal progenitor phenotype, but also potentially as a luminal progenitor. We were
230    able to further validate this in two ways. First, its expression was also higher in the more
231    immature luminal alveolar-like phenotype, in line with the fact that these alveolar luminal
232    cells should be more enriched for progenitors (**Fig.6B**). Second, using bulk expression data
233    from FACS sorted luminal progenitor and differentiated luminal cells (Shehata et al. 2012),
234    we found *YBX1* expression to be highest for the EpCAM+/CD49f+/ALDH+ population
235    (**Fig.6C, Wilcox test P=0.003**), which defines the most likely luminal progenitor phenotype,
236    or at least the one that gives rise to milk-producing alveolar cells (Shehata et al. 2012).
237    Of note, we obtained similar results if instead of *YBX1* we used the earlier 17-gene or
238    72-gene signatures marking the bipotent cells. Indeed, the great majority of these genes were
239    observed to be overexpressed in the EpCAM+/CD49f+/ALDH+ population compared to all
240    other cell populations, a result which was highly significant as assessed using 100,000
241    Monte-Carlo randomizations (P<1e-5, **SI fig.S11**).

242

243    *YBX1* **expression marks basal breast cancer**
244    Given that *YBX1* exhibited highest expression in the more potent single-cells, and that these
245    were enriched within the basal compartment, it is natural to posit that *YBX1* may mark the
246    cell of origin for basal breast cancer. If so, *YBX1* expression should be highest in basal breast
247    cancer compared to other breast cancer subtypes. We were able to confirm this with high
248    statistical significance within the METABRIC study (Curtis et al. 2012), which profiled
249    almost 2000 primary breast cancers (**Fig.6D**). Similar results were obtained if instead of
250    *YBX1* we used the complete 17-gene or 72-gene signatures marking the bipotent cells (**SI**
251    **fig.S12**). In terms of the integrative cluster (IC) subtypes, as defined by METABRIC, *YBX1*
252    expression was highest in IC-5 and IC-10 (**Fig.6E**). These two integrative cluster subtypes

253   exhibited the worst disease-specific 5-year survival rates among all IC subtypes (Curtis et al.
254   2012). In line with this, we observed that *YBX1* expression also correlated with a poor clinical
255   outcome (**HR=1.31, P=7e-9, Fig.6F**). However, the association with outcome was mainly
256   driven by ER-status, since in an analysis stratified by ER-status we did not observe any
257   significant association (HR=1.11, P=0.12 in ER+; HR=0.96, P=0.64 in ER-).

258

259

260

## 261   **Discussion**

262   Here we have demonstrated "proof-of-concept" that our signaling entropy rate measure can
263   be used to identify rare subpopulations of highly-potent cells, which may represent novel
264   candidate progenitor or stem-like cells. Indeed, application to almost 4,000 single cells from
265   the mammary epithelium identified a rare (5%) subpopulation of relatively high potency,
266   which is likely to represent a mammary progenitor-like state. We extensively validated the
267   potency assignments of the single-cells, and consistent with the prevailing view that most
268   mammary progenitors are basal cells, the highly potent cells were over-represented within the
269   basal compartment. The ability to stratify single cells into different potency states allowed us
270   to infer and compare the cell-density surface maps for all potency states, revealing that highly
271   potent cells exhibited a strikingly different landscape to those of lower potency, with the
272   region of maximum cellular density defining a distinctive bi-modal ridge between the basal
273   and alveolar luminal clusters, with the largest peak occurring within the basal compartment.
274   Thus, without the need for any prior assumptions, LandSCENT predicts that these highly
275   potent cells may represent a bi-potent subpopulation that gives rise not only to basal cells but
276   also to luminal progenitors. Supporting this view, we found that the main TF characterizing
277   these highly potent cells (*YBX1*) plays a key role in maintaining the self-renewal and
278   proliferative capacity of basal cells (Castro et al. 2016) and that it is also overexpressed in
279   FACS sorted luminal progenitor populations compared to luminal differentiated cells. In
280   addition, we found that among the top-ranked genes upregulated in these putative bipotent
281   cells, there was a clear and significant enrichment for genes that have been found to mark
282   quiescent mammary stem cells and stemness generally (Athanasiadis et al. 2017;
283   Teschendorff and Enver 2017). We stress that these independent validations using orthogonal
284   expression data from bulk samples clearly shows that our results are not technical artefacts of
285   single-cell data.

286   The significance of *YBX1* extends to the cancer-risk context. First, there is already substantial
287   evidence demonstrating that *YBX1* transforms mammary epithelial cells, via binding to the
288   *BMI1* promoter and chromatin remodeling, leading to basal breast cancer (Davies et al. 2014).

289   In line with this, *YBX1* is also more highly expressed in basal breast cancer compared to all
290   other breast cancer subtypes, consistent with it marking cells that give rise to basal breast
291   cancer. Second, *YBX1* expression also marks luminal progenitor cells, and a subset of basal
292   breast cancers, notably BRCA1 mutant ones, are thought to arise from a mis-programmed
293   luminal progenitor (Lim et al. 2009; Shehata et al. 2012). Indeed, the single-cell landscape
294   inferred with LandSCENT underscores the similarity of the highly potent cells within the
295   basal compartment with those in the immature luminal cluster, strongly suggesting that the
296   cell of origin for basal breast cancer may well be a bi-potent like cell that shares an
297   expression profile similar to that of luminal progenitors, including notably *YBX1*. Third,
298   *YBX1* has been shown to interact with *ESR1,* and via *FGFR2* signaling may contribute to
299   tamoxifen resistance (Campbell et al. 2018). Fourth, it has been observed that genes within
300   the *YBX1* regulon are strongly enriched for GWAS breast cancer eQTLs (Castro et al. 2016).
301   This is a highly significant observation, given the growing evidence that molecular alterations
302   (both inherited and somatic) affecting the adult stem/progenitor cells within the tissue is a
303   main risk factor for epithelial cancer development (Tomasetti and Vogelstein 2015b;
304   Tomasetti and Vogelstein 2015a; Yang et al. 2016; Zhu et al. 2016; Tomasetti et al. 2017).
305   Thus, we speculate that it is the genetic and epigenetic alterations that accumulate within the
306   bi-potent progenitor cell pool identified here, which may confer the risk of breast cancer,
307   especially basal breast cancer.
308   In future, it will be important to conduct more comprehensive and deeper sequencing of
309   single cells in the mammary epithelium in order to construct accurate expression profiles for
310   the bi-potent cell pool identified here. In this regard, we point out that we were here severely
311   limited by the relatively low coverage of the 10X Chromium data (an average of only
312   ~60,000 reads per cell), which did not allow us to fully determine the differential expression
313   landscape of the bi-potent cells. The identification of *YBX1* (and *ENO1*) is a promising start,
314   but we anticipate that other regulators will also play a key role in defining these bi-potent
315   cells. We envisage that the computational framework presented here will play an important
316   role as a means of identifying and characterizing the bi-potent cells in the larger and deeper
317   scRNA-Seq studies to be performed in the near future. Importantly, LandSCENT will be
318   equally applicable to future large-scale scRNA-Seq studies performed on cancer tissue which
319   aim to identify putative cancer-stem-cells (Tirosh et al. 2016a; Tirosh et al. 2016b;
320   Teschendorff and Enver 2017).
321   In summary, we have presented a novel 3-dimensional clustering algorithm for scRNA-Seq
322   data, which uses an unbiased and assumption-free approach to estimate cell potency, and
323   which is used to perform single-cell clustering within each potency state. Application of this
324   simple yet powerful approach to scRNA-Seq data from the mammary epithelium naturally
325   predicts a bipotent cluster, which as shown here is characterized by regulators that have been
326   shown to modulate breast cancer risk. This study therefore provides a link between the
327   progenitor and stem like cell population that controls homeostasis within a complex epithelial

328  tissue and regulatory factors implicated in cancer risk of that same tissue. Our algorithm and
329  findings may serve as a general paradigm for analogous studies in other tissue types.
330

## Methods

332

### Single cell data and preprocessing

334  The scRNA-Seq data analysed in this work derives from the study of Nguyen et al (Nguyen
335  et al. 2018), who used the 10X Genomics Chromium platform to sequence a total of 24,646
336  cells from reduction mammoplastic specimens from 4 separate nulliparous women (Ind4-7),
337  at an average read-depth of 60,000 reads per cell. Mapped read count data from the 4
338  individuals was downloaded from GEO (GSE113197), and further normalized as follows: for
339  each cell we counted the number of expressed genes ("coverage per cell"), and for each gene
340  we also counted the number of times it was expressed across all single cells ("coverage per
341  gene"). For each cell, we also computed the total read count mapping to mitochondrial genes,
342  which revealed low cell coverage for those cells having a high proportion of mitochondrial
343  gene read counts. Based on this, we selected all cells expressing at least 1000 genes and with
344  the proportion of mitochondrial read counts less than 0.05, leaving a total of 23,369 cells.
345  Mitochondrial genes were removed and the total read count per cell $c$ recomputed ($TRC_c$).
346  Denoting the maximum of $TRC_c$ by $maxC$, and the read count matrix by $RCM$, the latter was
347  normalized with the following transformation: $LSC_{gc}=log_2(\ RCM_{gc}*maxC/TRC_c\ +\ 1.1)$.
348  Finally, we only use Entrez gene ID annotated genes, which resulted in a log-normalized
349  single cells matrix of dimension 22049 genes and 23369 cells (3473 for Ind-4, 6811 for Ind-5,
350  5807 for Ind-6 and 7278 for Ind-7).
351

### The <u>Land</u>scape <u>S</u>ingle-<u>C</u>ell <u>E</u>ntropy and Cell-<u>T</u>ype (LandSCENT) algorithm

353  LandSCENT is a direct extension of the SCENT algorithm. There are three steps to the
354  LandSCENT algorithm: (1) Inference of potency states: estimation of the differentiation
355  potency of single cells via computation of the signaling entropy rate (SR) and subsequent
356  inference of the potency state distribution across the single cell population. (2) Inference of
357  cell-types: we perform t-SNE (van der Maaten 2008) followed by density-based spatial
358  clustering (dbscan) (Ester et al. 1996) on a suitably dimensionally reduced $LSC$ matrix. (3)
359  Construction of an integrated landscape defined over potency-states and cell-types using
360  cell-density surface maps to reveal cellular-states. We note that step-1 is the exact same
361  procedure as used in our original SCENT algorithm (Teschendorff and Enver 2017).
362

363  <u>Step-1 Inference of potency states:</u> We estimate differentiation potency of each single cell by
364  computing the signaling entropy using the same prescription as used in our previous
365  publications (Banerji et al. 2013; Teschendorff et al. 2014). Briefly, the normalized

366 genome-wide gene expression profile of a sample (this can be a single cell or a bulk sample)

367 is used to assign weights to the edges of a highly curated protein-protein interaction (PPI)

368 network. The construction of the PPI network itself is described in detail elsewhere (Banerji

369 et al. 2013), and is obtained by integrating various interaction databases which form part of

370 Pathway Commons (www.pathwaycommons.org) (Cerami et al. 2011). The weighting of the

371 network via the transcriptomic profile of the cell provides the biological context. The weight

372 of an edge between protein $i$ and protein $j$, denoted by $w_{ij}$, is assumed to be proportional to

373 the normalized expression levels of the coding genes in the cell, i.e. we assume that $w_{ij} \sim x_i x_j$.

374 We interpret these weights (if normalized) as interaction probabilities. The above

375 construction of the weights is based on the assumption that in a sample with high expression

376 of $i$ and $j$, that the two proteins are more likely to interact than in a sample with low

377 expression of $i$ and/or $j$. Viewing the edges generally as signaling interactions, we can thus

378 define a random walk on the network, assuming we normalize the weights so that the sum of

379 outgoing weights of a given node $i$ is 1. This results in a stochastic matrix, $P$, over the

380 network, with entries

$$p_{ij} = \frac{x_j}{\sum_{k \in N(i)} x_k} = \frac{x_j}{(Ax)_i}$$

381 where $N(i)$ denotes the neighbors of protein $i$, and where $A$ is the adjacency matrix of the PPI

382 network ($A_{ij}=1$ if $i$ and $j$ are connected, 0 otherwise, and with $A_{ii}=0$). The signaling entropy is

383 then defined as the entropy rate (denoted $Sr$) over the weighted network, i.e.

$$Sr(\vec{x}) = -\sum_{i=1}^{n} \pi_i \sum_{j \in N(i)} p_{ij} \log p_{ij}$$

384 where $\pi$ is the invariant measure, satisfying $\pi P = \pi$ and the normalization constraint $\pi^T \mathbf{1} = 1$.

385 The invariant measure, also known as steady-state probability, represents the relative

386 probability of finding the random walker at a given node in the network (under steady state

387 conditions i.e. long after the walk is initiated). Nodes with high values thus represent nodes

388 that are particularly influential in distributing signaling flux in the network. In the

389 steady-state we can assume detailed balance (conservation of signaling flux, i.e. $\pi_i p_{ij} =$

390 $\pi_j p_{ji}$ ), and it can be shown (Teschendorff et al. 2014) that $\pi_i = x_i (Ax)_i/(x^T Ax)$. Given a

391 fixed adjacency matrix $A$ (i.e. fixing the topology), it can also be shown (Teschendorff et al.

392 2014) that the maximum possible $Sr$ among all compatible stochastic matrices $P$, is the one

393 with $P = \frac{1}{\gamma} v^{-1} \otimes A \otimes v$ where $\otimes$ denotes product of matrix entries and where $v$ is the

394 dominant eigenvector of $A$, i.e. $Av=\lambda v$ with $\lambda$ the largest eigenvalue of $A$. We denote this

395 maximum entropy rate by $maxSr$, and define the normalized entropy rate (with range of

396 values between 0 and 1) as

$$SR(\vec{x}) = \frac{Sr(\vec{x})}{maxSr}$$

397 Since $SR$ is bounded between 0 and 1, we next transform the $SR$ value of each single cell into

398   their logit-scale value, i.e. $y(SR)=log_2(SR/(1-SR))$. Subsequently, we fit a mixture of
399   Gaussians to the $y(SR)$ values of the whole cell population, and use the Bayesian Information
400   Criterion (BIC) (as implemented in the *mclust* R-package) (Yeung et al. 2001) to estimate the
401   optimal number $K$ of potency states, as well as the state-membership probabilities of each
402   individual cell. Thus, for each single cell, this results in its assignment to a specific potency
403   state.

404

405   <u>Step-2 Inference of cell-types</u>: Cell-types are inferred as significant clusters using
406   cell-density in the two-dimensional t-SNE space as the main criterion. Preliminary
407   dimensional reduction is achieved by first selecting genes with a mean average expression
408   larger than 1, and also a standard deviation larger than 1. These thresholds were chosen after
409   inspection of the mean-variance plot, and in the case of Ind-4 this resulted in 4261 highly
410   variable and expressed genes. To map the high dimensional nature of the data matrix to a
411   two-dimensional subspace we used t-SNE with an initial dimension of 30, a perplexity
412   parameter of 30, 1000 maximum iterations and epoch parameter set to 100. We then used the
413   dbscan algorithm (density-based spatial clustering) with eps=5 and minPts=15 to identify
414   significant clusters. Thus, after steps-1 and 2, each cell is assigned to a unique potency state
415   and co-expression cluster (cell-type).

416   <u>Step-3 Inference and construction of an integrated landscape of cell-states:</u> Finally, we
417   construct cell-density surface maps for all single cells within each of the inferred potency
418   states. In these surface maps, the elevation is directly proportional to cell-density. By
419   comparing the resulting landscapes for each potency state, this may reveal novel cellular
420   states, defined by both potency and expression subtype.

421

422

423   **Estimation of cell-cycle and TPSC pluripotency scores**
424   To identify single cells in either the G1-S or G2-M phases of the cell-cycle we followed the
425   procedure described in (Tirosh et al. 2016a). Briefly, genes whose expression is reflective of
426   G1-S or G2-M phase were obtained from (Whitfield et al. 2002; Macosko et al. 2015). A
427   given normalized scRNA-Seq data matrix for a given individual is then z-score normalized
428   for all genes present in these signatures. Finally, a cycling score for each phase and each cell
429   is obtained as the average z-scores over all genes present in each signature. When adjusting
430   differential expression analyses for cell-cycle phase, we included the G1-S and G2-M scores
431   as covariates in the linear models.

432

433   **Bulk expression datasets**
434   In this study we used three mRNA expression datasets from bulk samples. One dataset

435 consists of 38 FACS sorted bulk samples (Illumina expression beadarrays), as profiled by
436 Shehata et al (Shehata et al. 2012). Of the 38 samples, 10 were categorized as luminal
437 non-clonogenic (L), i.e. terminally differentiated cells, with the rest (n=28) making up a
438 relatively differentiated (EpCAM+/CD49f+/ALDH-, n=17) and undifferentiated
439 (EpCAM+/CD49f+/ALDH+, n=11) luminal progenitor (LP) populations. The two
440 undifferentiated LP populations were further distinguished by expression or not of ERBB3.
441 mRNA expression data was generated using Illumina Beadarrays and we used the normalized
442 data, as described in (Shehata et al. 2012).
443 The second dataset is the METABRIC study, which profiled almost 2000 primary breast
444 cancers using Illumina expression beadarrays (Curtis et al. 2012). We used the assignment of
445 tumors to PAM50 intrinsic and integrative cluster (IC) subtypes as given by the METABRIC
446 study. We used the normalized data, as provided by the METABRIC consortium.
447 A third Affymetrix mRNA expression dataset derives from Pece et al (Pece et al. 2010). This
448 set consists of 3 separate pools of FACS sorted cell populations. Each pool contains a
449 quiescent putative mammary stem cell population, as well as a population of derived progeny,
450 consisting of transit-amplifying progenitor cells, thus a total of 6 bulk samples. We
451 normalized the HGU133 plus2 data using the affy BioC package, specifically, the rma
452 function. Only probes mapping to an Entrez gene ID were used, data was quantile normalized
453 using limma, and probes mapping to the same gene were averaged, resulting in a normalized
454 data matrix over 20186 genes and 6 samples.
455

456 **Differential Expression Analysis**
457 When performing differential expression analysis on single-cell data, for each gene we
458 always restrict to those cells where the gene is expressed. That is, we remove all dropouts and
459 don't impute data. When correlating to potency, we used a linear model between the
460 normalized expression profile and the potency estimates, optionally adjusting for the two
461 cell-cycle scores computed earlier. In the case of the Illumina beadarray datasets, we used the
462 normalized data from the respective publications (Curtis et al. 2012; Shehata et al. 2012) and
463 called DE using the empirical Bayes limma framework (Smyth 2004). We always use
464 Bonferroni-adjusted thresholds to call statistical significance unless there are too few hits, in
465 which case we relax the threshold using FDR<0.05 instead.
466

467

468 **Doublet score analysis**
469 We used two different simulation-based methods to derive doublet scores for each cell and to
470 identify those more likely to be doublets. One approach used the simulation method of Dahlin
471 et al (Dahlin et al. 2018) to obtain doublet scores for all single cells that passed QC and for
472 each individual separately. Specifically, we used the doubletCells function (using
473 approximate=TRUE option) from the *scran* R-package (version 1.10.1) (Lun et al. 2016). In

13

474    the second approach we used the Python package Scrublet (Wolock et al. 2018) (doi:

475    https://doi.org/10.1101/357368). Within Scrublet, the scrub_doublets function, which is

476    responsible for computing doublet scores and predicting doublets within a dataset, was run

477    using default parameters.

478

479    **Code Availability:** SCENT is freely available as an R-package from github:

480    https://github.com/aet21/SCENT

481

482    **Data Access:** Data analyzed in this manuscript is already publicly available from the

483    following GEO (www.ncbi.nlm.nih.gov/geo/) accession numbers: GSE113197, GSE35399,

484    GSE18931 or from the EGA (www.ebi.ac.uk/ega/) accession number EGAS00000000083.

485

486

## Acknowledgements

492    **Disclosure Declaration** The authors declare that they have no competing interests.

493

## References

495

496    Athanasiadis EI, Botthof JG, Andres H, Ferreira L, Lio P, Cvejic A. 2017. Single-cell RNA-sequencing uncovers
497            transcriptional states and fate decisions in haematopoiesis. *Nature communications* **8**: 2045.
498    Banerji CR, Miranda-Saavedra D, Severini S, Widschwendter M, Enver T, Zhou JX, Teschendorff AE. 2013.
499            Cellular network entropy as the energy potential in Waddington's differentiation landscape. *Scientific*
500            *reports* **3**: 3039.
501    Campbell TM, Castro MAA, de Oliveira KG, Ponder BAJ, Meyer KB. 2018. ERalpha Binding by Transcription
502            Factors NFIB and YBX1 Enables FGFR2 Signaling to Modulate Estrogen Responsiveness in Breast Cancer.
503            *Cancer research* **78**: 410-421.
504    Castro MA, de Santiago I, Campbell TM, Vaughn C, Hickey TE, Ross E, Tilley WD, Markowetz F, Ponder BA,
505            Meyer KB. 2016. Regulators of genetic risk of breast cancer identified by integrative network analysis.
506            *Nature genetics* **48**: 12-21.
507    Cerami EG, Gross BE, Demir E, Rodchenkov I, Babur O, Anwar N, Schultz N, Bader GD, Sander C. 2011. Pathway
508            Commons, a web resource for biological pathway data. *Nucleic acids research* **39**: D685-690.
509    Chen J, Schlitzer A, Chakarov S, Ginhoux F, Poidinger M. 2016. Mpath maps multi-branching single-cell

510      trajectories revealing progenitor cell progression during development. *Nature communications* **7**:
511           11988.
512      Costa F, Grun D, Backofen R. 2018. GraphDDP: a graph-embedding approach to detect differentiation pathways
513           in single-cell-data using prior class knowledge. *Nature communications* **9**: 3685.
514      Curtis C, Shah SP, Chin SF, Turashvili G, Rueda OM, Dunning MJ, Speed D, Lynch AG, Samarajiwa S, Yuan Y et al.
515           2012. The genomic and transcriptomic architecture of 2,000 breast tumours reveals novel subgroups.
516           *Nature* doi:10.1038/nature10983.
517      Dahlin JS, Hamey FK, Pijuan-Sala B, Shepherd M, Lau WWY, Nestorowa S, Weinreb C, Wolock S, Hannah R,
518           Diamanti E et al. 2018. A single-cell hematopoietic landscape resolves 8 lineage trajectories and
519           defects in Kit mutant mice. *Blood* **131**: e1-e11.
520      Davies AH, Reipas KM, Pambid MR, Berns R, Stratford AL, Fotovati A, Firmino N, Astanehe A, Hu K, Maxwell C et
521           al. 2014. YB-1 transforms human mammary epithelial cells through chromatin remodeling leading to
522           the development of basal-like breast cancer. *Stem Cells* **32**: 1437-1450.
523      Ester M, Kriegel HP, Sander J, Xu X. 1996. A Density-Based Algorithm for Discovering Clusters in Large Spatial
524           Databases with Noise. In *2nd International Conference on Knowledge Discovery and Data Mining*
525           *(KDD-96)*. Institute for Computer Science, University of Munich.
526      Grun D. 2018. Revealing routes of cellular differentiation by single-cell RNA-seq. *Curr Opin Syst Biol* **11**: 9-17.
527      Grun D, Muraro MJ, Boisset JC, Wiebrands K, Lyubimova A, Dharmadhikari G, van den Born M, van Es J, Jansen
528           E, Clevers H et al. 2016. De Novo Prediction of Stem Cell Identity using Single-Cell Transcriptome Data.
529           *Cell stem cell* **19**: 266-277.
530      Guo M, Bao EL, Wagner M, Whitsett JA, Xu Y. 2016. SLICE: determining cell differentiation and lineage based on
531           single cell entropy. *Nucleic acids research* doi:10.1093/nar/gkw1278.
532      Haber AL, Biton M, Rogel N, Herbst RH, Shekhar K, Smillie C, Burgin G, Delorey TM, Howitt MR, Katz Y et al.
533           2017. A single-cell survey of the small intestinal epithelium. *Nature* **551**: 333-339.
534      Haghverdi L, Buttner M, Wolf FA, Buettner F, Theis FJ. 2016. Diffusion pseudotime robustly reconstructs lineage
535           branching. *Nature methods* **13**: 845-848.
536      Han X, Wang R, Zhou Y, Fei L, Sun H, Lai S, Saadatpour A, Zhou Z, Chen H, Ye F et al. 2018. Mapping the Mouse
537           Cell Atlas by Microwell-Seq. *Cell* **173**: 1307.
538      Hon CC, Shin JW, Carninci P, Stubbington MJT. 2018. The Human Cell Atlas: Technical approaches and challenges.
539           *Briefings in functional genomics* **17**: 283-294.
540      Laurenti E, Gottgens B. 2018. From haematopoietic stem cells to complex differentiation landscapes. *Nature*
541           **553**: 418-426.
542      Lim E, Vaillant F, Wu D, Forrest NC, Pal B, Hart AH, Asselin-Labat ML, Gyorki DE, Ward T, Partanen A et al. 2009.
543           Aberrant luminal progenitors as the candidate target population for basal tumor development in
544           BRCA1 mutation carriers. *Nature medicine* **15**: 907-913.
545      Lun AT, McCarthy DJ, Marioni JC. 2016. A step-by-step workflow for low-level analysis of single-cell RNA-seq
546           data with Bioconductor. *F1000Res* **5**: 2122.
547      Macosko EZ, Basu A, Satija R, Nemesh J, Shekhar K, Goldman M, Tirosh I, Bialas AR, Kamitaki N, Martersteck EM
548           et al. 2015. Highly Parallel Genome-wide Expression Profiling of Individual Cells Using Nanoliter
549           Droplets. *Cell* **161**: 1202-1214.
550      Marco E, Karp RL, Guo G, Robson P, Hart AH, Trippa L, Yuan GC. 2014. Bifurcation analysis of single-cell gene
551           expression data reveals epigenetic landscape. *Proceedings of the National Academy of Sciences of the*
552           *United States of America* **111**: E5643-5650.
553      Nguyen QH, Pervolarakis N, Blake K, Ma D, Davis RT, James N, Phung AT, Willey E, Kumar R, Jabart E et al. 2018.

554          Profiling human breast epithelial cells using single cell RNA sequencing identifies cell diversity. *Nature*
555          *communications* **9**: 2028.
556 Patel AP, Tirosh I, Trombetta JJ, Shalek AK, Gillespie SM, Wakimoto H, Cahill DP, Nahed BV, Curry WT, Martuza
557          RL et al. 2014. Single-cell RNA-seq highlights intratumoral heterogeneity in primary glioblastoma.
558          *Science* **344**: 1396-1401.
559 Pece S, Tosoni D, Confalonieri S, Mazzarol G, Vecchi M, Ronzoni S, Bernard L, Viale G, Pelicci PG, Di Fiore PP.
560          2010. Biological and molecular heterogeneity of breast cancers correlates with their cancer stem cell
561          content. *Cell* **140**: 62-73.
562 Qiu X, Mao Q, Tang Y, Wang L, Chawla R, Pliner HA, Trapnell C. 2017. Reversed graph embedding resolves
563          complex single-cell trajectories. *Nature methods* **14**: 979-982.
564 Regev A, Teichmann SA, Lander ES, Amit I, Benoist C, Birney E, Bodenmiller B, Campbell P, Carninci P,
565          Clatworthy M et al. 2017. The Human Cell Atlas. *eLife* **6**.
566 Rozenblatt-Rosen O, Stubbington MJT, Regev A, Teichmann SA. 2017. The Human Cell Atlas: from vision to
567          reality. *Nature* **550**: 451-453.
568 Scialdone A, Tanaka Y, Jawaid W, Moignard V, Wilson NK, Macaulay IC, Marioni JC, Gottgens B. 2016. Resolving
569          early mesoderm diversification through single-cell expression profiling. *Nature* **535**: 289-293.
570 Shehata M, Teschendorff A, Sharp G, Novcic N, Russell A, Avril S, Prater M, Eirew P, Caldas C, Watson CJ et al.
571          2012. Phenotypic and functional characterization of the luminal cell hierarchy of the mammary gland.
572          *Breast cancer research : BCR* **14**: R134.
573 Shepherd MS, Li J, Wilson NK, Oedekoven CA, Li J, Belmonte M, Fink J, Prick JCM, Pask DC, Hamilton TL et al.
574          2018. Single-cell approaches identify the molecular network driving malignant hematopoietic stem
575          cell self-renewal. *Blood* **132**: 791-803.
576 Shi J, Teschendorff AE, Chen L, Li T. 2018a. Quantifying Waddington's epigenetic landscape: a comparison of
577          single-cell potency measures. *Briefings in bioinformatics* **In Press**.
578 Shi J, Teschendorff AE, Chen W, Chen L, Li T. 2018b. Quantifying Waddington's epigenetic landscape: a
579          comparison of single-cell potency measures. *Briefings in bioinformatics* doi:10.1093/bib/bby093.
580 Smyth GK. 2004. Linear models and empirical bayes methods for assessing differential expression in microarray
581          experiments. *Statistical applications in genetics and molecular biology* **3**: Article3.
582 Stingl J, Eirew P, Ricketson I, Shackleton M, Vaillant F, Choi D, Li HI, Eaves CJ. 2006. Purification and unique
583          properties of mammary epithelial stem cells. *Nature* **439**: 993-997.
584 Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, Paulovich A, Pomeroy SL, Golub TR,
585          Lander ES et al. 2005. Gene set enrichment analysis: a knowledge-based approach for interpreting
586          genome-wide expression profiles. *Proceedings of the National Academy of Sciences of the United*
587          *States of America* **102**: 15545-15550.
588 Teschendorff AE, Enver T. 2017. Single-cell entropy for accurate estimation of differentiation potency from a
589          cell's transcriptome. *Nature communications* **8**: 15599.
590 Teschendorff AE, Sollich P, Kuehn R. 2014. Signalling entropy: A novel network-theoretical framework for
591          systems analysis and interpretation of functional omic data. *Methods* **67**: 282-293.
592 Tirosh I, Izar B, Prakadan SM, Wadsworth MH, 2nd, Treacy D, Trombetta JJ, Rotem A, Rodman C, Lian C, Murphy
593          G et al. 2016a. Dissecting the multicellular ecosystem of metastatic melanoma by single-cell RNA-seq.
594          *Science* **352**: 189-196.
595 Tirosh I, Venteicher AS, Hebert C, Escalante LE, Patel AP, Yizhak K, Fisher JM, Rodman C, Mount C, Filbin MG et
596          al. 2016b. Single-cell RNA-seq supports a developmental hierarchy in human oligodendroglioma.
597          *Nature* **539**: 309-313.

598     Tomasetti C, Li L, Vogelstein B. 2017. Stem cell divisions, somatic mutations, cancer etiology, and cancer
599             prevention. *Science* **355**: 1330-1334.
600     Tomasetti C, Vogelstein B. 2015a. Cancer etiology. Variation in cancer risk among tissues can be explained by
601             the number of stem cell divisions. *Science* **347**: 78-81.
602     Tomasetti C, Vogelstein B. 2015b. Cancer risk: role of environment-response. *Science* **347**: 729-731.
603     Trapnell C. 2015. Defining cell types and states with single-cell genomics. *Genome research* **25**: 1491-1498.
604     Trapnell C, Cacchiarelli D, Grimsby J, Pokharel P, Li S, Morse M, Lennon NJ, Livak KJ, Mikkelsen TS, Rinn JL. 2014.
605             The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single
606             cells. *Nature biotechnology* **32**: 381-386.
607     Treutlein B, Brownfield DG, Wu AR, Neff NF, Mantalas GL, Espinoza FH, Desai TJ, Krasnow MA, Quake SR. 2014.
608             Reconstructing lineage hierarchies of the distal lung epithelium using single-cell RNA-seq. *Nature* **509**:
609             371-375.
610     Treutlein B, Lee QY, Camp JG, Mall M, Koh W, Shariati SA, Sim S, Neff NF, Skotheim JM, Wernig M et al. 2016.
611             Dissecting direct reprogramming from fibroblast to neuron using single-cell RNA-seq. *Nature* **534**:
612             391-395.
613     van der Maaten L. 2008. Visualizing Data using t-SNE. *Journal of machine learning research : JMLR* **9**:
614             2579-2605.
615     Whitfield ML, Sherlock G, Saldanha AJ, Murray JI, Ball CA, Alexander KE, Matese JC, Perou CM, Hurt MM, Brown
616             PO et al. 2002. Identification of genes periodically expressed in the human cell cycle and their
617             expression in tumors. *Molecular biology of the cell* **13**: 1977-2000.
618     Wolock SL, Lopez R, Klein AM. 2018. Scrublet: computational identification of cell doublets in single-cell
619             transcriptomic data. *bioRxiv* doi:https://doi.org/10.1101/357368.
620     Yang Z, Wong A, Kuh D, Paul DS, Rakyan VK, Leslie RD, Zheng SC, Widschwendter M, Beck S, Teschendorff AE.
621             2016. Correlation of an epigenetic mitotic clock with cancer risk. *Genome biology* **17**: 205.
622     Yeung KY, Fraley C, Murua A, Raftery AE, Ruzzo WL. 2001. Model-based clustering and data transformations for
623             gene expression data. *Bioinformatics* **17**: 977-987.
624     Yuan GC, Cai L, Elowitz M, Enver T, Fan G, Guo G, Irizarry R, Kharchenko P, Kim J, Orkin S et al. 2017. Challenges
625             and emerging directions in single-cell analysis. *Genome biology* **18**: 84.
626     Zhu L, Finkelstein D, Gao C, Shi L, Wang Y, Lopez-Terrada D, Wang K, Utley S, Pounds S, Neale G et al. 2016.
627             Multi-organ Mapping of Cancer Risk. *Cell* **166**: 1132-1146 e1137.
628
629

630

631     # Figure Legends

632

633     **Figure-1: Flowchart of the LandSCENT algorithm to construct an integrative landscape
634     of cell-states from scRNA-Seq data. A)** Left: Signaling entropy (SR) is applied to the
635     scRNA-Seq profile of each individual cell to estimate its differentiation potency and to infer
636     potency states. Clustering of single-cells is performed with t-SNE followed by density based
637     spatial clustering to identify clusters of high cell-density, which we call cell-types. Right:
638     Surface cell-density map representation in t-SNE space for all single cells, showing the main

639    cell-types, with the smoothed SR (potency) values projected at the bottom. **B)** An example of

640    an integrated layered landscape of cellular states, where surface cell-density maps are shown

641    for cells in each inferred potency state (low, medium and high potency), defining cell-states

642    within or between major cell-types. The integrated landscape can reveal cell-states not

643    discernable via standard two dimensional clustering (shown at the bottom of each landscape).

644

645    **Figure-2: Inferring cell-types and potency states in breast epithelium. A)** t-SNE

646    clustering diagram for single-cells derived from one individual (Ind-4). Single-clusters were

647    inferred with dbscan and are labeled with different colors. Of note, single-cells that mapped

648    to the periphery of clusters and therefore were not assigned to any cluster have been

649    suppressed. **B)** As A), but now with the single cells labeled by expression levels of *KRT14* (a

650    basal marker), *KRT18* (a luminal marker), *LTF* (lactotransferin) and mean expression of

651    *GATA3, FOXA1, ESR1* and *PGR*, as indicated. Different quantiles of expression levels of

652    each marker are indicated by color with brown indicating high expression and grey low

653    expression. **C)** As A), but now displaying all single cells (i.e. including those mapping to the

654    periphery of clusters) and with single-cells labeled by the inferred potency state (see D)). **D)**

655    **Left panel:** Gaussian mixture model fit to the logit transformed SR values (x-axis) from 3473

656    single cells infers 3 potency states. The density distributions for all cells (black line) and

657    those for the inferred mixture components (different shades of blue) are shown. The Bayesian

658    Information Criterion (BIC) was used to select the optimal number of potency states, which

659    in this case was found to be 3 (PS1, PS2, PS3). **Right panel:** Percentage barplot indicating

660    the fraction of single-cells assigned to each of the three potency states.

661

662    **Figure-3: Validation of potency assignments. A)** Boxplots of normalized log-expression

663    (y-axis) for known markers of luminal differentiated cells (*GATA3, FOXA1*) and hormone

664    receptor (*ESR1*) against inferred potency state (x-axis) for all single cells assigned to the two

665    main luminal clusters (L1 & L2) and further restricting to cells where these genes are

666    expressed. Numbers of single-cells assigned to each potency state is given. P-value is from a

667    (two-tailed) linear regression. **B)** As A), but for known basal differentiation markers (*KRT16,*

668    *KRT5, EGFR*) and restricting to cells that were assigned to the basal cluster.

669

670

671    **Figure-4: Integrated landscape reveals bi-potent state characterized by YBX1 and**

672    **ENO1 expression. A)** Percentage barplots displaying the relative distribution of breast

673    epithelial subtypes (as inferred from the clustering using t-SNE + DBSCAN) among inferred

674    potency states (Low, Medium, High). Single cells have been divided up into whether they

675    clustered into the basal compartment (B), into the luminal-1 cluster (L1), the luminal-2

676    cluster (L2), all other clusters (Other) or whether they were not assigned into any cluster,

677    defining peripheral cells (Periph). P-value is from a Kruskal-Wallis test to assess if the
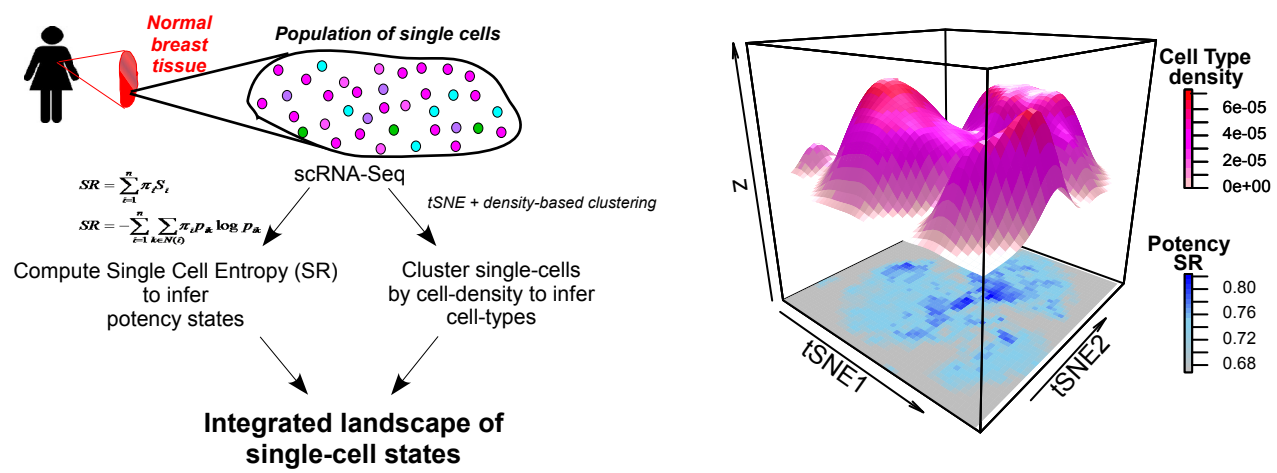
678　distribution of subtypes differs significantly within the high potency state. **B)** Surface

679　cell-density map of all single cells (magenta colored surfaces) with the corresponding surface

680　cell-density map of highly potent (PS3) cells superimposed (blueish colored surfaces). The x

681　and y-coordinates label the t-SNE1 and t-SNE2 axes. The height of the surfaces (z) is a

682　measure of cell-density in the x-y plane and is further indicated by different color tones. The

683　z-axis is therefore not a measure of cell potency. **C)** Volcano plot of differential expression

684　associated with potency, with x-axis labeling the t-statistic and y-axis labeling the statistical

685　significance. Horizontal bar denotes the Bonferroni threshold, and red points indicate

686　transcription factors (TFs). **D)** Boxplots of normalized log-expression (y-axis) for *YBX1* and

687　*ENO1* against inferred potency state (x-axis) for all single cells where these genes were

688　expressed. Numbers of single-cells assigned to each potency state is given. P-value is from a

689　two-tailed linear regression. All single-cell cells derive from one individual (Ind-4).

690

691　**Figure-5: Bipotent single-cell expression signature is enriched for mammary stem cell**

692　**genes. A)** Normalized relative expression heatmaps for 12 represented genes from the

693　17-genes upregulated in the putative bipotent single-cells and which overlap with a mammary

694　stem-cell signature, in 3 separate pools of FACS sorted quiescent mammary stem-cells (P)

695　and their derived proliferative non-stem like progeny (N). **B)** Average expression difference

696　between the P and N cells, averaged over the 3 separate pools. P-value is from a one-tailed

697　Wilcoxon rank sum test. **C)** Monte-Carlo randomization analysis, where in each of 100,000

698　random selections of 17 genes, the average difference over the 3 pools is computed (green

699　curve) and compared to the observed average difference (red, panel-B). Monte-Carlo P-value
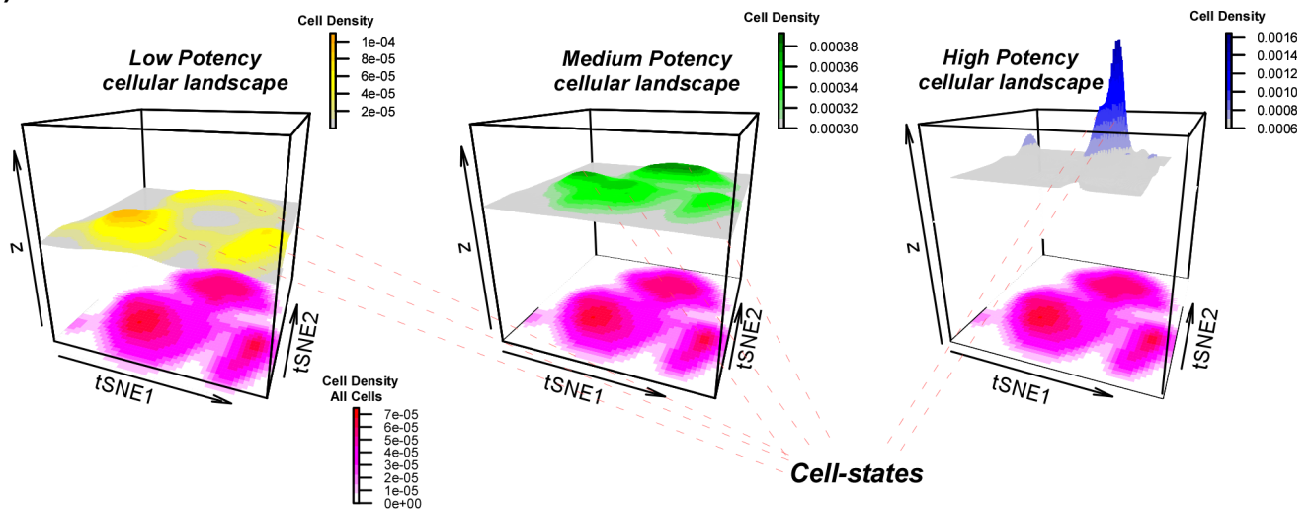
700　is given.

701

702

703　**Figure-6: YBX1 expression characterizes luminal progenitors and basal breast cancer. A)**

704　Boxplots of normalized log-expression (y-axis) for *YBX1* against inferred potency state

705　(x-axis) for all single cells assigned to luminal L1 and L2 clusters and where *YBX1* is

706　expressed. Numbers of single-cells in each group is given. P-value is from a linear regression.

707　**B)** Boxplots of normalized log-expression (y-axis) for *YBX1* against luminal cluster, using

708　only single cells where *YBX1* is expressed. Numbers of single-cells in each group is given.

709　P-value is from a one-tailed Wilcox test. **C)** Boxplots of Illumina normalized log-expression

710　(y-axis) for *YBX1* against luminal subtype as defined by FACS-sorting (x-axis):

711　L=differentiated non-clonogenic luminal, LP(ALDH-)=ALDH- luminal progenitor,

712　LP(ALDH+)=ALDH+ luminal progenitor. LP(ERBB3-)=ERBB3- and ALDH- luminal

713　progenitor. P-value is from a one tailed Wilcox-test comparing LP(ALDH+) to all others. **D)**

714　Boxplots of normalized Illumina log-expression (y-axis) for *YBX1* against PAM50 intrinsic

715　subtype in the full METABRIC cohort. P-value is from a Kruskal-Wallis test**. E)** As D), but

716　for the integrative cluster (IC) subtypes (available in discovery set only). **F)** Kaplan Meier
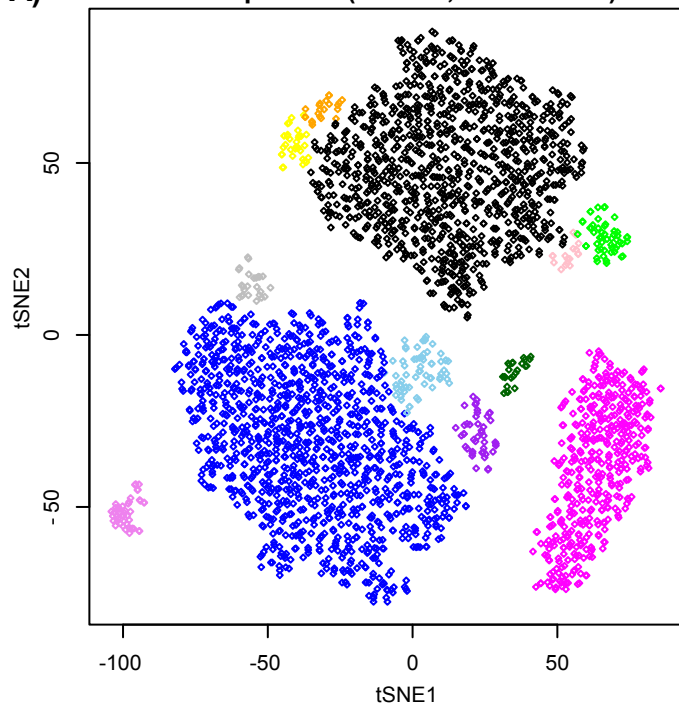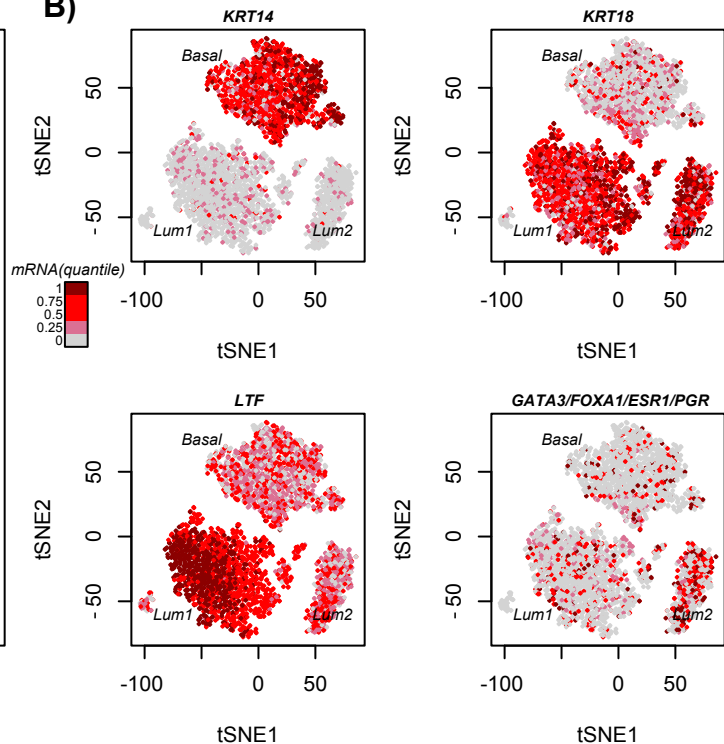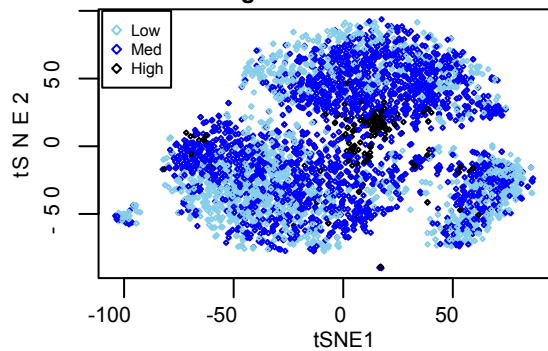
717     overall survival curves for *YBX1* expression, stratified by quantiles of *YBX1* expression, and

718     censored at 5 years after diagnosis. Hazard Ratio (HR), 95% CI and P-value are from a Cox

719     proportional hazards regression.

720

**A)** The LandSCENT algorithm

*Normal breast tissue*

*Population of single cells*

scRNA-Seq

*tSNE + density-based clustering*

$$SR = \sum_{i=1}^{n} \pi_i S_i$$

$$SR = -\sum_{i=1}^{n} \sum_{k \in N(i)} \pi_i p_{ik} \log p_{ik}$$

Compute Single Cell Entropy (SR) to infer potency states

Cluster single-cells by cell-density to infer cell-types

**Integrated landscape of single-cell states**

Cell Type density

6e-05
4e-05
2e-05
0e+00

Potency SR

0.80
0.76
0.72
0.68

tSNE1    tSNE2    Z

**B)**

*Low Potency cellular landscape*

Cell Density

1e-04
8e-05
6e-05
4e-05
2e-05

*Medium Potency cellular landscape*

Cell Density

0.00038
0.00036
0.00034
0.00032
0.00030

*High Potency cellular landscape*

Cell Density

0.0016
0.0014
0.0012
0.0010
0.0008
0.0006

Cell Density All Cells

7e-05
6e-05
5e-05
4e-05
3e-05
2e-05
1e-05
0e+00

Z    tSNE1    tSNE2

*Cell-states*

**A)** BrEpi Cells (n=3473, 1 individual)

**B)** *KRT14* *KRT18* *LTF* *GATA3/FOXA1/ESR1/PGR*

Basal Lum1 Lum2

mRNA(quantile)
1
0.75
0.5
0.25
0

**C)** Single Cell Potencies

Low
Med
High

**D)**

All
PS3 (High)
PS2 (Med)
PS1 (Low)

log2[SR/(1-SR)]

PS3 PS2 PS1

Fraction

**A)**

*GATA3: L1 & L2* — P=1e−17

*FOXA1: L1 & L2* — P=3e−07

*ESR1: L1 & L2* — P=5e−04

**B)**

*KRT16: Basal* — P=7e−06

*KRT5: Basal* — P=3e−06

*EGFR: Basal* — P=7e−28

**A)** Bar chart showing Fraction (y-axis) across Potency States (x-axis): Low(PS1), Medium(PS2), High(PS3). Legend: Periph (gray), B (black), Other (pink), L1 (blue), L2 (magenta). P<2e-16

**B)** 3D density plot with axes tSNE1, tSNE2, and Z. Cell Type Density All Cells (7e-05 to 0e+00, red-pink scale). Cell Type Density PS3 cells (5e-04 to 0e+00, blue scale). Labels: Lum-1, Lum-2, Basal, bi-potent like cells.

**C)** Volcano plot: -log10(P) (y-axis) vs t(DE:Potency) (x-axis). Legend: TF (red), Other (black). Labeled points: ENO1, YBX1, BTF3.

**D)** 
*YBX1* — Boxplot of mRNA across Potency State: Low n=1235, Med n=2019, High n=169. P=9e-27

*ENO1* — Boxplot of mRNA across Potency State: Low n=1268, Med n=2023, High n=169. P=1e-42

A) Pool-1, Pool-2, Pool-3

Gene labels (top to bottom): RPL10A, NACA, RPS3, RPS18, RPL8, RPS10, FAU, RPS7, RPS2, GAPDH, UBA52, TXN

N   P (Mammary Stem Cell Status)

Rel.Expr: 0.5, −0.5

B) MeanDiff, P=0.001

C) Density vs Av(MeanDiff); Null, Obs; P=0.0001

**A)** *YBX1: L1 & L2* — P=4e-09. mRNA vs Potency State (Low n=665, Med n=872, High n=33)

**B)** *YBX1: L1 & L2* — P=2e-08. mRNA vs Lum.Subtype (L1 n=1160, L2 n=410)

**C)** *YBX1: FACS luminals* — P=0.003. mRNA(Illum.) vs Luminal Subtype (L n=10, LP(ALDH-) n=10, LP(ALDH+) n=11, LP(ERBB3-) n=7)

**D)** *YBX1 in METABRIC (n=1980)* — P=2e-133. mRNA(Illum.) vs Pam50 subtype (LumB n=488, LumA n=718, Nor n=199, Basal n=329, HER2 n=240)

**E)** *YBX1 in METABRIC (n=997)* — P=2e-54. mRNA(Illum.) vs IC subtype (1 n=76, 2 n=45, 3 n=156, 4 n=167, 5 n=94, 6 n=44, 7 n=109, 8 n=143, 9 n=67, 10 n=96)

**F)** Survival vs YBX1 mRNA quantiles (n=1980). Prob(OS) vs years. 25%, 50%, 75%, 100%. HR=1.31 (1.19-1.43)  P=7e-9