# Transition to invasive breast cancer is associated with progressive changes in the structure and composition of tumor stroma

Tyler Risom[1], David R Glass[1], Candace C Liu[1], Belén Rivero-Gutiérrez[1], Alex Baranski[1], Erin F McCaffrey[1], Noah F Greenwald[1], Adam Kagel[1], Siri H Strand[1], Sushama Varma[1], Alex Kong[1], Leeat Keren[1], Sucheta Srivastava[1], Chunfang Zhu[1], Zumana Khair[1], Deborah J Veis[5], Katherine Deschryver[2], Sujay Vennam[1], Carlo Maley[4], E Shelley Hwang[3], Jefferey R Marks[3], Sean C Bendall[1], Graham A Colditz[2], Robert B West[1]*, Michael Angelo[1]*

[1]Stanford University School of Medicine, Department of Pathology; [2]Washington University School of Medicine, Department of Surgery; [3]Duke University, Department of Surgery; [4]Arizona State University, Biodesign institute; [5]Washington University School of Medicine, Departments of Pathology & Immunology

*co-corresponding authors

## Abstract

Ductal carcinoma *in situ* (DCIS) is a pre-invasive lesion that is thought to be a precursor to invasive breast cancer (IBC). To understand how the tumor microenvironment (TME) changes with transition to IBC, we used Multiplexed Ion Beam Imaging by time of flight (MIBI-TOF) and a 37-plex antibody staining panel to analyze 140 clinically annotated surgical resections covering the full spectrum of breast cancer progression. We compared normal, DCIS, and IBC tissues using machine learning tools for multiplexed cell segmentation, pixel-based clustering, and object morphometrics. Transition from DCIS to IBC was found to occur along a trajectory marked by coordinated shifts in location and function of myoepithelium, fibroblasts, and infiltrating immune cells in the surrounding stroma. Taken together, this comprehensive study within the HTAN Breast PreCancer Atlas offers insight into the etiologies of DCIS, its transition to IBC, and emphasizes the importance of the TME stroma in promoting these processes.

**Introduction**

Ductal Carcinoma in situ (DCIS) is a preinvasive lesion where tumor cells within the breast duct are isolated from the surrounding stroma by a near-continuous layer of myoepithelium and basement membrane proteins. This histologic feature is the central property that distinguishes it from invasive breast cancer (IBC), where this barrier has broken down and tumor cells have invaded the stroma (Figure 1A). DCIS comprises 20% of new breast cancer diagnoses, but unlike IBC, in itself is not a life-threatening disease. However, if left untreated, up to half of these patients will develop IBC within 10 years (Betsill et al., 1978; Erbas et al., 2006; Eusebi et al., 1994; Page et al., 1982; Ryser et al., 2019).

Sequencing-based approaches have been used extensively over the last decade to identify molecular features that could elucidate the connection between DCIS and IBC. Genomic profiling has identified recurrent copy number variants (CNV) that are more prevalent in high grade DCIS lesions (Afghahi et al., 2015; Buerger et al., 1999; Fujii et al., 1996). Meanwhile, comparison of paired DCIS and IBC lesions from the same patient has provided clues into the clonal evolution from *in situ* to invasive disease by revealing genomic alterations that are acquired during this transition (Ak et al., 2018; Kim et al., 2015; Newburger et al., 2013). To date, however, these findings have not been found to consistently explain this transition. Similarly, the utility of tumor phenotyping by single-plex immunohistochemical tissue staining has been limited as well.

In light of this uncertainty, clinical management has trended towards treating all patients presumptively as progressors with surgery, radiation therapy, and pharmacological interventions that carry risks for therapy-related adverse events. Consequently, this approach is likely to be overly aggressive for non-progressors. Thus, understanding the central biological features in DCIS that drive the transition to IBC is a critical unmet need.

Surprisingly, despite all the information now known about the genetic and functional state of tumor cells in DCIS, histopathology remains the only reliable way to diagnose it. Thus, DCIS is an intrinsically structured entity where the spatial orientation of tumor, myoepithelial, and stromal cells is the primary defining feature that distinguishes it from other forms of breast cancer.

69    To understand how DCIS structure and single cell function are interrelated, we use
70    new tools previously developed by our lab for highly multiplexed subcellular imaging to
71    analyze a large cohort of human archival tissue samples covering the spectrum of breast
72    cancer progression from *in situ* to invasive disease. In previous work, we used
73    Multiplexed Ion Beam Imaging by Time of Flight (MIBI-TOF) and a 36-plex antibody
74    staining panel to identify rule sets governing tumor microenvironment (TME) structure in
75    triple negative breast cancer that were highly predictive of the composition of immune
76    infiltrates, the expression of immune checkpoint drug targets, and 10-year overall survival
77    (Keren et al., 2018).

78    This effort provided a framework for how TME structure and composition could be
79    used more generally as a surrogate readout to understand the functional response to
80    neoplasia. With this in mind, we sought to determine to what extent similar features
81    involving myoepithelial, stromal, and immune cells in the DCIS TME might play a pivotal
82    role in breast cancer progression. Each of these have been implicated previously to
83    promote local invasion (Barsky and Karlin, 2005; Ibrahim et al., 2020), metastasis (Pelon
84    et al., 2020; Shani et al., 2020), and to correlate with clinical progression (Yang et al.,
85    2018; Zhou et al., 2018).

86    Here, we report the first systematic, high dimensional analysis of breast cancer
87    progression using the Washington University Resource Archival Human Breast Tissue
88    (RAHBT) cohort: a clinically annotated set of archival tissue from patients diagnosed with
89    DCIS and IBC. Because the DCIS patient population is complicated by differences in
90    age, parity status, tumor subtype, and treatment course, a well-conceived cohort design
91    is crucial for identifying meaningful features amidst these confounding variables. In light
92    of this, the RAHBT cohort was composed of primary DCIS tumors from women who later
93    progressed to invasive disease that were age and year-of-diagnosis matched with control
94    tissue from women with DCIS that did not recur.

95    We used MIBI-TOF and a 37-plex antibody staining panel to comprehensively
96    define the cellular composition and structural characteristics in 122 of these samples,
97    which included normal breast, DCIS, and recurrent IBC samples. We applied machine
98    learning tools for multiplexed cell segmentation and spatial analytics to enumerate 16 cell
99    populations and to quantify how these populations are spatially distributed relative to one

3

100  another.  Object morphometrics and high dimensional pixel clustering were used to

101  annotate the structure of stromal collagen and to discover new myoepithelial phenotypes

102  that track with disease progression.  These findings were corroborated by transcriptomic

103  data acquired on coregistered tissue regions isolated by laser capture microdissection.

104      We systematically compared these features to understand how different

105  phenotypic and structural properties of the DCIS TME change with progression to IBC.

106  BC progression was typified by a reduction in myoepithelial integrity, a shift in fibroblast

107  function towards proliferative cancer-associated states (CAFs), remodeling of collagen in

108  the extracellular matrix (ECM), and a compositional and spatial reorganization of the

109  immune microenvironment.  We used the 1,093 features quantified in these analyses to

110  build a random forest classifier for predicting which patients would later progress to

111  invasive disease based exclusively on the original diagnostic biopsy.  This classifier

112  demonstrated an AUC of 0.83 and was heavily weighted for stromal features that were

113  reliant on spatial information.  Taken together, this work provides new insight into potential

114  etiologies of DCIS progression that will guide development of future diagnostics and serve

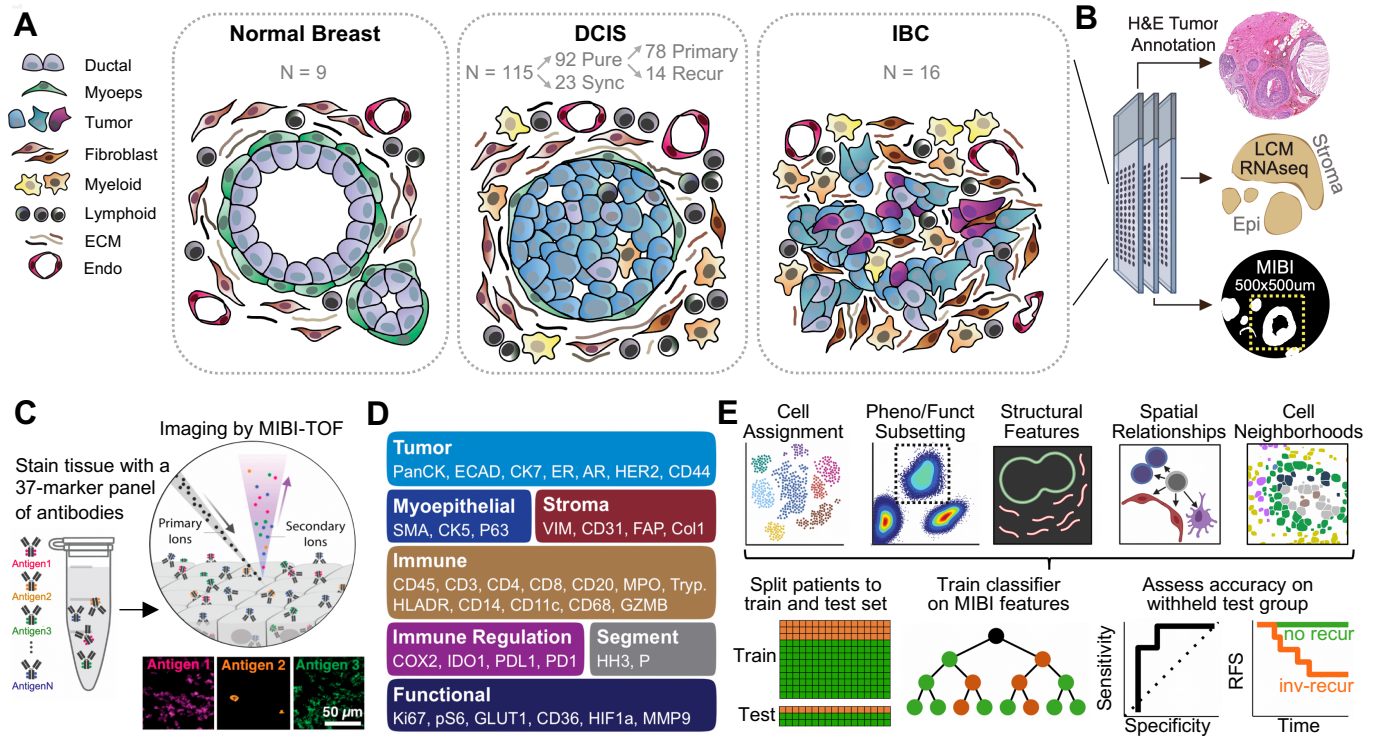115  as a template for how to carry out similar analyses of preinvasive cancers.

116

117

118  **Results**

119

120  **A multiplexed imaging interrogation of DCIS progression to invasive disease**

121  We examined the transition from DCIS to IBC by profiling accumulative changes in the

122  phenotype, structure, and spatial distribution of myoepithelium, tissue stroma, and

123  immune cells in archival formalin-fixed paraffin-embedded (FFPE) patient tissue of three

124  distinct progression groups: normal breast (n = 9), IBC (n = 16), and DCIS (n =

125  115).  These IBC samples were disease recurrences from women with a prior diagnosis

126  of DCIS.  Of the 115 DCIS samples, 78 were RAHBT patients with a new diagnosis and

127  no signs of IBC (pure, primary), while 14 were pure DCIS recurrences (pure, recur)(Figure

128  1A, Table S1).  The remaining 23 patients comprised a third group of synchronous lesions

4

**Figure 1. A multiplexed Imaging Interrogation of DCIS Progression to Invasive Disease**

**A.** Schematic depicting the tumor stages and patient sample numbers profiled in this study, including normal breast, pure DCIS (primary or recurrent), synchronous DCIS (Sync), and invasive breast carcinoma (IBC). **B.** Depiction of the parallel tissue analysis methods used in this study including H&E, laser capture microdissection (LCM) RNAseq, and MIBI-TOF. **C.** Overview of the MIBI-TOF workflow. **D.** Markers used in the MIBI-TOF panel are displayed, grouped by target cell type or protein class. **E.** Workflow showing feature types extracted from the MIBI-TOF analysis that were used to train a random forest classifier to differentiate DCIS samples with or without risk of recurrence.

129    procured at Stanford Hospital where both DCIS and IBC were identified in different parts

130    of the tissue at the time of diagnosis (Sync).  For this set of patients, only the *in situ*

131    component was analyzed.

132        1.5 mm cores of each tumor were arranged in tissue microarrays (TMAs).  Three

133    adjacent sections were then used for 1) H&E staining and annotation by a pathologist, 2)

134    RNA transcriptome analysis of ductal and stromal regions isolated using laser-capture

135    microdissection (LCM-Smart-3SEQ)(Foley et al., 2019), and 3) highly multiplexed

136    imaging by MIBI-TOF of a 500x500μm field-of-view (FOV)(Figure 1B).  By ensuring that

137    each of these analyses were spatially coregistered with one another, the proteomic and

138    transcriptomic features revealed by MIBI-TOF and LCM-RNAseq could be directly

139    correlated to understand the interplay between single cell composition and global
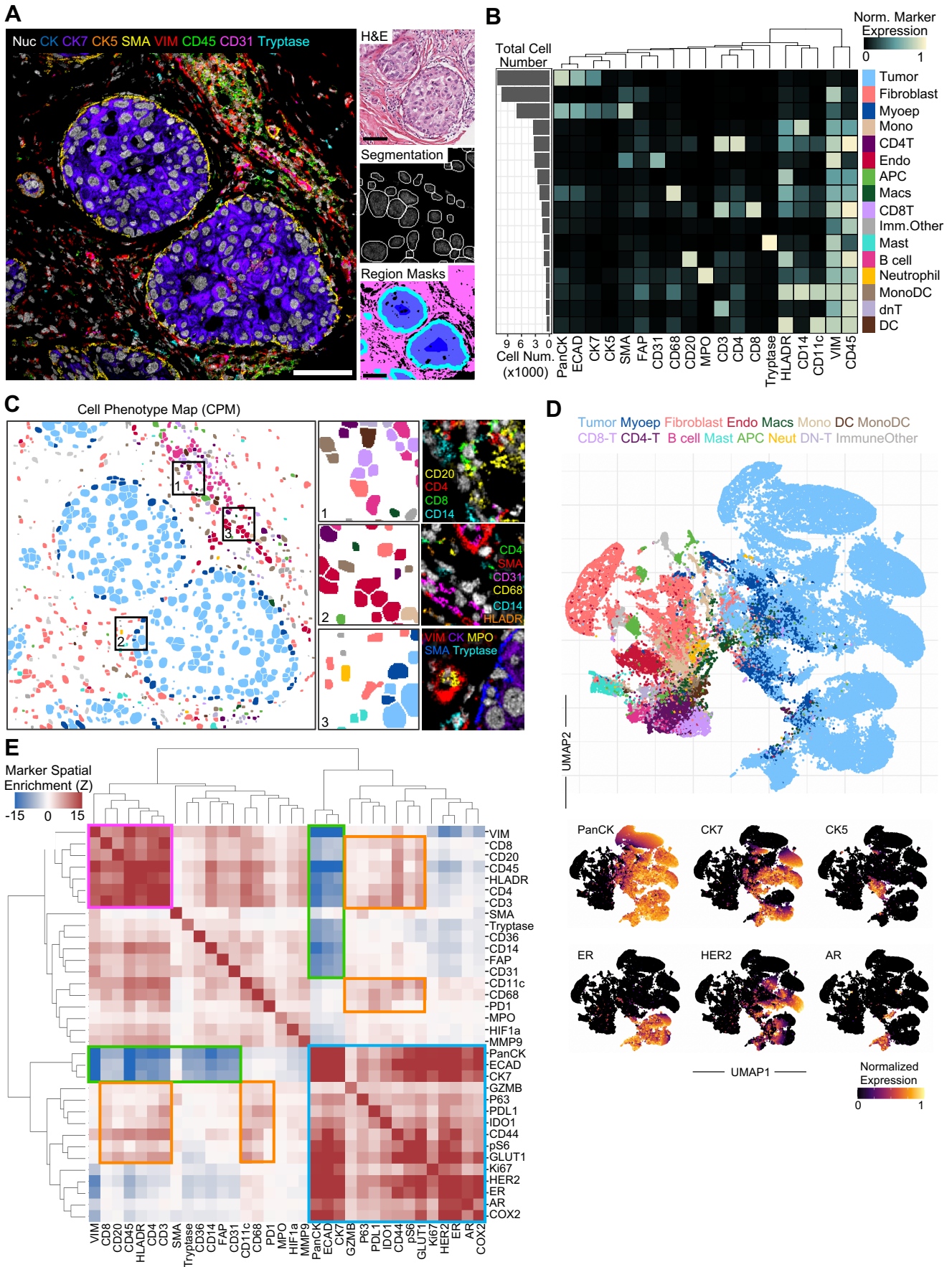
140    transcriptional programs.

141       For MIBI-TOF, we constructed a 37-plex staining panel of metal-conjugated

142    antibodies that would permit us to: 1) map the lineage and spatial location of every cell,

143    2) identify lineage subsets of tumor, fibroblasts, and immune cells previously implicated

144    in BC progression, and 3) characterize the composition, integrity, and morphology of

145    myoepithelium and collagen (Figure 1D, Table S2). The panel also included 11 functional

146    markers for annotating proliferation, activation, hypoxic signaling, as well as markers

147    implicated in cancer immunoregulation, including PD-L1, IDO1, COX2 and PD1 (Figure

148    S1). The features extracted in this analysis were then used to train a random forest

149    classifier for predicting long term outcome (Figure 1E).

150

## A single cell phenotypic and spatial atlas of DCIS

151

152    The workflow outlined in Figure 1 enabled high-dimensional, subcellular imaging of

153    dozens of proteins that recapitulated the tissue architecture observed in H&E (Figure

154    2A). Multiplexed imaging data were processed with a low-level pipeline prior to single-

155    cell segmentation (Figure 2B, Figure S2B)(Keren et al., 2018; McCaffrey et al., 2020;

156    Moen et al., 2019; Valen et al., 2016), which identified on average ~924 cells in each

157    FOV (sd = 317). To determine cell location with respect to canonical histological features,

158    we demarcated duct, stroma, and myoepithelial regions of each image based on

159    combinatorial marker expression (Figure 2B bottom-right). Importantly, throughout this

160    work we will be presenting cellular data either as the frequency of a parental lineage

161    across the entire image (e.g., macrophages as % of total immune cells) or as a cell density

162    within a particular compartment of the image (e.g., 50 fibroblasts/mm$^2$ of stroma).

163       Hierarchical application of the FlowSOM algorithm (Van Gassen et al., 2015) was

164    employed to identify 16 unique cell subsets in the dataset amongst the epithelial, stromal,

165    and immune lineages (Figure 2B, S2B). Altogether, we assigned 95% (n = 127,451 single

166    cells) of cells to one of these subsets that in aggregate ranged in frequency from 0.7-

167    56%. These data were used to generate cell phenotype maps (CPM) where each cell is

168    colored according to its subset assignment. CPM images illustrated focal enrichment of

169    lymphocytes (Figure 2C "1"), endothelial-associated immune phenotypes (Figure 2C, "2")

170    and sparser subsets of periductal granulocytes that included neutrophils and mast cells

171    (Figure 2C, "3").

**Figure 2. A single cell phenotypic and spatial atlas of DCIS**

**A.** Representative MIBI image overlay of a DCIS tumor with a 9-marker overlay of major cell lineage markers (left) and the corresponding H&E image (top right), example of cell segmentation (middle right), and example of region masks marking stroma (pink), myoepithelial (cyan) and ductal (blue) area, scale bars = 100μm. **B.** Cell lineage assignments based on normalized expression of lineage markers (heatmap columns), rows are ordered by absolute abundance shown in the bar plot (left), while columns are hierarchically clustered (euclidean distance, average linkage). **C.** A cell phenotype map (CPM) showing cell identity by color, as defined in *F,* overlaid onto the segmentation mask. Zoomed insets with adjacent MIBI overlays show diverse lymphoid rich regions (1), endothelial-associated immune cells (2) and rare subsets like neutrophils and mast cells near ducts (3). **D.** UMAP visualization of all cell type populations in DCIS tumors (top), colored by cell type as in *F,* with additional plots overlaid with the normalized expression of tumor lineage and functional markers used to delineate tumor subsets (bottom).

172   Tumor cells were the most abundant cell type in DCIS samples (60% ± 20 of all

173   cells) and were comprised of multiple subsets that were defined by variable expression

174   of the luminal and basal lineage markers (CK7 and CK5, respectively), as well as ER,

175   AR, and HER2 (Figure 2D). Since these cells are isolated by a layer of myoepithelium,

176   by definition the tissue structure of DCIS is highly compartmentalized. In order to

177   determine if our analyses were capturing this fundamental facet, we used an unbiased

178   computational approach to identify sets of proteins that colocalize or avoid one another

179   more frequently than would be expected by chance. Consistent with the

180   compartmentalized nature of DCIS, tumor cell markers were spatially enriched (PanCK,

181   ECAD, CK7, HER2, ER, AR, Figure 2E, blue box) and segregated from vascular,

182   fibroblast, and immune markers (Figure 2E, green box). With respect to the latter,

183   lymphoid markers demonstrated the most prominent spatial enrichment (Figure 2E,

184   magenta box). These analyses also revealed moderate preferential enrichment in tumor

185   positive regions for pS6, COX2, and Ki67, while immunoregulatory markers were more

186   evenly dispersed between tumor and immune-enriched regions (Figure 2E, orange box).

187

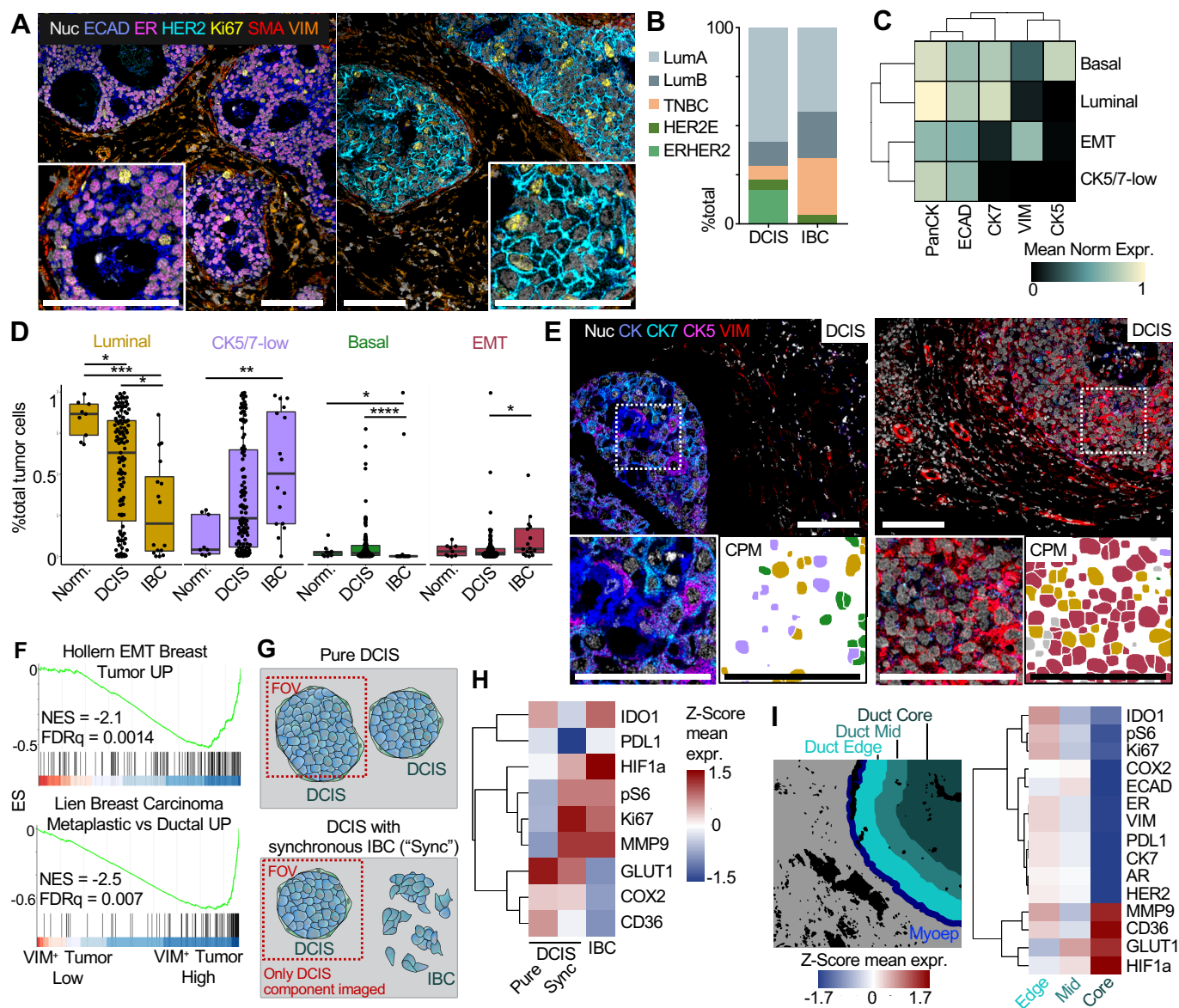188   **A tumor cell phenotypic switch marks invasive transition**

189   Tumor heterogeneity in breast cancer can manifest as variations in the level of hormone

190   receptor expression and the degree of luminal, basal, and mesenchymal differentiation.

191   DCIS has been shown to vary across the full spectrum of both of these axes, which can

192   confound identification of conserved features correlating with clinical outcome. In order

193   to understand how this heterogeneity manifests in pure DCIS and throughout the

194   transition to invasive disease, we first examined the distribution of DCIS subtypes with

195   respect to hormone receptor status (ER, AR), HER2, and Ki67 proliferation index. These

196   markers were robustly expressed in DCIS tumors (Figure 3A) and showed expected inter-

7

197   patient variability. Using clinical cutoffs as a guide (Figure S3A), we subtyped tumors as

198   Luminal A (ER$^+$, HER2$^-$, Ki67$^-$), Luminal B (ER$^+$, HER2$^-$, Ki67$^+$), HER2E (ER$^-$, HER2$^+$),

199   ERHER2 (ER$^+$, HER2$^+$), and TNBC (ER$^-$, HER2$^-$) based on the frequency of positive cells

200   for each marker. All subtypes were present in both DCIS and IBC, with similar numbers

201   of luminal samples in each progression group (Figure 3B). HER2$^+$ tumors were more

202   predominant in DCIS, while TNBC was more prevalent in IBC (Figure S3B-C).

203       On comparing epithelial differentiation states in each progression group, we

204   identified a consistent trend towards reduced luminal cell identity throughout tumor

205   progression. Distinct phenotypic subsets of luminal (CK7$^+$), basal (CK5$^+$), EMT-like

206   (VIM$^+$), and CK5/7-low cells were observed in the epithelial lineage (Figure 3C). While

207   the majority of ductal cells in normal breast were consistently luminal (84% $\pm$ 11) (Figure

208   3D), the composition in DCIS varied widely between being predominantly luminal or

209   CK5/7-low (57% $\pm$ 33, 36% $\pm$ 33 respectively). In comparison to normal tissue and IBC,

210   these lesions were also enriched with a minority fraction of basal cells (6.1% $\pm$ 11.9). With

211   progression to IBC, CK5/7-low cells predominate more frequently and were accompanied

212   by a relative increase in EMT-like cells that express vimentin (Figure 3E). We further

213   examined a subset of patients with high frequencies of vimentin-positive tumor cells by

214   LCM-RNAseq. Consistent with the shift to a mesenchymal phenotype captured by MIBI-

215   TOF, geneset enrichment analysis (GSEA) revealed upregulation of signaling pathways

216   relating to mesenchymal breast tumor histology and tumor invasion in patients with high

217   vs low frequencies of VIM$^+$ tumor cells (Hollern et al., 2018; Lien et al., 2007; Poola et al.,

218   2005)(Figure 3F, Figure S3D).

219       The coordinated changes in tumor phenotype illustrate how cell differentiation

220   during BC progression may follow an orderly trajectory. To further explore this possibility,

221   we compared tumor cell functional states in pure, DCIS synchronous DCIS, and

222   IBC. Synchronous DCIS describes lesions where distinct areas of tissue contained either

223   fully encapsulated tumor cells (i.e., DCIS) or areas of local invasion (i.e., IBC) were both

224   present at the time of diagnosis, but in different areas of tissue (Figure 3G). Consistent

225   with their more aggressive behavior, DCIS tumor cells from synchronous lesions

**Figure 3. A tumor cell phenotypic switch marks invasive transition**

**A.** Representative MIBI image overlays showing an ER$^+$HER2$^-$ tumor (left) and ER$^-$HER2$^+$ (right), scale bars = 100μm. **B.** Stacked barplot showing the distribution of intrinsic breast cancer subtypes in DCIS and IBC tumors, as defined by receptor expression. **C.** Tumor phenotype assignments based on normalized expression of markers related to markers of tumor differentiation (heatmap columns). **D.** Frequency of tumor differentiation states across normal breast, DCIS, and IBC. **E.** Representative MIBI image overlays of DCIS tumors with basal and mesenchymal features, respectively. Zoomed insets (left) with paired cell phenotype maps (right) colored by tumor phenotype identity as in *D*, scale bars = 100μm. **F.** Geneset enrichment analysis comparing VIM-high and low tumors with genesets related to mesenchymal tumor differentiation. **G.** Schematic showing the imaging FOV location in pure and synchronous DCIS tumors, which only included the DCIS component. **H.** Heatmap of z-score normalized functional marker expression between tumor progression groups. **I.** Heatmap of z-score normalized functional marker expression in DCIS tumors comparing tumor cells on the outer duct edge, tumor cells in the duct middle (duct mid), and tumor cells in the duct core.

226   demonstrated an intermediate functional profile, with features overlapping between pure

227   DCIS (GLUT1, CD36, COX2) and IBC (Ki67, pS6, HIF1α, MMP9) (Figure 3H).

9

228    It is not well understood how these functional states are affected by the location of
229    tumor cells within the duct of carcinoma *in situ*, where interior tumor cells far from the duct
230    edge may have limited access to nutrients and oxygen.  Interestingly, we found almost all
231    proliferative and cell signaling molecules to be enriched in tumor cells on the duct edge,
232    whereas HIF1$\alpha$ and metabolite import receptors GLUT1 and CD36 were enriched in cells
233    in the duct core, consistent with an adaptation to a low nutrient, hypoxic environment
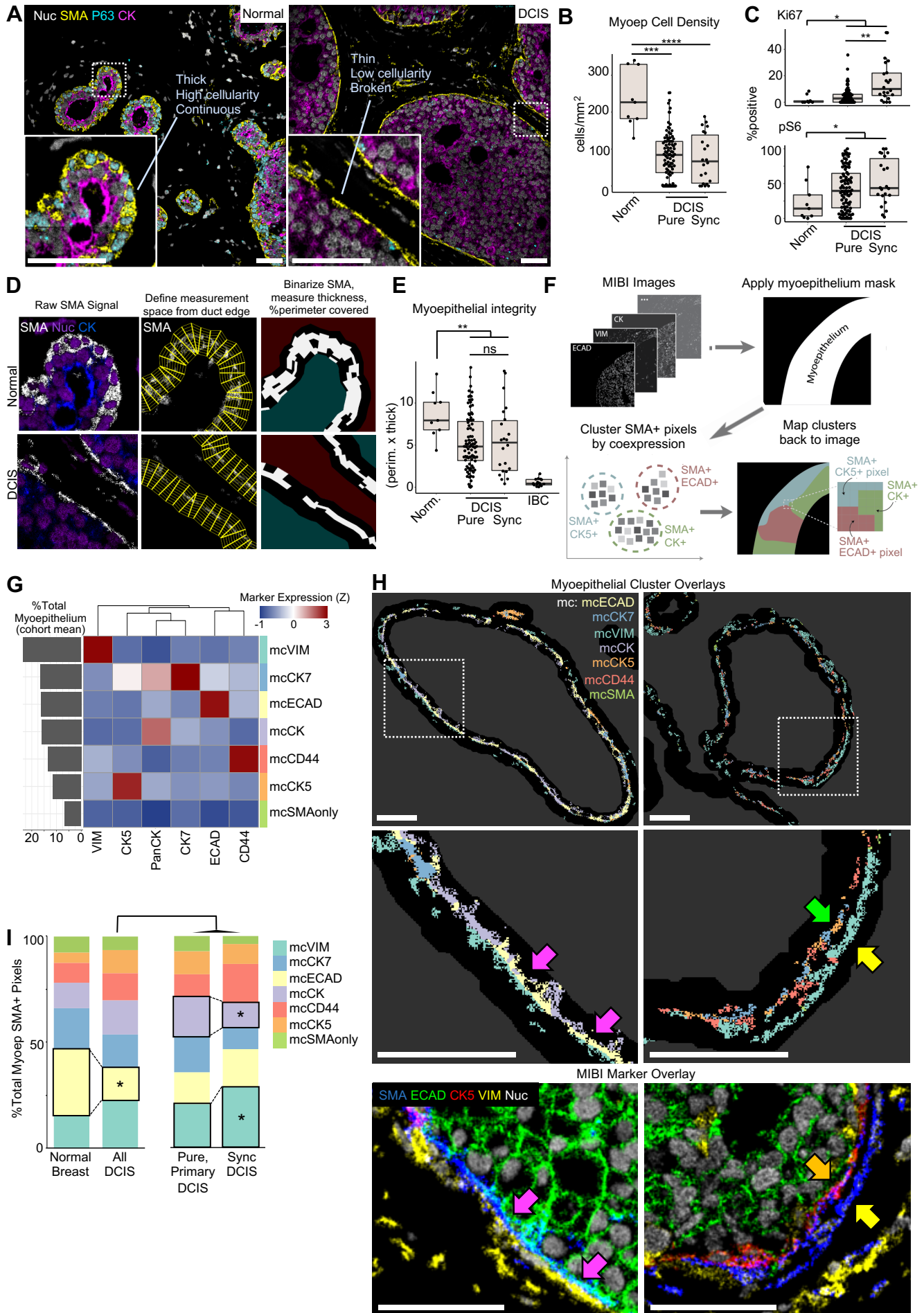234    (Figure 3I).

235

236

237    **Myoepithelial breakdown and phenotypic change during DCIS progression**
238    To understand how the structure and function of this key cellular barrier changes with
239    progression to IBC, we next performed a targeted analysis characterizing myoepithelial
240    cells which circumscribe both normal breast ducts and tumor cells in DCIS. Breast
241    myoepithelium in normal tissue is a thick, highly cellular layer between the stroma and
242    ductal cells (Figure 4A).  In DCIS, the myoepithelium is notably thinned out and reduced
243    in cellular density (Figure 4A-B).  The remaining myoepithelial cells in DCIS tumors were
244    found to have higher proliferation relative to normal tissue, with synchronous tumors
245    having the highest levels of the Ki67 positivity of these three groups (Figure 4C).

246    Given these findings, we hypothesized that loss of myoepithelial integrity
247    (thickness x percentage of duct-perimeter covered) in synchronous DCIS lesions would
248    also be greater than in pure DCIS.  To explore this question, we developed a new image
249    analysis tool to quantify myoepithelial thickness and percent coverage of the duct edge
250    (Figure 4D, see *Myoepithelial Coverage and Thickness Analysis* in Methods).  This
251    analysis revealed significant loss in myoepithelial integrity in DCIS tumors relative to
252    normal tissue.  To our surprise, however, no significant difference was observed between
253    pure and synchronous disease.  Thus, *in situ* tumorigenesis is accompanied by a
254    reduction of myoepithelial cell density and myoepithelial integrity independent of the
255    presence of a neighboring invasive component.

256    After quantifying these changes in myoepithelial structure, we next sought to
257    determine how the function of this regulatory barrier is altered with disease progression.

10

**Figure 4. Myoepithelial breakdown and phenotypic change during DCIS progression**

**A.** Representative MIBI image overlays showing SMA (yellow), p63 (cyan), and PanCK (magenta) expression in myoepithelium in normal breast (left) and DCIS (right), scale bars = 50μm. **B.** Myoepithelial cell density (cell/mm$^2$) was quantified in periductal regions is shown for normal breast, pure DCIS, and synchronous DCIS samples. **C.** The frequency of Ki67 (top) and pS6 (bottom) positivity is compared between groups as in *B*. **D.** Illustration of workflow for quantifying myoepithelial thickness and continuity. **E.** Boxplot showing myoepithelial integrity (percent coverage x average thickness) for normal tissue and patients with pure or synchronous DCIS. **F.** Workflow schematic for pixel-based clustering of myoepithelial phenotype. **G.** Heatmap showing frequency and average marker expression for 7 myoepithelial pixel clusters (mc) with a bar plot (left) of mc abundance out of total identified myoepithelium in the cohort. **H.** *Top.* Pseudo-colored image illustrating the spatial distribution of myoepithelial pixel clusters defined in *G* for a pure (left) and synchronous (right) DCIS tumor, scale bars = 50μm. *Middle.* Magnified periductal region with mcECAD (pink arrows), mcCK5 (orange arrow), and mcVIM (yellow arrow) areas denoted. *Bottom.* Coregistered color overlays showing variations in coexpression of SMA, ECAD, CK5, and VIM corresponding to pixel cluster assignments, scale bars = 50μm. **I.** Area plots comparing the frequency of each myoep cluster across normal breast, pure, and synchronous DCIS.

258    Due to their thin, elongated, and non-spherical cell bodies, myoepithelial cells are

259    inherently challenging to profile with classical nuclear-based segmentation approaches

260    which have been optimized for more conventional, ovoid cell shapes. Consequently,

261    outlines for myoepithelial cells predicted by these methods often extend significantly

262    beyond the true cellular border to erroneously include pixels from neighboring epithelial

263    and stromal cells. These errors propagate in downstream cell clustering analyses to

264    result in inaccurate phenotypic descriptions that are biased by what proteins are

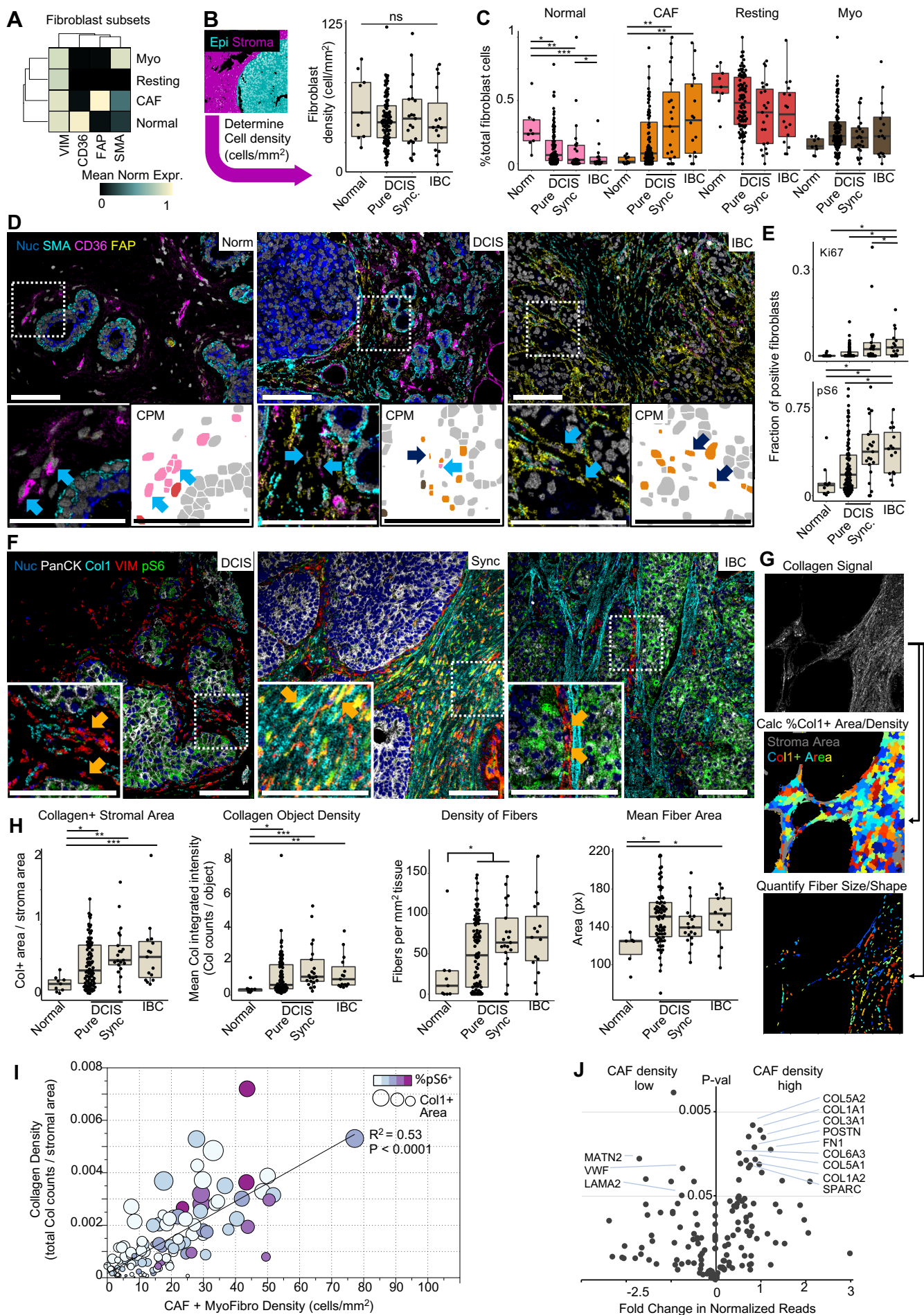265    expressed by closely approximated neighboring cells.

266    To avoid this pitfall, we created a new computational approach that assigns

267    phenotypes at the level of single pixels, rather than for whole cells (Figure 4F, see

268    *Myoepithelial Pixel Clustering Analysis* in Methods). This strategy yielded 7 distinct,

269    SMA$^+$ myoepithelial pixel clusters (mc) defined by coexpression of PanCK, ECAD, CK7,

270    CK5, VIM, or CD44, with SMA (Figure 4G). Mapping these pixel clusters back onto the

271    original images revealed that multiple expressional states can exist along the perimeter

272    of a single duct, from ECAD$^+$ and CK5$^+$ expression states often observed with apical

273    preference (Figure H, pink and green arrows), and more mesenchymal states that

274    exhibited a basal preference (e.g., VIM$^+$, CD44$^+$, yellow arrows). Notably, this analysis

275    also revealed a transition from a more luminal-like state in normal samples to a more

276    mesenchymal-like state in synchronous DCIS that aligned with analogous shifts in tumor

277    cell differentiation and function (Figure 4I).

278

279    **Fibroblast transition and collagen architecture remodeling during DCIS**

280    **tumorigenesis and progression**

281    In light of previous studies revealing a functional and structural interdependence between
282    myoepithelium and the surrounding stroma (Jones et al., 2003; Morsing et al., 2020), we
283    next sought to determine if the progressive loss of myoepithelial integrity observed here
284    correlated with changes in fibroblast function and extracellular matrix remodeling (ECM).
285    Single cell clustering revealed four fibroblast populations that included normal (CD36
286    high), resting (VIM-only), myofibroblast (SMA$^+$), and CAF (FAP$^+$) subsets (Figure 5A). No
287    significant differences in stromal cell density between progression groups were identified
288    when treating fibroblasts as a single cell population (Figure 5B). However, on comparing
289    the frequency of fibroblast subsets in normal tissue and DCIS, CAFs were found to
290    significantly increase across tumor progression as resting fibroblasts decreased (Figure
291    5C), with pure DCIS tumors having a heterogeneous mixture of these two states (Figure
292    5D, normal fibroblasts with light blue arrows, CAFs with dark blue arrows). A
293    corresponding increase in Ki67$^+$ fibroblasts suggests that this shift in identity is driven in
294    part by CAF proliferation (Figure 5E), which is accompanied by an increase in protein
295    translation (high pS6). We confirmed this relationship by comparing the CAF frequency
296    in samples with high and low pS6 and Ki67 (Figure S4A-B).

297    Given these findings, and that dense fibrillar collagen often appeared to be
298    juxtaposed with pS6$^+$ fibroblasts in progressed tumors (Figure 5F, orange arrows), we
299    next sought to determine how collagen remodeling was related to CAF location,
300    frequency, and phenotype. To achieve this, we developed new computational tools for
301    collagen morphometrics that were used to determine the shape, length, and density of
302    individual fibers (Figure 5G, see *Collagen Morphometrics* in Methods). These analyses
303    revealed that DCIS and IBC tumors had higher collagen density and longer fiber length
304    compared to normal breast (Figure 5H), suggesting that collagen deposition and fibrillar
305    remodeling were coordinated with the phenotypic shift to CAFs. Indeed, direct
306    comparison of collagen density and collagen-positive area to the density of CAFs and
307    myofibroblasts in the stroma revealed a strong correlation (Figure 5I). Furthermore, pS6$^+$
308    fibroblasts were also enriched in these collagen and CAF-dense tumors. Together these
309    data suggest a direct relationship between CAF activation and collagen deposition and
310    remodeling.

**Figure 5. Fibroblast transition and collagen architecture remodeling during DCIS tumorigenesis and progression**
**A.** Heatmap showing normalized marker expression for four fibroblast cell subsets: myofibroblasts (Myo), resting fibroblasts (Resting), cancer-associated fibroblasts (CAFs) and normal fibroblasts (Normal). **B.** *Left.* Example epithelial (cyan) and stromal (magenta) masks used to quantify stromal fibroblast density. *Right.* Boxplot of fibroblast density between tumor progression groups. **C.** Boxplots of fibroblast subset frequency across tumor progression groups. **D.** Representative MIBI image overlays showing normal, pure DCIS, and sync DCIS tumors with fibroblast markers. Zoomed insets (*left*) have paired cell phenotype maps (CPM, *right*) colored by fibroblast identity as in *C,* scale bars = 100µm. **E.** The frequency of Ki67 and pS6 positivity in fibroblasts is shown across progression groups. **F.** Representative MIBI image overlays showing VIM+ fibroblasts (red) with varying levels of pS6 expression (green) and nearby collagen 1 (Col1, cyan) deposition, scale bars = 100µm. **G.** Schematic showing the quantitation of MIBI collagen signal to identify %collagen+ stromal area, collagen density, and collagen fiber morphometrics. **H.** Collagen+ stromal area, collagen density, collagen fiber density (fibers/mm$^2$) and fiber area are quantified across tumor progression groups. **I.** Scatterplot comparing summed density of CAFs and myofibroblasts versus collagen density. Size and color of points are proportional to collagenized area and fibroblast pS6 positivity, respectively. **J.** Volcano plot of ECM-related gene expression for the top and bottom CAF-enriched DCIS tumors.

311  Finally, to identify which specific collagen isoforms correlate with this activity and

312 to determine if additional ECM proteins are involved, we compared ECM transcript levels

313 in stroma of CAF-high- and low-density tumors using LCM RNAseq. We found the

314 majority of collagen species were upregulated in CAF-high tumors with COL5A2 and

315 COL1A1 being the most significant of these, consistent with MIBI-TOF quantitation of

316 COL1A1 protein (Figure 5J). In addition, CAF-dense tumors showed increased deposition

317 of fibronectin (FN1), SPARC and periostin (POSTN), indicative of CAF-remodeling and a

318 shift towards a pro-invasive stroma (Barth et al., 2005; Malanchi et al., 2012).

319

320 **Characterizing the preinvasive immune microenvironment and its compartmental**

321 **evolution throughout progression**

322 Having identified coordinated shifts in tumor differentiation, myoepithelial integrity, and

323 fibroblast function, we next sought to understand how immune composition changed with

324 disease progression. We found monocytes, mast cells, and HLA-DR+ antigen presenting

325 cells (APCs) to be the most abundant immune cells in pure DCIS (Figure 6A). Immune

326 cells were typically found in the stroma and were occasionally embedded in ducts (Figure

327 B, orange arrow). To quantify the spatial distribution of immune cells in these

328 compartments, we interrogated cell density in epithelial and stromal mask regions (Figure

329 6C). This analysis identified a clear stromal preference when treating immune cells as a

330 single population (Figure 6D, S5A). To understand if this preference remained valid when

331 considering specific subsets of lymphoid and myeloid cells, we compared the local

332 frequency within stromal and ductal regions for each cell type. CD4+ T cells, B cells,

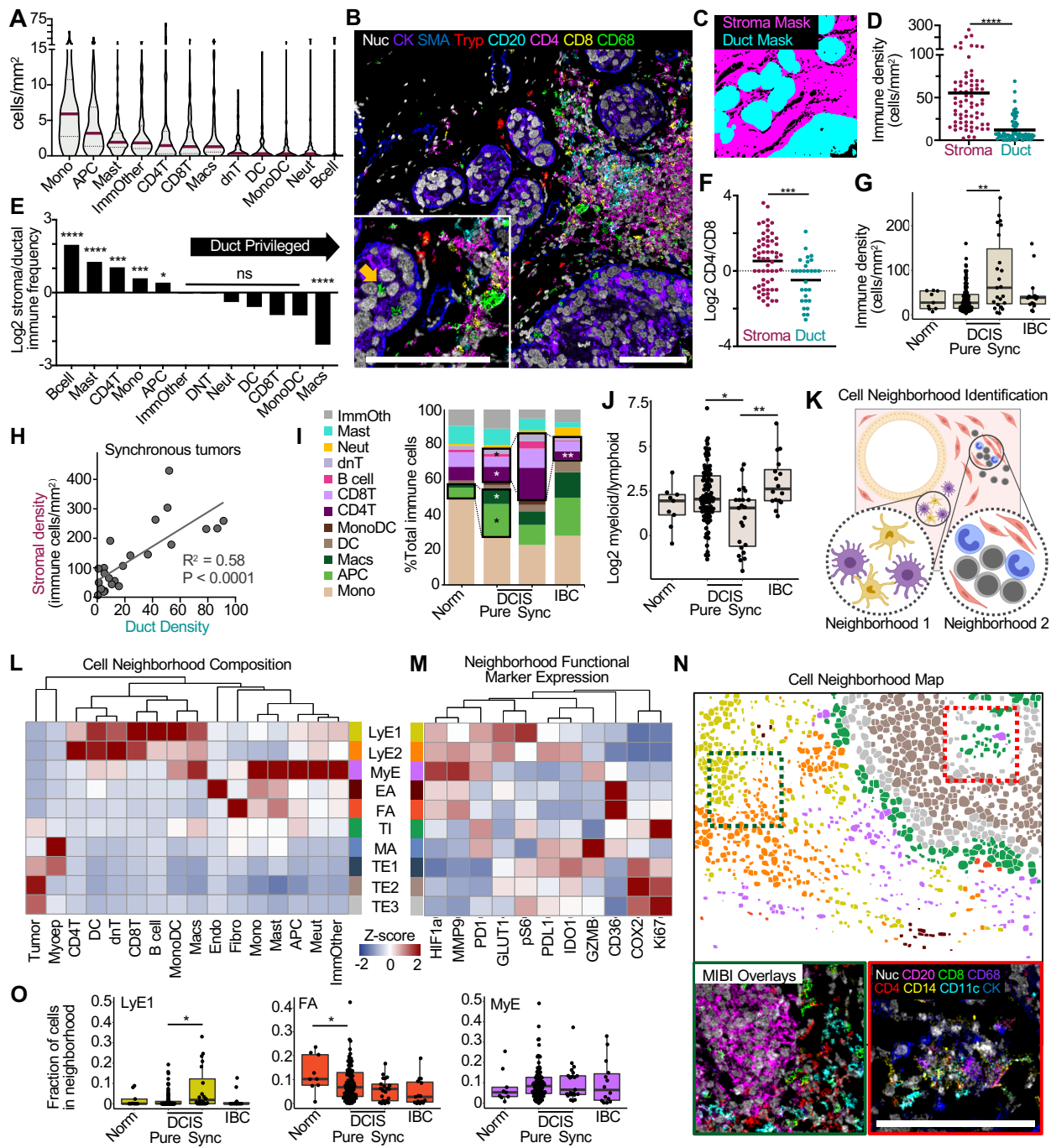333 monocytes, APCs and mast cells all demonstrated a statistically significant stromal

334    preference, while macrophages were significantly enriched in ductal regions (Figure

335    6E). Interestingly, differential enrichment of CD4$^+$ and CD8$^+$ T cells resulted in a

336    CD4/CD8 ratio that skewed towards CD8$^+$ T cells in ducts and CD4$^+$ T cells in stroma

337    (Figure 6F).

338        We next investigated how immune cell prevalence and spatial enrichment evolves

339    with transition from *in situ* tumorigenesis to invasive disease by comparing pure DCIS

340    with synchronous lesions and IBC. Immune cell density was significantly increased in

341    synchronous lesions compared to all other groups (Figure 6G). Notably, this increase in

342    immune infiltrate was present in both the stroma and ducts of these lesions (Figure 6H),

343    suggesting a coordinated influx into the ducts during increased stromal immune

344    infiltration. By comparing the cell density for each immune cell subset with respect to

345    disease stage, we observed an increase in effector myeloid cells (Macs, APC) in pure

346    DCIS compared to normal breast (Figure 6I). Importantly, this also revealed the increase

347    in immune infiltrate in synchronous tumors to be driven primarily by an influx of B and T

348    lymphocytes (Figure 6I, S5B), resulting in an immune microenvironment more skewed

349    towards lymphocytes (Figure 6J). Subsequently, both T cell frequency and myeloid to

350    lymphoid ratio in IBC tumors return to values similar to pure DCIS.

351        In order to better understand how this feature and other immune programs were

352    spatially organized, we applied a K-means clustering approach to identify distinct cellular

353    neighborhoods (CNs), where a CN is defined by a set of cell types found to spatially co-

354    occur across the cohort (Figure 6K, see *Protein and Cellular Spatial Enrichment Analyses*

355    in Methods). Through this approach, we identified 10 CNs that we categorized as being

356    lymphocyte-enriched (LyE1, LyE2), myeloid-enriched (MyE), endothelial-associated

357    (EA), fibroblast-associated (FA), myoepithelial-associated (MA), tumor-interface (TI), and

358    tumor-enriched (TE1-3, Figure 6L-N).

359        Interestingly, single cell expression of functional markers was found to be

360    correlated with CN, even though these parameters were not included in the K-means

361    neighborhood assignment analysis. For example, HIF1$\alpha$ and MMP9 expressing cells

362    were enriched in MyE, while the frequency of pS6$^+$ cells was highest in LyE1 (Figure

363    6L). Macrophages were a constituent of numerous CNs and showed functional state

364    distinction based on neighborhood association, including increased PDL1 expression

**Figure 6. Characterizing the preinvasive immune microenvironment and its compartmental evolution throughout progression**

**A**. Violin plot examining immune cell density in pure DCIS, ranked by median density per patient. **B**. Representative MIBI image overlay of a pure DCIS tumor with major immune cell type markers, inset and arrow highlighting intraductal immune phenotypes. **C.** Mask overlay showing delineation of stroma and duct regions in *B,* scale bars = 100μm. **D**. Scatterplot comparing immune cell density between the stroma and duct compartments per patient. **E**. Column plot showing the ratio (Log2) of immune cell type frequency between stroma and ductal compartments, ranked from high (stromal preference) to low (duct preference). Asterisks denote significance comparing compartment frequency of a given cell type across all pure DCIS patients. **F**. Log2 ratio of CD4+ to CD8+ T cells is displayed per patient for the stroma and duct compartments. **G**. Whole image immune density is compared across tumor progression groups. **H**. Scatterplot comparing stromal and ductal immune density per patient in synchronous tumors. **I**. Area plot showing the change in immune subset frequency across progression groups. Effector-myeloid cell subsets are boxed and compared between normal breast and pure DCIS tumors; asterisks denote significant differences in frequency. Lymphocyte subsets are boxed and compared between pure DCIS, synchronous DCIS, and IBC, asterisks denote significance vs the synchronous group. **J.** Boxplots showing the log2 ratio of myeloid to lymphoid cells in tumor progression groups. **K.** Illustration depicting different spatially-enriched cellular neighborhoods. **L.** Heatmap showing z-score normalized cell type frequency for each cellular neighborhood: lymphocyte-enriched (LyE1, LyE2), myeloid-enriched (MyE), endothelial-associated (EA), fibroblast-associated (FA), tumor-interface (TI), myoepithelial-associated (MA), and tumor-enriched (TE1-3). **M.** Heatmaps showing z-score normalized mean expression for functional markers in each cellular neighborhood. **N.** *Top.* Cell neighborhood map showing the spatial localization of distinct neighborhoods, denoted by color as in *M. Bottom.* Color overlays for lymphocyte-enriched (green dotted line) or tumor-interface (red dotted line), scale bar = 100μm. **O.** Boxplot showing frequency of cells assigned to LyE1 (yellow), Fibroblast-associated (Red) and MyE (purple) cell neighborhoods across tumor progression groups.
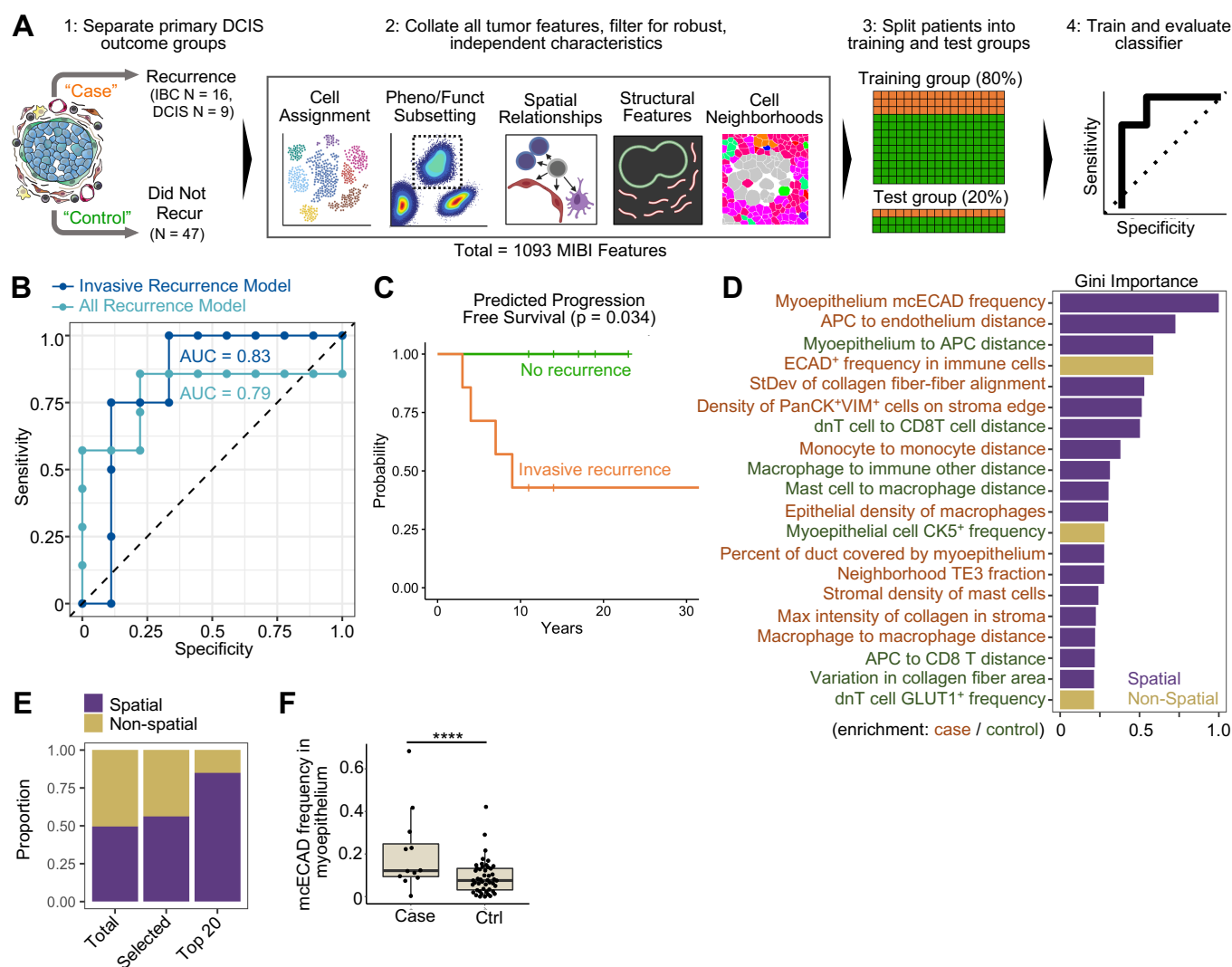
365 within LyE neighborhoods, in addition to pS6 (Figure S5C-D). Notably, the LyE1

366 neighborhood was also enriched for T and B cells, consistent with tertiary lymphoid

367 structure formation (see Figure 6N, bottom left). In line with the trends observed for T cell

368 infiltrates, we found the frequency of cells belonging to LyE1 to be increased in

369 synchronous lesions (Figure 6O). Taken together, these findings indicate that early

370 stromal invasion in synchronous tumors triggers an influx of T cells and formation of TLS

371 structures. We find that by IBC, however, the tumor immune microenvironment has

372 reverted to a myeloid-skewed, immunosuppressed state with diminished T cell

373 presence.

374

375 **Identifying DCIS features correlated with recurrence outcomes**

376 Having extensively quantified the multi-compartmental cellular and structural elements of

377 DCIS tumors, we leveraged these data to identify features associated with the risk of

378 recurrence following primary DCIS resection. We selectively examined these features in

379 diagnostic tissue procured at the time of initial presentation in two sets of patients. The

380 first set, referred to as "case", consisted of 31 patients who had a recurrence (DCIS or

381 IBC) within 2-15 years of being treated for newly diagnosed pure DCIS. The second set,

382 referred to as "control", consisted of 47 patients with pure DCIS that did not recur within

383 11+ years.

384 Using these outcome groups and 1,093 phenotypic, functional, spatial, and

385 morphologic features extracted from our MIBI-TOF analyses (Table S3), we trained two

386 random forest classifier models. The first was an all-recurrence model for predicting

387 which patients would have a recurrence of DCIS *or* IBC. The second was an invasive

388 recurrence model for predicting IBC recurrence *exclusively* (Figure 7A). Low observation

389 and overly correlated features were removed from the dataset and the patient population

390 was randomly split 80/20 to training and test groups. We evaluated classifier accuracy in

391 the withheld test set, where the all-recurrence and invasive models achieved an AUC of

392 0.79 (CI 0.51:1) and 0.83 (CI 0.59:1), respectively (Figure 7B). When stratifying patients

393 by their predicted labels, we found a significant difference in recurrence probability over

394 time (Fig. 7C, Figure S6A), with no recurrence events in the patients predicted by the

395 invasive model to be non-progressors. Although sample size precluded us from being

**Figure 7. Identifying DCIS features correlated with recurrence outcomes**

**A.** Schematic illustrating the different outcome groups of primary DCIS including "cases" that recurred either as IBC or DCIS, and "controls" with no recurrence in >11yr follow-up. 1,093 MIBI features of numerous tumor metrics were used to train a random forest classifier to differentiate case and control samples. Classifier specificity was then tested on a withheld 20% of patients. **B.** AUC plot showing classifier sensitivity and specificity. **C.** Predicted survival of patients identified in the test set of the invasive-recurrence model as case or control. **D.** MIBI features with top classifier importance for the IBC recurrence model are shown, ranked by Gini importance. Features are colored based on enrichment either in cases (orange) or controls (green), importance bars are colored based on the feature utilizing spatial information (purple) or not (gold). **E.** The distribution of spatial vs non-spatial features are shown for all features identified (total), those used by the model (selected), and those in the top 20 most important features (top 20). **F.** Boxplot showing the frequency of the mcECAD myoepithelial phenotype between invasive cases and controls.

396     able to eliminate patient demographics and differences in clinical therapy as a confounder

397     in this analysis, treatment regimens known to affect recurrence rates (i.e., mastectomy,

398     radiation, tamoxifen) were well distributed between the case and control patients (Figure

399     S6B). Likewise, no significant difference in classifier predictions were identified with

400     respect to these variables (Figure S6C).

401 To understand the biology being leveraged by this classifier to accurately

402 discriminate pre-invasive from indolent DCIS tumors, we ranked the top 20 features

403 based on Gini importance.  These features primarily consisted of metrics related to the

404 phenotype of myoepithelium, the structure of collagen fibers in the extracellular matrix,

405 and the spatial distribution of multiple immune cell subsets (Figure 7C).  Notably, spatial

406 metrics describing cell densities, cell neighborhoods, pairwise cell distances, collagen

407 structure, and multiplexed subcellular features were overrepresented and accounted for

408 17 of the top 20 metrics in the invasive model (Figure 7D, Table S3).  Immune cell metrics

409 comprised about half of these and were myeloid skewed (Figure S6D, with 9 relating

410 specifically to myeloid subsets and 3 to lymphoid subsets.  Similarly, enrichment for

411 spatial metrics related to myoepithelium, collagen, and myeloid cells were observed in

412 the all-recurrence model as well (Figure S6D-F).  Stromal density of  $PanCK^+VIM^+$ cells

413 ranked in the top 20 features.  These cells were rare (median of 0 in case and controls)

414 and on manual inspection appeared to represent fibroblasts where PanCK expression

415 from closely neighboring epithelial cells was misassigned.  Interestingly, both models

416 identified pixel-level, $ECAD^+$ myoepithelial expression as the most predictive metric

417 (mcECAD, see Figure 4).  When comparing case and control samples, we found the

418 frequency of this feature to be significantly different between these outcome groups,

419 independent of the classifier model, and to be readily identifiable on targeted inspection

420 of the original imaging data ($p < 0.001$, Figure 7F).

421

422

423 **Discussion**

424 Here, we report the first multicompartmental atlas of the single cell composition and

425 structure of DCIS.  The central focus of this study was to characterize the changes

426 undergone with progression to IBC where tumor cells breach the duct to invade the

427 surrounding stroma.  Previous work examining BC progression have attempted to

428 attribute this transition either to tumor-intrinsic factors or to specific features of stromal

429 cells in the surrounding TME.  By simultaneously mapping both tumor and stromal cell

430 identity and function in intact human tissue, we sought to treat the DCIS TME as a single

431 ecosystem where progression to invasive disease depends on the spatial distribution and
432 function of multiple cell types, rather than on any single cell subset.

433       Meeting this goal required first assembling a large, well-annotated, and diversified
434 pool of human DCIS tissue: the RAHBT cohort. This effort was motivated in part by the
435 success of similar work investigating invasive disease (i.e. METABRIC) that have
436 provided deep insights into breast tumor composition and have served as authoritative
437 resources in breast cancer research (Curtis et al., 2012). To achieve this, the Breast
438 PreCancer Atlas constructed a unique set of archival human surgical resections that
439 captured the full spectrum of breast cancer progression, from normal tissue, to pure DCIS
440 and IBC. Assembling all of these cases into TMAs has enabled a one-of-a-kind workflow
441 for multiomics analyses where genomic, transcriptomic, and proteomic techniques are
442 performed not only on the same samples, but on coregistered serial sections of the same
443 local region of tissue.

444       Here, we describe the first major analysis of the RAHBT cohort where high
445 dimensional imaging was used to characterize BC progression. We used MIBI-TOF for
446 subcellular imaging of 140 tumor and normal breast samples using a 37-marker staining
447 panel (122 and 23 samples from RAHBT and Stanford cohorts, respectively). Tumor cell
448 differentiation and function were found to transition along a continuum from pronounced
449 luminal features in normal breast to a more undifferentiated, cytokeratin-low state in
450 invasive disease that had increased mesenchymal features. This shift was accompanied
451 by an upregulation of HIF1$\alpha$, MMP9, and IDO in tumor cells, which have been shown to
452 directly elicit EMT, promote invasion, and drive immune tolerance, respectively (Kolijn et
453 al., 2018; Lin et al., 2011; Peng et al., 2018; Zhang et al., 2015, 2019). With transition to
454 DCIS, the frequency of an E-cadherin-high myoepithelial phenotype that predominated
455 normal breast tissue decreased, as a more mesenchymal, CD44- and VIM-high state
456 increased. Interestingly, no difference in myoepithelial cell density or structural integrity
457 was found when comparing DCIS in pure and synchronous lesions. Given that the
458 invasive and *in situ* components of synchronous tumors are closely related on a genomic
459 level (Ak et al., 2018; Kim et al., 2015; Newburger et al., 2013) these findings suggest
460 that transition to invasion disease is regulated at least in part by the local
461 microenvironment.

18

462    These epithelial changes were accompanied by a stromal transition towards higher

463    numbers of activated, proliferating CAFs and densely aligned fibrillar collagen (Conklin et

464    al., 2011; Esbona et al., 2018).  Although the total immune density was comparable to

465    normal breast tissue, DCIS tumors exhibited a shift from a monocyte-predominant

466    environment to one enriched for APCs and intraductal macrophages.  In line with recent

467    findings by other groups (Alcazar et al., 2017; Kim et al., 2020) synchronous DCIS/IBC

468    tumors were marked by a stromal spike in T and B cells and formation of tertiary lymphoid

469    structures.   This feature distinguishes them from the myeloid-skewed IBC samples

470    profiled in this study.  Taken together, these findings support a model for breast cancer

471    progression where invasive disease occurs through multiple coordinated, dynamic

472    interactions of the surrounding stroma, myoepithelium, and tumor.

473

474    Given the urgent need to better stratify DCIS patients based on risk of progression,

475    we tested to see if these spatial and phenotypic features could be used to predict IBC

476    recurrence based exclusively on diagnostic DCIS tissue.  Using 1,093 features, we

477    trained a random forest classifier model for identifying patients that would later progress

478    to IBC that achieved an AUC of 0.83 on withheld test samples.  Although the performance

479    was impressive, certain caveats should be taken into account when considering how

480    generalized this model might be. Given the complexity of breast cancer subtypes and the

481    impact of patient demographics on outcome (Alaeikhanehshir et al., 2020; Liu et al.,

482    2019), the sample size in this study may not have been sufficient to fully account for the

483    confounding effects of these variables.  Lastly, since all patients in the RAHBT cohort

484    received one or more therapeutic interventions, the features leveraged by this model to

485    identify non-progressors might not be valid when applied to patient populations where

486    therapy is omitted.

487    With these considerations in mind however, these results do offer three compelling

488    and overarching insights.  First, spatial metrics relating phenotype to structure and

489    morphology were significantly over-represented relative to non-spatial metrics,

490    accounting for almost 85% of the top 20 features identified by the classifier model.

491    Second, the most influential features were primarily related to the stroma rather than the

492    tumor cells themselves.   This included a previously unreported E-cadherin high

493    myoepithelial phenotype as well as collagen fiber size and alignment with respect to the

494    duct.  Third, high ranking immune features more often related to myeloid than to lymphoid

495    subsets, particularly those in close proximity with myoepithelium or residing inside the

496    duct.   This skewing underscores the need to better understand how macrophages

497    promote TME immune suppression, tumor proliferation, and local invasion (Esbona et al.,

498    2018; Goswami et al., 2005; Linde et al., 2018; Ruffell et al., 2012).

499         Taken together, this study offers a comprehensive, multi-compartmental atlas of

500    preinvasive breast cancer that illustrates the full continuum of tissue structure and

501    function starting from a homeostatic state in normal breast through *in situ* and invasive

502    disease. Combining this comprehensive data set with extensive patient follow-up has

503    enabled identification of tumor features that are associated with DCIS recurrence and

504    offers a framework for exciting follow-on efforts. With this in mind, we are actively planning

505    a larger study that will further evaluate the biological significance of spatial features

506    relating to myoepithelium, collagen, and myeloid cells and to determine if they can be

507    used to prospectively risk stratify patients with a new DCIS diagnosis.

508

509

510    **Methods**

511

512    **Patient Cohort**

513    We utilized a retrospective study cohort of patients from the Washington University

514    Resource of Archival Tissue (RAHBT) that contained two outcome groups: controls

515    ("Ctrl") composed of patients with DCIS who had no recurrence and cases ("Case")

516    composed of patients with DCIS who had either a DCIS or an IBC recurrence.  For each

517    case, we matched two controls who remained free from recurrent lesions, based on age

518    at diagnosis (+/- 5 years), and type of definitive surgery (mastectomy or lumpectomy). For

519    each DCIS diagnosis we retrieved primary and recurrent tumor slides and blocks for

520    pathology review, secured a whole slide image of each sample, marked for TMA cores,

521    and generated TMA blocks with 84 1.5mm cores, including additional tonsil and normal

522    breast controls.

523       Supplemental table 1 summarizes the data for the cases in the cohort. Median age
524    at diagnosis was 54, year of diagnosis was 1986 to 2017, and time to recurrence with
525    was 8.8 years for invasive lesions, and 5.3 years for premalignant lesions. For women in
526    the cohort with no recurrence, follow up extended to 132 months, on average. Treatment
527    of initial DCIS ranged from lumpectomy with radiation (approximately half of cases), and
528    lumpectomy with no radiation (20%) and mastectomy with no radiation for 30%. The
529    RAHBT cohort is composed of African American women (26%) and white women (74%).
530    We also profiled a supplemental cohort of patients from the Stanford Hospital with
531    synchronous ("Sync") DCIS and IBC tumors from 2007-2009. A 216-core TMA block was
532    generated with 1mm tumor cores, with additional tissue controls.

533       5µm serial sections of each TMA slide were cut onto glass slides for hematoxylin
534    and eosin (H&E) staining, onto laser-capture slides for LCM-RNAseq (SMART-3SEQ)
535    and cut onto gold- and tantalum-sputtered slides for MIBI-TOF imaging. H&E slides were
536    inspected by a breast cancer pathologist to address DCIS purity and demarcate regions
537    of DCIS to guide MIBI imaging and laser dissection of epithelial and stromal area. The
538    Stanford Hospital cohort was without paired LCM-RNAseq analysis.

539

540    **Antibody Preparation**
541    Antibodies were conjugated to isotopic metal reporters as described previously (Keren et
542    al., 2018; McCaffrey et al., 2020). Following conjugation antibodies were diluted in Candor
543    PBS Antibody Stabilization solution (Candor Bioscience). Antibodies were either stored
544    at $4^O$C or lyophilized in 100 mM D-(+)-Trehalose dehydrate (Sigma Aldrich) with ultrapure
545    distilled H2O for storage at -20oC. Prior to staining, lyophilized antibodies were
546    reconstituted in a buffer of Tris (Thermo Fisher Scientific), sodium azide (Sigma Aldrich),
547    ultrapure water (Thermo Fisher Scientific), and antibody stabilizer (Candor Bioscience) to
548    a concentration of 0.05 mg/mL. Some metal-conjugated antibodies in this study were
549    used as secondary antibodies, targeting hapten groups on hapten-conjugated primary
550    antibodies, this included the pairs PDL1-Biotin and Anti-Biotin[149Sm], and ER-Alexa488 and
551    Anti-Alexa488[142Nd]. Information on the antibodies, metal reporters, and staining
552    concentrations is located in Table S2.

553

**Tissue Staining**

Tissues were sectioned (5µm section thickness) from tissue blocks on gold and tantalum-sputtered microscope slides. Slides were baked at 70ºC overnight followed by deparaffinization and rehydration with washes in xylene (3x), 100% ethanol (2x), 95% ethanol (2x), 80% ethanol (1x), 70% ethanol (1x), and ddH2O with a Leica ST4020 Linear Stainer (Leica Biosystems). Tissues next underwent antigen retrieval by submerging sides in 3-in-1 Target Retrieval Solution (pH 9, DAKO Agilent) and incubating at 97ºC for 40 minutes in a Lab Vision PT Module (Thermo Fisher Scientific). After cooling to room temperature slides were washed in 1x PBS IHC Washer Buffer with Tween 20 (Cell Marque) with 0.1% (w/v) bovine serum albumin (Thermo Fisher). Next, all tissues underwent two rounds of blocking, the first to block endogenous biotin and avidin with an Avidin/Biotin Blocking Kit (Biolegend). Tissues were then washed with wash buffer and blocked for 1 hour at room temperature with 1x TBS IHC Wash Buffer with Tween 20 with 3% (v/v) normal donkey serum (Sigma-Aldrich), 0.1% (v/v) cold fish skin gelatin (Sigma Aldrich), 0.1% (v/v) Triton X-100, and 0.05% (v/v) Sodium Azide. The first antibody cocktail was prepared in 1x TBS IHC Wash Buffer with Tween 20 with 3% (v/v) normal donkey serum (Sigma-Aldrich) and filtered through a 0.1µm centrifugal filter (Millipore) prior to incubation with tissue overnight at 4ºC in a humidity chamber. Following the overnight incubation slides were washed twice for 5 minutes in wash buffer. The second day antibody cocktail was prepared as described and incubated with the tissues for 1 hour at 4ºC in a humidity chamber. Following staining, slides were washed twice for 5 minutes in wash buffer and fixed in a solution of 2% glutaraldehyde (Electron Microscopy Sciences) solution in low-barium PBS for 5 minutes. Slides were washed in PBS (1x), 0.1 M Tris at pH 8.5 (3x), ddH2O (2x), and then dehydrated by washing in 70% ethanol (1x), 80% ethanol (1x), 95% ethanol (2x), and 100% ethanol (2x). Slides were dried under vacuum prior to imaging.

**MIBI-TOF Imaging**

Imaging was performed using a MIBI-TOF instrument with a Hyperion ion source. $Xe^+$ primary ions were used to sequentially sputter pixels for a given FOV. The following imaging parameters were used: Acquisition setting: 80 kHz, Field size: 500 µm$^2$, 1024 x

22

585    1024 pixels, dwell time: 5ms, median gun current on tissue: 1.45nA Xe⁺, ion dose: 4.23

586    nAmp hours / mm$^2$ for 500 μm$^2$ FOVs.

587

**Low-level Image Processing and Single Cell Segmentation**

589    Multiplexed image sets were extracted, slide background-subtracted, denoised, and

590    aggregate filtered as previously described (Keren et al., 2018; McCaffrey et al., 2020).

591    Nuclear segmentation was performed using an adapted version of the DeepCell CNN

592    architecture (McCaffrey et al., 2020; Valen et al., 2016). To more effectively capture the

593    range of cell shapes and morphologies present in DCIS, we generated two distinct

594    segmentations for each image. The first used a radial expansion of three pixels and a

595    stringent threshold for splitting cells (See Figure S2A, *Stroma Parameters*). The second

596    used a radial expansion of one pixel and lenient threshold for splitting cells (*Epithelial*

597    *Parameters*). We combined these masks together using a post-processing step which

598    gave preference to the epithelial segmentation mask, overriding and stromal-mask-

599    detected objects in the same area. Smaller cells identified by the stromal settings and

600    missed in the epithelial settings were combined to the final cell mask. A cell nuclei ("Nuc")

601    channel combining HH3 and endogenous phosphorous (P) signal was made to increase

602    signal robustness for nuclei detection.

603

**Single Cell Phenotyping and Composition**

605    Single cell data was extracted for all cell objects and area normalized. Single cell data

606    was linearly scaled by average cell area across the cohort and asinh-transformed with a

607    co-factor of 5. All mass channels were scaled to 99.9th percentile. In order to assign each

608    cell to a lineage, the FlowSOM clustering algorithm was used in iterative rounds with the

609    Bioconductor "FlowSOM" package in R (Van Gassen et al., 2015). The first clustering

610    round separated cells into 100 clusters that were subsequently merged into one of five

611    major cell lineages (tumor, myoepithelial, fibroblast, endothelial, immune) based on the

612    clustering nodes. Proper lineage assignments were ensured by overlaying Flowsom

613    cluster identity with lineage-specific markers. Supervised lineage reassignment was

614    performed where needed. Immune cells were subclustered again to delineate B cells,

615    CD4⁺ T cells, CD8⁺ T cells, monocytes, MonoDC cells, DC cells, macrophages,

616    neutrophils, mast cells, double-negative CD4$^-$CD8$^-$ T cells (dnT cells), and HLADR$^+$ APC

617    cells. CD45$^+$-only immune cells were annotated as 'immune other.' Tumor and fibroblast

618    cells were similarly clustered again to reveal phenotypic subsets, as shown in Figure S2.

619    Altogether, we assigned 94% (n = 127,451 of 134,631) of cells to 16 subsets, with the

620    remaining nucleated cells with absent or very low levels of lineage markers assigned as

621    "other". The relative abundance of all major lineages was determined out of total cells per

622    FOV and the relative frequency of cell subsets were determined out of total cells of a

623    given lineage, per FOV.

624

625    **Region Masking**

626    Region masks were generated to define histologic regions of each FOV including the

627    epithelium, stroma, myoepithelial (periductal) zone, and duct, which was further

628    subdivided into the duct edge, duct mid, and duct core. We removed gold-positive area

629    which marked regions of bare slide from holes in the tissue, providing an accurate

630    measurement of tissue area.  This area measurement could be used to calculate cellular

631    density in specific histologic regions, e.g., fibroblast density in the stroma, which was

632    critical to normalize the observed cell abundances by how much tissue of a specific type

633    was sampled, and prevent bias based on how much tumor vs stroma the FOV covered.

634    The epithelial mask was first generated though merging ECAD and PanCK signal and

635    applying smoothing and radial expansion to incorporate the myoepithelial zone, and the

636    inside of ducts were filled.  The stromal mask included all image area outside of the

637    epithelial mask. Duct masks were generated through the erosion of the epithelial masks

638    by 25 pixels.  The myoepithelial mask was generated by subtracting the duct mask from

639    the epithelial mask.  Duct edge, duct mid, and duct core masks (Figure 3I) were generated

640    by eroding the duct mask by subsequent 100-pixel increments.

641

642    **Protein and Cellular Spatial Enrichment Analyses**

643    A spatial enrichment approached was used as previously described (Keren et al., 2018,

644    2019; McCaffrey et al., 2020) to identify patterns of protein enrichment or exclusion across

645    all protein pairs. HH3 was excluded from the analysis. For each pair of markers, X and Y,

646    the number of times cells positive (normalized expression >0.25) for protein X was within

647  a ~50 um radius of cells positive for protein Y was counted. A null distribution was

648  produced by performing 100 bootstrap permutations where the locations of cells positive

649  for protein Y were randomized. A z-score was calculated comparing the number of true

650  cooccurrences of cells positive for protein X and Y relative to the null distribution.

651  Importantly, symmetry is assumed: the values of when calculating the spatial enrichment

652  of protein X close to protein Y are the same as with protein Y close to protein X. For each

653  pair of proteins X and Y the average z-score was calculated across all DCIS FOVs.

654  To analyze cellular associations with the myoepithelium, the distances between all

655  cell centroids to the nearest perimeter location of the myoepithelium mask (described

656  above) were calculated. To quantity cell type spatial interactions, the mean distances

657  between cell centroids for all cell phenotype pairs (self-self pairs excluded) were

658  calculated per region.

659  Cell neighborhoods were produced by first generating a cell neighbor matrix,

660  where each row represents an index cell, and the columns indicate the relative frequency

661  of each cell phenotype within an 36um radius of the index cell. Next the neighbor matrix

662  was clustered to 10 clusters using k-means clustering. Neighborhood cellular profile was

663  determined by assessing the mean prevalence of each cell phenotype in the index cells'

664  36um radius, while functional marker expression was determined by assessing mean

665  marker expression by the index cells assigned to each neighborhood cluster.

666

667  **DCIS UMAP Visualization**

668  UMAP embeddings were determined for all DCIS tumors (pure, synchronous, primary

669  and recurrent) using the R implementation (McInnes et al., 2020) with the following

670  parameters: n_neighbors = 15, min dist = 0.1 and the following markers: PanCK, CK7,

671  CK5, ECAD, VIM, ER, HER2, AR, CD31, SMA, CD45, HLADR, CD68, CD11c, CD14,

672  CD20, CD3, CD4, CD8, MPO, Tryptase.

673

674  **EMT GSEA**

675  To identify genes and pathways associated to EMT, MIBI-identified DCIS vimentin high

676  vs low samples were selected, and the epithelial fraction of an adjacent tissue section

677  was analyzed by LCM-RNAseq (Vim high, n = 26; Vim low, n = 32). DESeq2 R package

678   (version 1.30.0) was used for data normalization and differential expression analysis.

679   Results were sorted by decreasing log fold change and the ranked list was subjected to

680   GSEA against C2 curated dataset of molecular signature database

681   (MSigDB)(Subramanian et al., 2005). P values were corrected for multiple comparisons

682   by using Benjamini-Hochberg method and terms with p adj < 0.05 were considered.

683

684   **ECM Gene Analysis**

685   To analyze extracellular matrix components by gene expression, an extracellular matrix

686   gene signature (GO extracellular matrix structural constituent, GO:0030021) was

687   downloaded from GSEA website and used to compare MIBI-identified samples with the

688   top and bottom quartiles of cancer associated fibroblast density in the stroma. Stromal

689   LCM-RNAseq samples were used for this analysis. Raw reads were normalized with

690   DESeq2 R package (version 1.30.0)(Anders and Huber, 2010) and a paired T-test was

691   compared to the log2 ratio of group means to generate the volcano plot.

692

693   **Myoepithelial Continuity and Thickness Analysis**

694   To define a window of myoepithelial signal quantitation, we used a topology-preserving

695   operation to define a curve 5 pixels out from the epithelial mask edge (see *Region*

696   *Masking*) and a curve 30 pixels in from the epithelium mask edge, and we defined those

697   pixels in between these two curves as the myoepithelium mask. We subdivided the outer

698   curve into 5-pixel long arc-segments, and for each point on the outer edge in between

699   two segments, found the nearest point on the inner edge, dividing the myoepithelium into

700   a string of quadrilaterals or "wedges". Wedges are then subdivided each wedge along the

701   in-out (of the epithelium) axis into 10 segments. Wedges are merged when both their

702   combined inner and outer edges has an arc-length less than 15 pixels.

703       We took pre-processed (background subtracted, de-noised) SMA pixels within the

704   mesh and smoothed them with a Gaussian blur of radius of 1. We then calculated the

705   density of SMA signal within each mesh-segment as the mean pixel value of smoothed

706   SMA within that mesh-segment. This density was then binarized to create a SMA-

707   positivity mesh, using a threshold of 0.5 (density > 0.5 as positive).

26

708    The percentage of duct perimeter covered by myoepithelium was calculated by

709    assigning an "SMA-present" variable to each wedge, "0" if no mesh-segments in the

710    wedge were positive for SMA, and "1" otherwise. Each wedge is weighted by its area

711    relative to the myoepithelium area. The sum over all wedges of the product of the "SMA-

712    present" variable and the weight was defined as the percent perimeter SMA positivity.

713    The average (non-zero) thickness of the myoepithelium for each duct was calculated by

714    finding the weighted average "wedge thickness" for SMA-positive wedges ("SMA-

715    present" was 1). The wedge thickness was calculated as the distance between the inner-

716    most and outer-most positive mesh-segments. The positive wedges were weighted by

717    their area relative to the total area of positive wedges.

718    The percent myoepithelial-covered perimeter and average myoepithelial thickness

719    metrics were waited over meshes (ducts) in a given image by assigning a weight to each

720    duct equal to the total area of the duct myoepithelium divided by the sum of the total areas

721    of all myoepithelium in the image that met a minimum size filter of 7500 pixels.

722

723    **Myoepithelial Pixel Clustering Analysis**

724    Pre-processed (background subtracted, de-noised) images were first subset for pixels

725    within the myoepithelium mask. Pixels within the myoepithelium mask were then further

726    subset for pixels with SMA expression greater than 0. For all SMA$^+$ pixels within the

727    myoepithelium mask, a Gaussian blur was applied using a standard deviation of 1.5 for

728    the Gaussian kernel. Pixels were normalized by their total expression, such that the total

729    expression of each pixel was equal to 1. A 99.9% normalization was applied for each

730    marker. Pixels were clustered into 100 clusters using FlowSOM (Van Gassen et al., 2015)

731    based on the expression of 6 markers: PanCK, CK5, Vimentin, ECAD, CD44, and CK7.

732    The average expression of each of the 100 pixel clusters was found and the z-score for

733    each marker across the 100 pixel clusters was computed. All z-scores were capped at 3,

734    such that the maximum z-score was 3. Using these z-scored expression values, the 100

735    pixel clusters were hierarchically clustered using Euclidean distance into 6 metaclusters.

736    SMA$^+$ pixels that were negative for the 6 markers used for FlowSOM were annotated as

737    the SMA-only metacluster, resulting in a total of 7 metaclusters. These metaclusters were

738    mapped back to the original images to generate overlay images colored by pixel

739    metacluster.

740

741    **Collagen Morphometrics**

742    To identify collagen fibers the background-removed Col1 images are first preprocessed:

743    Col1 pixel intensities were capped at 5 and gamma transformed (1 of 2), and contrast

744    enhanced. Images are then blurred via gaussian with sigma of 2. While this enhances

745    fidelity, it gives less clear '0-borders'. This is mitigated by generating a '0-region' mask

746    and setting all values to 0 in that region. Then, highly localized contrast enhancement is

747    applied. Raw fiber signal intensity can vary greatly within a FOV, so this step helps to

748    enhance locally recognizable, but globally dim fiber candidates. After this process,

749    contrast is globally enhanced via a reverse gamma transformation (2 of 2).

750         Collagen fiber objects are generated by watershed segmentation on the

751    preprocessed images. An adaptive thresholding method was developed to appreciate

752    variability in total image intensities across the large dataset. A dilated and eroded version

753    of each preprocessed image was produced and subjected to multiotsu thresholding. For

754    thin fibers, the higher watershed region is set to everywhere where the eroded image has

755    greater intensity than the highest multiotsu threshold for the eroded image, while the lower

756    watershed region is set to everywhere where the dilated image has lower intensity than

757    the highest multiotsu threshold for the eroded image. For thick fibers, the same procedure

758    is performed, except the lower watershed region uses the middle multiotsu threshold for

759    the dilated image. Elevation maps for watershed are generated via the sobel gradient of

760    a blurred version of the preprocessed images. Once objects are extracted and

761    segmented, length, global orientation, perimeter, and width are computed for each object.

762    Objects which cover low intensity regions of the image are treated as preprocessing

763    artifacts and are not included in averaging.

764         For fiber alignment scoring, fibers are filtered for elongated shape (length >

765    2*width), and alignment is scored as the normalized total paired square difference over

766    its k nearest neighbors (k = 4 was chosen). To accommodate for the elongated shape of

767    these object, K-nearest neighbors were computed with the 'ellipsoidal membrane

768  distance' (EM distance), which is the Euclidean centroid distance minus the portion of

769  said distance that lies within the ellipse representation of the object.

770

771  **Cibersort Analysis**

772  CIBERSORTx (CSx)(Newman et al., 2019) was used to infer the immune fraction in LCM-

773  SMART3SEQ samples. We first generated a tissue resident immune cell signature matrix

774  by using a published breast cancer scRNAseq dataset, downloaded from Gene

775  Expression Omnibus database (GEO data repository accession numbers GSE114727,

776  GSE114725)(Barrett et al., 2013). Normalized counts were obtained by using Seurat R

777  package (version 3.2.0). The resultant signature matrix contained 3484 genes and

778  allowed to resolve different immune cell types, including B, CD8 T, CD4 T, NKT, NK, mast

779  cells, neutrophils, monocytes, macrophages and dendritic cells. The signature matrix was

780  first in-silico validated. In order to test the accuracy of the signature matrix, a set of

781  samples from the same scRNAseq dataset was reserved to build a synthetic matrix of

782  bulk RNAseq data. By mixing different proportion of single cells transcripts, the synthetic

783  bulk was used to analyze the correlation between known vs obtained cell proportions by

784  CSx. Pearson's coefficient was above 0.75 in all of the cases, most of them above 0.9.

785  Therefore, we used the aforementioned matrix to deconvolve the LCM-RNAseq samples

786  and to compare CSx-estimated cell abundance with MIBI-identified cell types.

787

788  **Prediction of recurrence**

789  To predict recurrence, we identified patients in the cohort with follow-up data

790  demonstrating carcinoma recurrence (n=12), invasive recurrence (n=19), or at least 11

791  years without recurrence (n=47). For each patient, a vector of summary statistics was

792  generated from MIBI data using only images derived from the original lesion. The cohort

793  was split into training and test sets (80/20%); all model optimization and predictor

794  selection used only the training set. Any missing values were replaced with the set's

795  predictor mean. Predictors with <12 unique values in the training set were dropped from

796  the analysis. Two-class random forest probability models (ranger package)(Wright and

797  Ziegler, 2017) were trained to discriminate recurrence versus non-recurrence, and

798  invasive recurrence versus non-recurrence. Hyperparameters were tuned to minimize

29

799   out-of-bag error. One tuned hyperparameter was predictor subset selection by

800   correlation thresholding: predictors were ranked in importance by performing a KS test

801   between recurrence and non-recurrence. Greater importance was placed on predictors

802   with lower p-values, with ties broken by weighting predictors with greater coefficients of

803   variance (CV). All predictors were correlated (Spearman method) and correlations were

804   thresholded (*invasive* r>|0.5|, *all recurrence* r>|0.6|). For each group of correlated

805   predictors above a given threshold, only the highest-ranked predictor was used in the

806   model. The optimized random forest model was evaluated on the test set and a receiver

807   operating characteristic (ROC) curve was generated (pROC package)(Robin et al.,

808   2011) using the model's assigned probability scores. Area under the curve (AUC) was

809   calculated with 95% confidence intervals, determined by bootstrapping. Each predictor's

810   importance was evaluated in the model by its Gini index. Similarly, two-class random

811   forest probability models were also trained using only clinical parameters as predictors

812   (age, mammograph density, tumor grade, and tumor necrosis) without subset selection.

813   For the MIBI-based predictions, an optimal probability threshold was selected by the

814   Youden method to assign predicted class to the test set, and Kaplan-Meier curves were

815   calculated (survival package)(Therneau and Grambsch, 2000).

816

817   **Statistical Analysis**

818   All statistical analyses were performed using GraphPad Prism software or in R. Grouped

819   data is presented with individual sample points throughout, and where not applicable,

820   data is presented as a mean with standard deviation. For determining significance,

821   grouped data was first tested for normality with the D'Agostino & Pearson omnibus

822   normality test. Normally distributed data was compared between two groups with the two-

823   tailed Student's T-test. Non-normal data was compared between two groups using the

824   Mann–Whitney Test. Multiple groups were compared using the Dunn's Multiple

825   Comparison Test.

826

827   **Software**

828   Image processing was conducted with Matlab 2016a and Matlab 2019b. Statistical

829   analysis was conducted in Graphpad Prism. Data visualization and plots were generated

830 in R with ggplot and pheatmap packages, in Graphpad Prism, and in Python using the
831 scikitimage, matplotlib, and seaborn packages. Representative images were processed
832 in Adobe Photoshop. Schematic visualizations were produced with Biorender. R
833 packages for GSEA: AnnotationDbi, 1.52.0 & org.Hs.eg.db, 3.12.0, clusterProfiler,
834 version 3.19.0, for GSEA msigdbr, version '7.2.1', for C2 curated datasets. Python
835 packages for spatial enrichment analysis and collagen morphometrics: sckikit-image,
836 pandas, numpy, xarray, scipy, statsmodels.

837

## Data and Code Availability

839 All custom code used to analyze data will be made available through a Github repository
840 and all processed images and annotated single cell data will be made available on a
841 Human Tumor Atlas Network public repository.

842

843

## Author Contributions

845 TR conceived the study design, performed experiments, analyzed data, and wrote the
846 manuscript with MA. DG developed the classifier model and performed related analyses.
847 CCL developed the myoepithelial pixel clustering approach and performed related
848 analyses. SHS. processed the LCM-RNAseq data and BRG performed all RNAseq
849 analyses. EFM assisted with data analysis with AK, LK, and SV. N.F.G. assisted with
850 image segmentation. AB developed and performed the myoepithelial morphology
851 analyses and AK performed the collagen morphology analyses. GAC, DJV, KD assisted
852 with cohort design and patient sample preparation, and SS performed pathological
853 review, and SV and ZK assisted with immunohistochemistry. SEH, SCB, RBW and MA
854 supervised the work.

855

## Acknowledgements

867

## Conflicts of Interest

869   M.A. and S.C.B. are inventors on patent US20150287578A1. M.A. and S.C.B. are board

870   members and shareholders in IonPath Inc. T.R. and E.F.M. have previously consulted for

871   IonPath Inc.

872

873

## References

875   Afghahi, A., Forgó, E., Mitani, A.A., Desai, M., Varma, S., Seto, T., Rigdon, J., Jensen, K.C., Troxell,
876   M.L., Gomez, S.L., et al. (2015). Chromosomal copy number alterations for associations of
877   ductal carcinoma in situ with invasive breast cancer. Breast Cancer Res. *17*, 108.

878   Ak, C., A, S., R, G., E, S., A, L., W, P., T, C., F, M.-B., Me, E., and Ne, N. (2018). Multiclonal
879   Invasion in Breast Tumors Identified by Topographic Single Cell Sequencing (Cell).

880   Alaeikhanehshir, S., Engelhardt, E.G., van Duijnhoven, F.H., van Seijen, M., Bhairosing, P.A.,
881   Pinto, D., Collyar, D., Sawyer, E., Hwang, S.E., Thompson, A.M., et al. (2020). The impact of
882   patient characteristics and lifestyle factors on the risk of an ipsilateral event after a primary
883   DCIS: A systematic review. Breast Edinb. Scotl. *50*, 95–103.

884   Alcazar, C.R.G.D., Huh, S.J., Ekram, M.B., Trinh, A., Liu, L.L., Beca, F., Zi, X., Kwak, M., Bergholtz,
885   H., Su, Y., et al. (2017). Immune Escape in Breast Cancer During In Situ to Invasive Carcinoma
886   Transition. Cancer Discov. *7*, 1098–1115.

887   Anders, S., and Huber, W. (2010). Differential expression analysis for sequence count data.
888   Genome Biol. *11*, R106.

889   Barsky, S.H., and Karlin, N.J. (2005). Myoepithelial Cells: Autocrine and Paracrine Suppressors of
890   Breast Cancer Progression. J. Mammary Gland Biol. Neoplasia *10*, 249–260.

891   Barth, P.J., Moll, R., and Ramaswamy, A. (2005). Stromal remodeling and SPARC (secreted
892   protein acid rich in cysteine) expression in invasive ductal carcinomas of the breast. Virchows
893   Arch. *446*, 532–536.

894     Betsill, W.L., Rosen, P.P., Lieberman, P.H., and Robbins, G.F. (1978). Intraductal carcinoma.
895     Long-term follow-up after treatment by biopsy alone. JAMA *239*, 1863–1867.

896     Buerger, H., Otterbach, F., Simon, R., Poremba, C., Diallo, R., Decker, T., Riethdorf, L.,
897     Brinkschmidt, C., Dockhorn-Dworniczak, B., and Boecker, W. (1999). Comparative genomic
898     hybridization of ductal carcinoma in situ of the breast-evidence of multiple genetic pathways. J.
899     Pathol. *187*, 396–402.

900     Conklin, M.W., Eickhoff, J.C., Riching, K.M., Pehlke, C.A., Eliceiri, K.W., Provenzano, P.P., Friedl,
901     A., and Keely, P.J. (2011). Aligned Collagen Is a Prognostic Signature for Survival in Human
902     Breast Carcinoma. Am. J. Pathol. *178*, 1221–1232.

903     Curtis, C., Shah, S.P., Chin, S.-F., Turashvili, G., Rueda, O.M., Dunning, M.J., Speed, D., Lynch,
904     A.G., Samarajiwa, S., Yuan, Y., et al. (2012). The genomic and transcriptomic architecture of
905     2,000 breast tumours reveals novel subgroups. Nature *486*, 346–352.

906     Erbas, B., Provenzano, E., Armes, J., and Gertig, D. (2006). The natural history of ductal
907     carcinoma <Emphasis Type="BoldItalic">in situ</Emphasis> of the breast: a review. Breast
908     Cancer Res. Treat. *97*, 135–144.

909     Esbona, K., Yi, Y., Saha, S., Yu, M., Doorn, R.R.V., Conklin, M.W., Graham, D.S., Wisinski, K.B.,
910     Ponik, S.M., Eliceiri, K.W., et al. (2018). The Presence of Cyclooxygenase 2, Tumor-Associated
911     Macrophages, and Collagen Alignment as Prognostic Markers for Invasive Breast Carcinoma
912     Patients. Am. J. Pathol. *188*, 559–573.

913     Eusebi, V., Feudale, E., Foschini, M.P., Micheli, A., Conti, A., Riva, C., Di Palma, S., and Rilke, F.
914     (1994). Long-term follow-up of in situ carcinoma of the breast. Semin. Diagn. Pathol. *11*, 223–
915     235.

916     Foley, J.W., Zhu, C., Jolivet, P., Zhu, S.X., Lu, P., Meaney, M.J., and West, R.B. (2019). Gene
917     expression profiling of single cells from archival tissue with laser-capture microdissection and
918     Smart-3SEQ. Genome Res. *29*, 1816–1825.

919     Fujii, H., Szumel, R., Marsh, C., Zhou, W., and Gabrielson, E. (1996). Genetic progression,
920     histological grade, and allelic loss in ductal carcinoma in situ of the breast. Cancer Res. *56*,
921     5260–5265.

922     Goswami, S., Sahai, E., Wyckoff, J.B., Cammer, M., Cox, D., Pixley, F.J., Stanley, E.R., Segall, J.E.,
923     and Condeelis, J.S. (2005). Macrophages promote the invasion of breast carcinoma cells via a
924     colony-stimulating factor-1/epidermal growth factor paracrine loop. Cancer Res. *65*, 5278–
925     5283.

926     Hollern, D.P., Swiatnicki, M.R., and Andrechek, E.R. (2018). Histological subtypes of mouse
927     mammary tumors reveal conserved relationships to human cancers. PLoS Genet. *14*, e1007135.

928   Ibrahim, A.M., Moss, M.A., Gray, Z., Rojo, M.D., Burke, C.M., Schwertfeger, K.L., dos Santos,
929   C.O., and Machado, H.L. (2020). Diverse Macrophage Populations Contribute to the
930   Inflammatory Microenvironment in Premalignant Lesions During Localized Invasion. Front.
931   Oncol. *10*.

932   Jones, J.L., Shaw, J.A., Pringle, J.H., and Walker, R.A. (2003). Primary breast myoepithelial cells
933   exert an invasion-suppressor effect on breast cancer cells via paracrine down-regulation of
934   MMP expression in fibroblasts and tumour cells. J. Pathol. *201*, 562–572.

935   Keren, L., Bosse, M., Marquez, D., Angoshtari, R., Jain, S., Varma, S., Yang, S.-R., Kurian, A., Van
936   Valen, D., West, R., et al. (2018). A Structured Tumor-Immune Microenvironment in Triple
937   Negative Breast Cancer Revealed by Multiplexed Ion Beam Imaging. Cell *174*, 1373-1387.e19.

938   Kim, M., Chung, Y.R., Kim, H.J., Woo, J.W., Ahn, S., and Park, S.Y. (2020). Immune
939   microenvironment in ductal carcinoma in situ: a comparison with invasive carcinoma of the
940   breast. Breast Cancer Res. BCR *22*.

941   Kim, S.Y., Jung, S.-H., Kim, M.S., Baek, I.-P., Lee, S.H., Kim, T.-M., Chung, Y.-J., and Lee, S.H.
942   (2015). Genomic differences between pure ductal carcinoma in situ and synchronous ductal
943   carcinoma in situ with invasive breast cancer. Oncotarget *6*, 7597–7607.

944   Kolijn, K., Verhoef, E.I., Smid, M., Böttcher, R., Jenster, G.W., Debets, R., and van Leenders,
945   G.J.L.H. (2018). Epithelial-Mesenchymal Transition in Human Prostate Cancer Demonstrates
946   Enhanced Immune Evasion Marked by IDO1 Expression. Cancer Res. *78*, 4671–4679.

947   Lien, H.C., Hsiao, Y.H., Lin, Y.S., Yao, Y.T., Juan, H.F., Kuo, W.H., Hung, M.-C., Chang, K.J., and
948   Hsieh, F.J. (2007). Molecular signatures of metaplastic carcinoma of the breast by large-scale
949   transcriptional profiling: identification of genes potentially related to epithelial-mesenchymal
950   transition. Oncogene *26*, 7859–7871.

951   Lin, C.-Y., Tsai, P.-H., Kandaswami, C.C., Lee, P.-P., Huang, C.-J., Hwang, J.-J., and Lee, M.-T.
952   (2011). Matrix metalloproteinase-9 cooperates with transcription factor Snail to induce
953   epithelial-mesenchymal transition. Cancer Sci. *102*, 815–827.

954   Linde, N., Casanova-Acebes, M., Sosa, M.S., Mortha, A., Rahman, A., Farias, E., Harper, K.,
955   Tardio, E., Reyes Torres, I., Jones, J., et al. (2018). Macrophages orchestrate breast cancer early
956   dissemination and metastasis. Nat. Commun. *9*, 21.

957   Liu, Y., West, R., Weber, J.D., and Colditz, G.A. (2019). Race and risk of subsequent aggressive
958   breast cancer following ductal carcinoma in situ. Cancer *125*, 3225–3233.

959   Malanchi, I., Santamaria-Martínez, A., Susanto, E., Peng, H., Lehr, H.-A., Delaloye, J.-F., and
960   Huelsken, J. (2012). Interactions between cancer stem cells and their niche govern metastatic
961   colonization. Nature *481*, 85–89.

962  McCaffrey, E.F., Donato, M., Keren, L., Chen, Z., Fitzpatrick, M., Jojic, V., Delmastro, A.,
963  Greenwald, N.F., Baranski, A., Graf, W., et al. (2020). Multiplexed imaging of human
964  tuberculosis granulomas uncovers immunoregulatory features conserved across tissue and
965  blood. BioRxiv 2020.06.08.140426.

966  Moen, E., Bannon, D., Kudo, T., Graf, W., Covert, M., and Van Valen, D. (2019). Deep learning
967  for cellular image analysis. Nat. Methods *16*, 1233–1246.

968  Morsing, M., Kim, J., Villadsen, R., Goldhammer, N., Jafari, A., Kassem, M., Petersen, O.W., and
969  Rønnov-Jessen, L. (2020). Fibroblasts direct differentiation of human breast epithelial
970  progenitors. Breast Cancer Res. *22*, 102.

971  Newburger, D.E., Kashef-Haghighi, D., Weng, Z., Salari, R., Sweeney, R.T., Brunner, A.L., Zhu,
972  S.X., Guo, X., Varma, S., Troxell, M.L., et al. (2013). Genome evolution during progression to
973  breast cancer. Genome Res. *23*, 1097–1108.

974  Newman, A.M., Steen, C.B., Liu, C.L., Gentles, A.J., Chaudhuri, A.A., Scherer, F., Khodadoust,
975  M.S., Esfahani, M.S., Luca, B.A., Steiner, D., et al. (2019). Determining cell type abundance and
976  expression from bulk tissues with digital cytometry. Nat. Biotechnol. *37*, 773–782.

977  Page, D.L., Dupont, W.D., Rogers, L.W., and Landenberger, M. (1982). Intraductal carcinoma of
978  the breast: follow-up after biopsy only. Cancer *49*, 751–758.

979  Pelon, F., Bourachot, B., Kieffer, Y., Magagna, I., Mermet-Meillon, F., Bonnet, I., Costa, A., Givel,
980  A.-M., Attieh, Y., Barbazan, J., et al. (2020). Cancer-associated fibroblast heterogeneity in
981  axillary lymph nodes drives metastases in breast cancer through complementary mechanisms.
982  Nat. Commun. *11*, 404.

983  Peng, J., Wang, X., Ran, L., Song, J., Luo, R., and Wang, Y. (2018). Hypoxia-Inducible Factor 1α
984  Regulates the Transforming Growth Factor β1/SMAD Family Member 3 Pathway to Promote
985  Breast Cancer Progression. J. Breast Cancer *21*, 259–266.

986  Poola, I., DeWitty, R.L., Marshalleck, J.J., Bhatnagar, R., Abraham, J., and Leffall, L.D. (2005).
987  Identification of MMP-1 as a putative breast cancer predictive marker by global gene
988  expression analysis. Nat. Med. *11*, 481–483.

989  Robin, X., Turck, N., Hainard, A., Tiberti, N., Lisacek, F., Sanchez, J.-C., and Müller, M. (2011).
990  pROC: an open-source package for R and S+ to analyze and compare ROC curves. BMC
991  Bioinformatics *12*, 77.

992  Ruffell, B., Affara, N.I., and Coussens, L.M. (2012). Differential Macrophage Programming in the
993  Tumor Microenvironment. Trends Immunol. *33*, 119–126.

994  Ryser, M.D., Weaver, D.L., Zhao, F., Worni, M., Grimm, L.J., Gulati, R., Etzioni, R., Hyslop, T., Lee,
995  S.J., and Hwang, E.S. (2019). Cancer Outcomes in DCIS Patients Without Locoregional
996  Treatment. JNCI J. Natl. Cancer Inst. *111*, 952–960.

997 Shani, O., Vorobyov, T., Monteran, L., Lavie, D., Cohen, N., Raz, Y., Tsarfaty, G., Avivi, C.,
998 Barshack, I., and Erez, N. (2020). Fibroblast-derived IL-33 facilitates breast cancer metastasis by
999 modifying the immune microenvironment and driving type-2 immunity. Cancer Res.

1000 Subramanian, A., Tamayo, P., Mootha, V.K., Mukherjee, S., Ebert, B.L., Gillette, M.A., Paulovich,
1001 A., Pomeroy, S.L., Golub, T.R., Lander, E.S., et al. (2005). Gene set enrichment analysis: a
1002 knowledge-based approach for interpreting genome-wide expression profiles. Proc. Natl. Acad.
1003 Sci. U. S. A. *102*, 15545–15550.

1004 Therneau, T.M., and Grambsch, P.M. (2000). Modeling Survival Data: Extending the Cox Model
1005 (New York: Springer-Verlag).

1006 Valen, D.A.V., Kudo, T., Lane, K.M., Macklin, D.N., Quach, N.T., DeFelice, M.M., Maayan, I.,
1007 Tanouchi, Y., Ashley, E.A., and Covert, M.W. (2016). Deep Learning Automates the Quantitative
1008 Analysis of Individual Cells in Live-Cell Imaging Experiments. PLOS Comput. Biol. *12*, e1005177.

1009 Van Gassen, S., Callebaut, B., Van Helden, M.J., Lambrecht, B.N., Demeester, P., Dhaene, T., and
1010 Saeys, Y. (2015). FlowSOM: Using self-organizing maps for visualization and interpretation of
1011 cytometry data. Cytom. Part J. Int. Soc. Anal. Cytol. *87*, 636–645.

1012 Wright, M.N., and Ziegler, A. (2017). ranger: A Fast Implementation of Random Forests for High
1013 Dimensional Data in C++ and R. J. Stat. Softw. *77*, 1–17.

1014 Yang, M., Li, Z., Ren, M., Li, S., Zhang, L., Zhang, X., and Liu, F. (2018). Stromal Infiltration of
1015 Tumor-Associated Macrophages Conferring Poor Prognosis of Patients with Basal-Like Breast
1016 Carcinoma. J. Cancer *9*, 2308–2316.

1017 Zhang, W., Shi, X., Peng, Y., Wu, M., Zhang, P., Xie, R., Wu, Y., Yan, Q., Liu, S., and Wang, J.
1018 (2015). HIF-1α Promotes Epithelial-Mesenchymal Transition and Metastasis through Direct
1019 Regulation of ZEB1 in Colorectal Cancer. PloS One *10*, e0129603.

1020 Zhang, W., Zhang, J., Zhang, Z., Guo, Y., Wu, Y., Wang, R., Wang, L., Mao, S., and Yao, X. (2019).
1021 Overexpression of Indoleamine 2,3-Dioxygenase 1 Promotes Epithelial-Mesenchymal Transition
1022 by Activation of the IL-6/STAT3/PD-L1 Pathway in Bladder Cancer. Transl. Oncol. *12*, 485–492.

1023 Zhou, J., Wang, X.-H., Zhao, Y.-X., Chen, C., Xu, X.-Y., sun, Q., Wu, H.-Y., Chen, M., Sang, J.-F., Su,
1024 L., et al. (2018). Cancer-Associated Fibroblasts Correlate with Tumor-Associated Macrophages
1025 Infiltration and Lymphatic Metastasis in Triple Negative Breast Cancer Patients. J. Cancer *9*,
1026 4635–4641.

1027