

UCSF

UC San Francisco Previously Published Works

Title

The breast pre-cancer atlas illustrates the molecular and micro-environmental diversity of ductal carcinoma in situ

Permalink

<https://escholarship.org/uc/item/4r10231d>

Journal

npj Breast Cancer, 8(1)

ISSN

2374-4677

Authors

Nachmanson, Daniela

Officer, Adam

Mori, Hidetoshi

et al.

Publication Date

2022

DOI

10.1038/s41523-021-00365-y

Peer reviewed

ARTICLE OPEN



The breast pre-cancer atlas illustrates the molecular and micro-environmental diversity of ductal carcinoma in situ

Daniela Nachmanson¹, Adam Officer^{1,2}, Hidetoshi Mori³, Jonathan Gordon^{4,5}, Mark F. Evans^{4,6}, Joseph Steward⁷, Huazhen Yao⁸, Thomas O'Keefe⁹, Farnaz Haste^{7,10}, Gary S. Stein^{4,5}, Kristen Jepsen⁸, Donald L. Weaver^{4,6}, Gillian L. Hirst¹¹, Brian L. Sprague^{4,12}, Laura J. Esserman¹¹, Alexander D. Borowsky¹³, Janet L. Stein¹³ and Olivier Harismendy¹³✉

Microenvironmental and molecular factors mediating the progression of Breast Ductal Carcinoma In Situ (DCIS) are not well understood, impeding the development of prevention strategies and the safe testing of treatment de-escalation. We addressed methodological barriers and characterized the mutational, transcriptional, histological, and microenvironmental landscape across 85 multiple microdissected regions from 39 cases. Most somatic alterations, including whole-genome duplications, were clonal, but genetic divergence increased with physical distance. Phenotypic and subtype heterogeneity was frequently associated with underlying genetic heterogeneity and regions with low-risk features preceded those with high-risk features according to the inferred phylogeny. B- and T-lymphocytes spatial analysis identified three immune states, including an epithelial excluded state located preferentially at DCIS regions, and characterized by histological and molecular features of immune escape, independently from molecular subtypes. Such breast pre-cancer atlas with uniquely integrated observations will help scope future expansion studies and build finer models of outcomes and progression risk.

npj Breast Cancer (2022)8:6; <https://doi.org/10.1038/s41523-021-00365-y>

INTRODUCTION

Increasing adoption of breast cancer screening and advances in imaging capabilities have improved our ability to identify breast ductal carcinoma in situ (DCIS). Rarely diagnosed 40 years ago, DCIS now comprises nearly 20% of all breast cancer-related diagnoses^{1,2}. Unfortunately, this progress has not resulted in decreased breast cancer mortality. Standard treatment, involving surgical excision often complemented with radiation therapy (in the setting of breast-conserving surgery) and endocrine recurrence risk reduction (particularly with ER + DCIS), therefore constitutes overtreatment, and not without treatment-related consequences for many^{2,3}. DCIS progression is particularly difficult to study longitudinally due to the current standard of surgical excision of the lesion and the infrequent progression and/or occurrence of new primary lesions over a long timespan (5–10% after 10 years)⁴. Clinicopathological risk factors such as large size, dense breast, younger age, high pathological grade, presence of comedo necrosis, or Her2 positivity have been associated with increased risk of recurrence, but the resulting predictive models, or those relying on gene expression signatures, are currently insufficient to safely distinguish patients to watch from patients to treat⁵.

Contrary to models of progression in other tissue types, there is little evidence for the sequential accumulation of somatic alterations during progression from in situ to invasive breast cancer (IBC), but rather all IBC intrinsic subtypes and known driver

mutations have been identified in DCIS, albeit at variable prevalence^{6–12}. Moreover, both single-cell and bulk studies have shown similar clonal make-up of synchronous invasive and in situ lesions, convoluting the idea that clonal selection drives invasion^{11,13}. The role of the immune environment has also been investigated, highlighting the higher lymphocyte infiltration in Her2+ or Triple Negative DCIS, or specific immunological make-up of samples at higher risk of progression^{12,14–19}. Similarly, the role of the basal layer, fibroblasts, adipocytes, other stromal cells, or overall extracellular matrix has identified features that are different between DCIS and IBC, likely mediated by chemokine signaling, and can be associated with known progression risk factors^{20–23}. Their active participation in the malignant transformation of the breast epithelium remains to be established as similar mechanisms are typically involved in normal development, activity, and aging of the mammary gland^{24,12}.

Progress in our understanding of the processes mediating DCIS onset and progression has been considerably hindered by technical and logistical limitations. Indeed, pure DCIS lesions are commonly small in size, formalin and paraffin embedded (which damages nucleic acids), and can display significant histological heterogeneity²⁵. As a consequence, comprehensive molecular and cellular assays and their integrated analysis have seldom been performed in pure DCIS cohorts. Capturing evidence of phenotypic, genetic, and cellular heterogeneity, and how they relate to each other is necessary to develop a better spatial, temporal, and functional understanding of the mechanisms at play. Recent

¹Bioinformatics and Systems Biology Graduate Program, University of California San Diego, 9500 Gilman Drive, San Diego, CA 92093, USA. ²Division of Biomedical Informatics, Department of Medicine, University of California San Diego, 9500 Gilman Drive, San Diego, CA 92093, USA. ³Department of Pathology and Laboratory Medicine, Center for Immunology and Infectious Diseases, School of Medicine, University of California Davis, 2315 Stockton Blvd, Sacramento, CA 95817, USA. ⁴University of Vermont Cancer Center, 111 Colchester Avenue Main Campus, Main Pavillion, Level, 2, Burlington, VT 05401, USA. ⁵Department of Biochemistry, University of Vermont, Burlington, VT 05405, USA. ⁶Department of Pathology and Laboratory Medicine, University of Vermont, Burlington, VT 05405, USA. ⁷Moores Cancer Center, University of California San Diego, 3855 Health Science Drive, San Diego, CA 92093, USA. ⁸Institute for Genomic Medicine, University of California San Diego, 9500 Gilman Drive, San Diego, CA 92093, USA. ⁹Department of Surgery, University of California San Diego, 9500 Gilman Drive, San Diego, CA 92093, USA. ¹⁰Department of Pathology, University of California San Diego, 9500 Gilman Drive, San Diego, CA 92093, USA. ¹¹Helen Diller Family Comprehensive Cancer Center, University of California San Francisco, 1450 3rd St, San Francisco, CA 94158, USA. ¹²Department of Surgery, University of Vermont, Burlington, VT 05405, USA. ✉email: oharismendy@ucsd.edu

advances in genome-wide assays, becoming compatible with ever more challenging samples^{26–29}, have improved our ability to connect histological and molecular observations and enabled such application even to individual microbiopsies from a histological slide of pure DCIS.

Here we describe the combined, parallel histological, molecular, and immunological profiling of premalignant lesions from 39 patients diagnosed with DCIS, including multiple epithelial microbiopsies within a subset of samples. The dissection of specific epithelial lesions provided a detailed assessment of the association of their histological architecture with intrinsic subtypes, mutational landscape, driver mutations, and immunological states. Multi-region profiling resulted in the inference of clonal relationships, illustrating how genotypes related to phenotypes within a specimen. We, therefore, report multi-modal and sub-histological profiling of a cohort of pure DCIS, illustrating spatial heterogeneity and placing diverse states of immune activity observed in their specific molecular and histological context.

RESULTS

Histological and molecular characterization

We collected a total of 43 specimens (referred to as samples) from 39 patients diagnosed with pure DCIS, including three samples from subsequent DCIS diagnosed between 14 and 70 months after the index DCIS (Fig. 1a, b, Table 1, Supplementary Table 1). Sixty-nine percent (29/42) of the samples were positive for estrogen receptor (ER) expression and 40% (16/40) had *ERBB2*

gene overexpression or amplification (Supplementary Fig. 1a). Each sample was further annotated for grade and histological architecture and the annotations were used to identify regions of interest, guide the microbiopsies of the epithelial areas and the immuno-histological analysis. On the basis of their studied regions, the cohort consisted of 32 high or intermediate grade DCIS (HG-DCIS), nine low-grade DCIS (LG-DCIS), and two low-grade atypical ductal hyperplasia (ADH). The DCIS regions could be further annotated according to their dominant histological architecture (17 cribriform, 19 solid, three mixed, two micropapillary) and the presence of necrosis (ten comedo necrosis, 17 other). LG-DCIS were more frequently of cribriform architecture (8/9), while HG-DCIS were frequently necrotic (25/32). The relative area of adipose tissue in each sample varied between four and 91% as estimated by segmental classification of the whole slide digital image (Fig. 1c, Methods). The lower adipose fraction was associated with higher mammographic breast density ($p = 0.0067$) suggesting the sample histology was representative of the whole breast texture. Interestingly, solid DCIS were associated with a higher adipose fraction (median 69% vs 40%, $p = 0.008$), suggesting a contribution of the breast microenvironment to the growth architecture. Overall, the cohort represents a diverse set of pure in situ lesions identified in absence of any detectable invasive component. The studied samples are enriched for DCIS lesions and specifically annotated for their histological architecture. Each sample was profiled using multiple assays, performed on sequential histological sections (4–7 μm) used for whole transcriptome, whole exome, and spatial immune profiling. Whenever

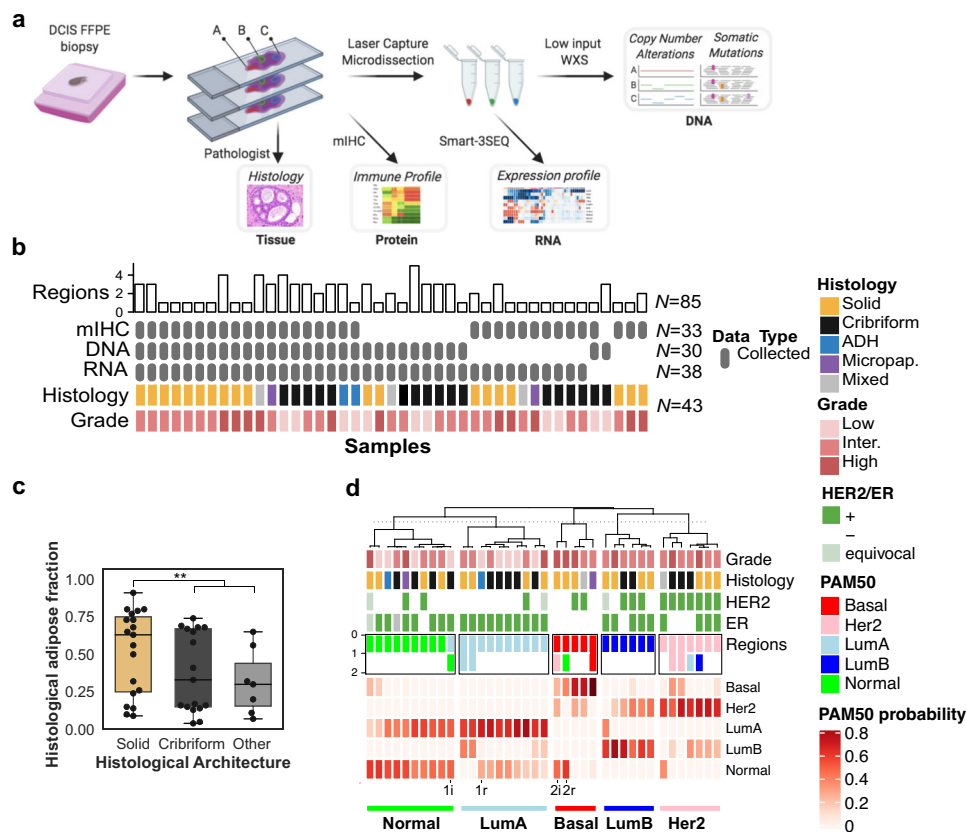


Fig. 1 Study design and cohort overview. **a** Archival sample processing and analysis workflow including histology (H&E and miHC) and microbiopsy-derived whole exome (WXS) and whole transcriptome (Smart-3SEQ) profiling. **b** Study cohort overview including histological characteristics (colored rows), data type (gray rows), and number of histological regions (bar chart) investigated. **c** Estimate of the fractional characteristics in epithelium of the miHC images according to each histological architecture, ** $p < 0.01$, Mann-Whitney U -test. **d** Sample classification according to the probabilities of each PAM50 expression subtype. For ten eligible samples, the intrinsic subtype of a spatially distinct region is indicated. Two patients with recurrence (r) and index (i) samples are indicated at the bottom. Error bars in the box and whiskers plots represent 1.5 fold the interquartile range above (resp. below) the first (resp. third) quartile of the distribution.

Table 1. Clinical and pathological features of the patient and specimen studied.

Patient ID	Block ID	Age at Index	Size (cm)	Laterality	Grade	Architecture	ER	HER2 ^a	N. of Regions	Diagnosis Order ^b
MCL76_044	12800	56	0.9	Left	Low	Cribriform	+	-	1	Index
MCL76_049	18100	50	6	Left	Low	Cribriform	+	-	1	Index
MCL76_060	16100	47	17	Left	Low	Cribriform	+	-	3	Index
MCL76_061	16200	34	8	Left	Low	Cribriform	+	-	4	Index
MCL76_064	15200	78	0.4	Left	Low	Cribriform	+	-	3	Recur. (+18 mos.)
MCL76_066	14400	70	1.1	Right	Low	Cribriform	+	-	3	Index
	16500	70	0.3	Right	Low	ADH	+	-	1	Recur. (+14 mos.)
MCL76_076	15700	45	4.1	Right	Low	Cribriform	+	-	5	Index
MCL76_078	15500	68	1.4	Left	Low	Solid	+	-	3	Index
MCL76_080	15800	59	3.7	Left	Low	ADH	+	-	3	Index
MCL78_020	10001	59	0.3	Right	Low	Cribriform	+	+	1	Index
MCL76_012	11600	50	3.6	Right	Inter.	Solid	+	-	2	Index
MCL76_048	13100	51	3.8	Right	Inter.	Cribriform	+	-	3	Index
MCL76_064	14600	78	NA	Left	Inter.	Solid	+	NA	1	Index
MCL76_067	16600	54	6	Left	Inter.	Solid	+	+	3	Index
MCL76_070	16400	69	8	Right	Inter.	Solid	+	-	3	Index
MCL76_071	14800	68	5.8	Right	Inter.	Micropapillary	-	-	3	Index
MCL76_074	14700	45	14	Right	Inter.	Cribriform	+	-	3	Index
MCL76_077	15300	70	1.2	Left	Inter.	Cribriform	-	+	2	Index
MCL76_079	15400	62	3.4	Right	Inter.	Cribriform	-	+	3	Index
MCL78_001	10001	50	2.5	Right	Inter.	Cribriform	NA	-	1	Index
MCL78_002	10001	48	2	Left	Inter.	Solid	+	-	1	Index
MCL78_006	10001	75	4	Left	Inter.	Cribriform	+	+	1	Index
MCL78_007	10001	43	1.6	Right	Inter.	Cribriform	+	+	1	Index
MCL78_008	10001	66	1.5	Right	Inter.	Solid	+	+	1	Index
MCL78_009	10001	78	2.4	Right	Inter.	Solid	+	-	1	Index
MCL78_010	10001	67	1.1	Left	Inter.	Cribriform	-	+	1	Index
MCL78_011	10001	59	0.6	Right	Inter.	Solid	+	+	1	Index
MCL78_013	10001	63	2.2	Right	Inter.	Mixed	-	+	1	Index
MCL78_016	10001	65	4.5	Left	Inter.	Solid	-	+	1	Index
MCL78_017	10001	52	2.5	Left	Inter.	Solid	+	+	1	Index
MCL78_018	10001	65	2	Right	Inter.	Mixed	+	equ	2	Index
MCL76_007	11000	78	3.5	Left	High	Solid	-	-	3	Index
	11100	78	2.6	Right	High	Solid	-	-	4	Recur. (+39 mos.)
MCL76_016	11800	35	5	Left	High	Mixed	+	+	4	Index
MCL76_025	16800	75	1.2	Right	High	Solid	-	NA	2	Index
MCL76_068	14900	59	9.5	Left	High	Cribriform	-	+	3	Index
MCL78_003	10001	43	5	Left	High	Solid	+	equ	1	Index
MCL78_005	10001	81	0.5	Right	High	Solid	+	-	1	Index
MCL78_012	10001	54	4	Right	High	Solid	-	+	1	Index
	10014	54	3	Right	High	Solid	-	NA	1	Synchronous
MCL78_015	10001	57	0.5	Left	High	Micropapillary	+	+	1	Index
MCL78_019	10001	57	1.9	Left	High	Solid	-	equ	1	Index

^aInferred from ERBB2 copy number and expression (Fig. S1), equ equivocal.

^bRecur. Recurrence, mos. months. All recurrence DCIS were in different quadrants than the index.

possible, the investigated regions were matched across assays to preserve the spatial information in the analysis and limit the variation due to spatial heterogeneity. Spatial heterogeneity was further addressed in 21 samples for which multiple sub-regions were profiled independently.

The expression of genes was measured using high-throughput sequencing of RNA-seq libraries directly prepared from the microbopsied regions²⁷ (Supplementary Table 2). The samples

were classified according to the PAM50 intrinsic subtypes used for invasive breast cancer (IBC), which identified Basal ($N = 5$), Luminal A ($N = 10$), Luminal B ($N = 6$), Her2-like ($N = 7$), and Normal-like ($N = 10$) samples (Fig. 1d). Consistent with IBC classification, Luminal A and B were enriched for samples from ER + cases, while Her2-like were enriched for Her2+ cases. Similarly, Luminal B and Her2-like were enriched in HG-DCIS, while Luminal A was almost exclusively composed of cribriform LG-DCIS. Luminal A and

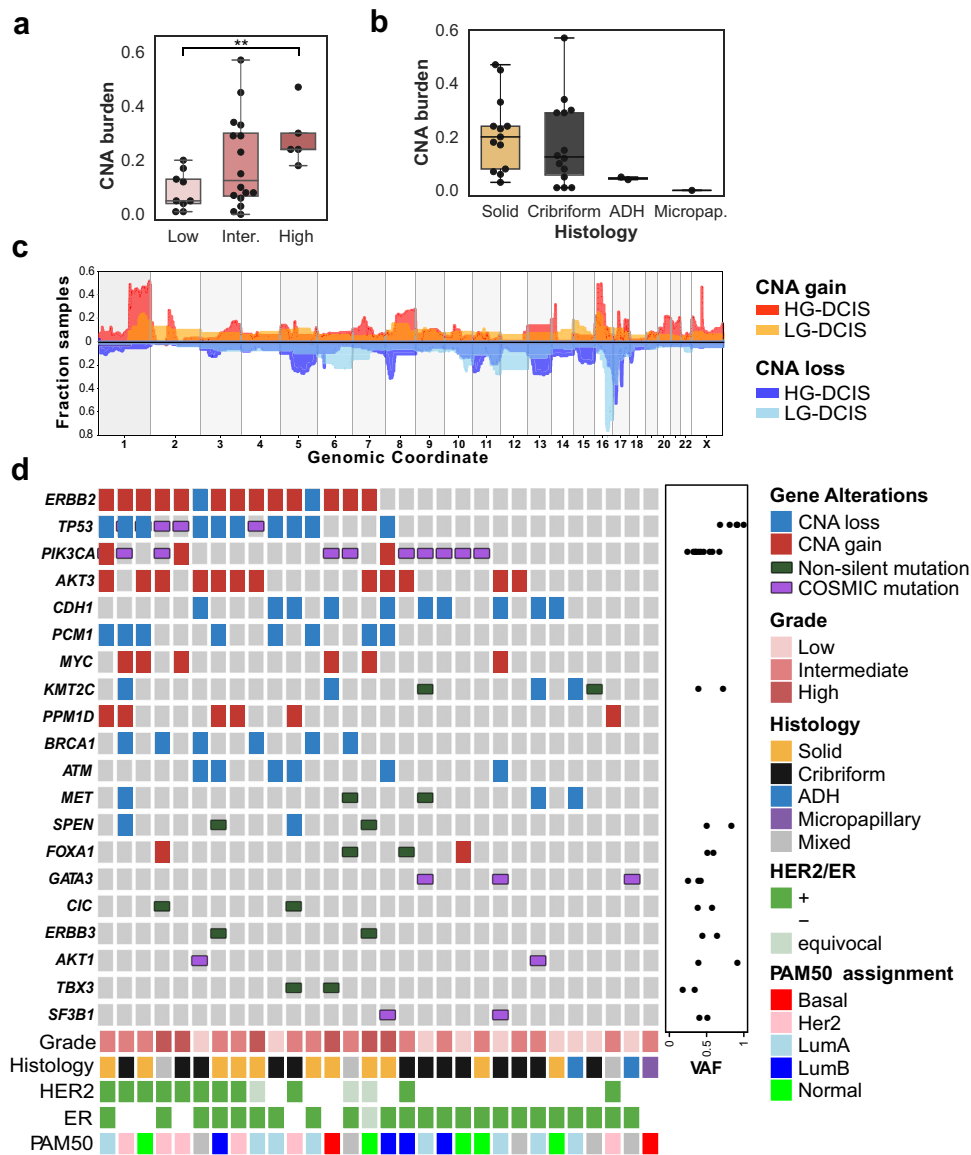


Fig. 2 Pure DCIS genomic landscape. **a**, **b** CNA burden (fraction of base pairs involved in copy number gain or loss) as a function of grade (**a**) and **b** histological architecture, $**p < 0.05$, ANOVA. **c** Smoothed frequency (y-axis) of CNA gains (top) and losses (bottom) smoothed along the genome (x-axis) for HG-DCIS ($N = 22$ —dark colors) and LG-DCIS ($N = 5$ light colors). **d** Oncoprint diagram displaying the mutational status of driver genes commonly altered in breast cancer. Genes were included if they were mutated in at least two patients or located in a CNA segment present in at least six patients, and ordered by frequency of alteration. The variant allele fraction (VAF) of mutations (right panel) and histological characteristics (bottom panel) are indicated. Error bars in the box and whiskers plots represent 1.5 fold the interquartile range above (resp. below) the first (resp third) quartile of the distribution.

Normal-like represented closely related classes and together comprised the majority of the samples (20/38), which is not unexpected given the higher fraction of low-grade and pure in situ lesions in the cohort, in contrast with IBC and previous DCIS expression profiling studies^{18,30}. The PAM50 subtype of two independent sub-regions with matching histology and grade was determined in ten samples and observed to be discordant in five samples (Supplementary Table 2B), which was associated with larger distances between the regions (Mann–Whitney, $p = 0.005$, Supplementary Fig. 1b). Interestingly, matched index and recurrent samples from two patients had at least one region with a concordant subtype. Across all samples, the distribution of probabilities for each PAM50 subtype likely captures such heterogeneity. Normal-like were truly a mix of Normal and Luminal A, while Her2-like tended to have two main subsets: Her2/Basal and Her2/Luminal B. This suggests that subtypes

inferred from bulk analysis, even after epithelial microdissection, are frequently the result of a variable mixture of pure subtypes.

Subtype differences in the mutational landscape

To determine whether any of the histological or molecular subtypes described above were associated with specific genetic alterations, we characterized their mutational landscape. Whole exome sequencing was carried out on microbiopsies from 30 samples using a procedure specifically optimized for a low amount of damaged DNA²⁶. Mutations and copy number alterations (CNA) were identified in 27 and 30 samples, respectively (Supplementary Table 3). The median copy number burden—or fraction of the genome involved in CNA—was 0.14 and was 2.5 fold higher in HG-DCIS (Mann–Whitney, $p = 0.017$, Fig. 2a, b). Whole-genome doubling (WGD) events were detected in 3/8 eligible samples, all of which were low or intermediate grade

Table 2. Frequency of *PIK3CA*, *TP53*, and *GATA3* driver mutations in previously reported DCIS studies and pure DCIS in this study.

Gene	Pang et al. 2017 (N = 20)	Lin et al. 2019 (N = 65)	Nagasawa et al. 2021 (N = 72)	Pareja et al. 2020 (N = 7)	This study					
					All ^a	Grade		Histology		
						Low	Inter.-high	Cribriform	Solid	Other
<i>PIK3CA</i>	55%	40%	50%	0%	43% (10/23)	29% (2/7)	50% (8/16)	55% (6/11)	40% (4/10)	0% (0/2)
<i>TP53</i>	30%	13.8%	21%	14.3%	31.3% (5/16)	0% (0/4)	41.7% (5/12)	33.3% (3/9)	40% (2/5)	0% (0/2)
<i>GATA3</i>	45%	13.8%	56%	28.6%	20% (3/15)	75% (3/4)	0% (0/13)	33.3% (2/6)	0% (0/9)	50% (1/2)

^aThe denominator represents samples with at least 20x coverage across the targeted regions.

cribriform DCIS consistent with its early timing in breast carcinogenesis³¹ (Supplementary Table 3). Consistent with previous studies, loss of 16q (13/30) and 17p (12/30) or gain of 1q (12/30) were among the most frequent chromosomal alterations and, while many events were more frequent in HG-DCIS (Fig. 2c, Supplementary Table 4), these hallmarks were also observed in low-grade or benign lesions, including ADH: (1q gain: 1/9, 16q loss: 7/9, 17p loss: 3/9).

We identified between 74 and 207 coding mutations per sample. The mutational burden was higher in HG-DCIS (Mann–Whitney, $p = 0.003$) and Her2-like subtypes (Mann–Whitney, $p = 0.025$), recognizing that these categories are overlapping. The HG-DCIS burden (4.4 mut/Mb) was higher than previous reports, possibly due to residual germline variants in our study^{11,13}. We identified aging-associated mutational signatures (SBS1 and SBS5) in all samples eligible for analysis ($N = 13$), APOBEC signature (SBS2 and SBS13) on one intermediate grade solid DCIS and mismatch repair signature (SBS15 or SBS21) in three DCIS of variable grade and architecture (Supplementary Table 5). The APOBEC signature is, therefore, rarer in DCIS than IBC (~8% vs >75%), but can be present in premalignant lesions. Interestingly, this sample also displayed clustered mutations ($N = 3$ within 1,416 bp) in chromosome 17q (Supplementary Fig. 2), an APOBEC-driven kataegis site frequently seen in IBC³². The most recurrently mutated genes were *PIK3CA* (44%), *TP53* (31%), and *GATA3* (20%), and were all affected by known somatic mutations in breast cancer at similar rates to previous studies of pure DCIS^{6,7,9–11} (Fig. 2d and Table 2). *TP53* mutations were only found in HG-DCIS and associated with high CNA burden (Mann–Whitney, $p = 0.018$), while *GATA3* mutations were only found in cribriform or ADH histologies and associated with LG-DCIS (Fisher Exact, $p = 0.005$). Interestingly, *GATA3* mutations were identified in larger lesions (Mann–Whitney, $p = 0.038$), consistent with a similar observation in invasive cancer and the larger tumor size of *GATA3* mutated xenograft models^{33,34}. Another nine selected genes known to be mutated in IBC were recurrently affected by 18 mutations predicted to be deleterious, four of which are known somatic mutations⁷. The result suggests that oncogenic driver mutations are already present at the premalignant stage, including in LG-DCIS (e.g. *SF3B1* c.2098 A > G) or ADH (*GATA3* c.925-3_925-2del). This is consistent with previous reports and reports of field-effect mutations in normal ducts or benign lesions^{35,36}, though the contribution of these mutations to the lesion progression remains to be determined.

Genetic heterogeneity and clonal diversity

The histologic assessment and expression profiling have revealed variable levels of phenotypic heterogeneity across the samples. In order to determine whether such heterogeneity is present at the genetic level, we measured genetic heterogeneity in two distinct ways: (1) divergence, which measures the genetic distance between regions of a sample and, (2) clonal relationships, which uses phylogenetic tree construction to establish evolutionary

order to genetic alterations (Supplementary Table 3). We measured divergence by computing a CNA-based score on 19 pairs of histologically matching regions in 11 samples (Supplementary Table 3, Methods). With no pairs completely independent, the spatial distance separating dissected DCIS regions was correlated with the extent of their genetic divergence ($R^2 = 0.65$, $p = 0.00017$, Supplementary Fig. 3), while this could simply be a result of local proliferation, it could also be a consequence of selective pressures of the microenvironment, migratory capacity or genomic instability of particular clones. Interestingly, one ADH had the lowest divergence despite a large distance, suggesting either a different pattern in ADH or a distance threshold for the extent of the correlation. Divergence was not associated with grade, Her2/ER status, or adipose fraction suggesting that local genetic heterogeneity is not associated with progression risk factors.

More precise clonal relationships between regions were evaluated using phylogenetic analysis in 12 samples, comparing CNA, and mutations when available (Fig. 3, Supplementary Fig. 4, Methods). While the majority (88.4%) of CNA were shared across all regions of a sample, 11.6% were private to some regions, as observed in 7/12 samples. Multiple samples (3/12) contained mutations in putative cancer driver genes that were private to one region only. These included known and likely pathogenic mutations in *ATR*, *PIK3CA*, *MET*, *KDM5C*, suggesting that not all driver mutations are acquired early. Interestingly, the three samples with the most private CNA displayed discordant histological architecture or discordant PAM50 subtypes between regions, suggesting that within a sample, genetic and phenotypic differences are linked. Furthermore, in 4/5 samples containing regions with discordant histology and 3/4 with discordant PAM50 subtypes, features historically associated with low-risk of progression (benign histology, Normal or LumA subtype), appeared earlier than regions with high-risk features (Her2 or Basal subtype, presence of necrosis). Overall, these results illustrate that in these samples, regions evolved to acquire distinct histological and molecular features, and in particular, regions with low-risk features can precede regions with high-risk features.

Substantial heterogeneity and evolutionary patterns are evident in samples like MCL76_061_16200 (Fig. 3a–c), where a region of benign columnar alterations preceded two cribriform regions. While all regions shared a WGD event as well as several arm-level CNA and pathogenic mutations in *GATA3* and *SF3B1*, the cribriform region A acquired private 5q and 8q gains and necrotic features. While this example shows tandem genetic and histological changes as seen across the cohort, it also illustrates that despite occurring earlier, the benign region shares many “driver-like” alterations with both cribriform regions. Furthermore, in another example, despite homogeneous cribriform histologies in regions of MCL76_077_15300 (Fig. 3d–f) only one cribriform region lost a copy number of chromosome 8 and presented with Her2 PAM50 subtype as opposed to its Luminal A predecessor. Notably, bulk studies have shown chromosome 8 loss to be more frequent in Her2 vs Luminal A breast cancers³⁷. Taken all together we

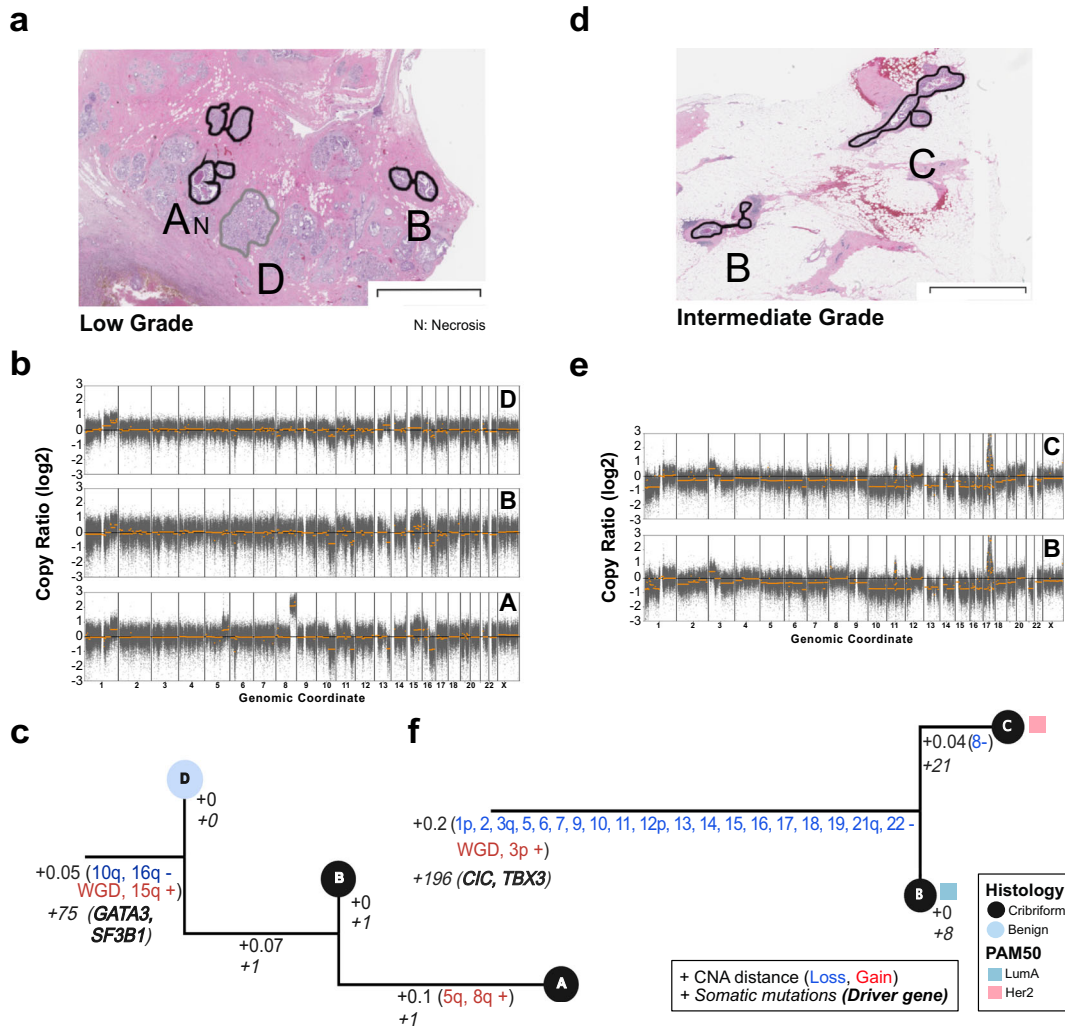


Fig. 3 Clonal relationships of multi-region DCIS. Multi-region phylogenetic reconstruction using both CNA and somatic mutations for MCL76_061_16200 (**a–c**) and MCL76_077_15300 (**d–f**). For each case, the spatial annotation of the microdissected regions on the H&E images (**a, d**), corresponding copy number profiles (**b, e**), and phylogenetic trees (**c, f**) are displayed. Copy number profile plots show bins (gray dots) and segment (orange) log₂ copy number ratio (y-axis). The phylogenetic tree leaves (single dissected region) are colored according to histological type and the branches (hamming distances based on CNA segments) are annotated with corresponding specific somatic alterations or their total number (CNA: regular, genes: italic font). The tree root corresponds to an inferred normal diploid ancestor. PAM50 subtype of the region is indicated when available. Annotations and trees are available for ten additional samples in Supplementary Fig. 4. The scale bars in panels **a** and **b** correspond to a size of 3 mm.

illustrate abundant genetic heterogeneity in pure DCIS of all histologies and grades that parallels the levels of phenotypic heterogeneity and often accompanies it, even in regions that are millimeters apart.

Regional differences in the immune microenvironment

To measure the diversity of the immune landscape and to investigate its potential association with molecular or histological features, we used multiplex immunohistochemistry (mIHC) to measure the number and density of four cell types—T-cells (CD3+), B-cells (CD20+), T-regs (CD3+/FOXP3+) and epithelial cells (PanCK+)—according to their proliferative status (Ki67+). Both epithelial (PanCK+) and adjacent stromal (PanCK- proximal to epithelium) areas from premalignant ($N=36$ regions across 32 samples) or normal ($N=21$ across 21 samples) histologies were evaluated. Among premalignant regions, the high-grade epithelial areas had lower cell density due to larger cell sizes and frequent central necrosis (median 3.8 vs 6.4 10^3 cells/mm² $p < 0.03$ —Mann–Whitney). Solid lesions had the highest fraction of

proliferating epithelial cells (median 11.5% vs 2.8% $p < 0.02$ —Mann–Whitney, Supplementary Fig. 5a), and interestingly 3/10 HG-DCIS cribriform lesions (two Her2-like, one Luminal B) had markedly higher proliferation. Consistent with previous findings, we observed higher lymphocyte infiltration in ER- and Her2+ samples compared to ER+ ones (Supplementary Fig. 5b, Mann–Whitney, $p < 0.001$). We next classified all regions using non-negative matrix factorization of the stromal and epithelial cell densities, resulting in three immune states characterized by their dominant meta-markers (MM; Fig. 4a, b Supplementary Table 7a–c): (1) “Active”—ubiquitous high T-cells (high MM2), including a subset with elevated T-cell proliferation (high MM1), (2) “Suppressed”—ubiquitous low T-cells (low MM1 and MM2), high B-cells and T-regs (high MM3), and (3) “Excluded”—high stromal, low epithelial densities (high MM4). To further confirm differences between immune states, we compared the total T-cell, B-cell and T-regs densities in epithelial and stromal compartments. While overall lymphocyte densities were much higher in stroma than in epithelium across all examined regions (median ratio 9.8, Supplementary Table 7a), the skew was a distinguishing feature

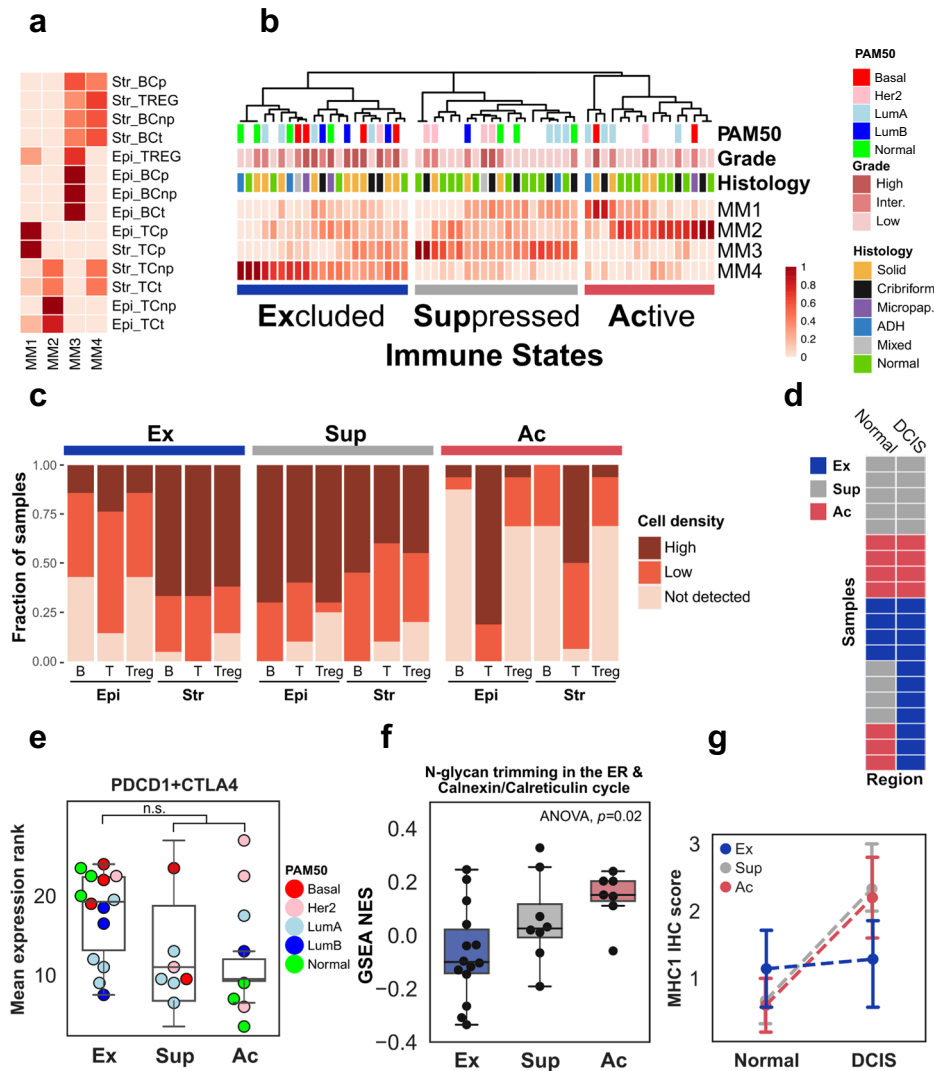


Fig. 4 Characterization of the immune landscape. Decomposition of immune cell density scores by non-negative matrix factorization (NMF) into **a** W-matrix which shows the composition of the Meta-Markers (columns MM1-4) according to the densities scores (red scale) of each cell type (BC: B-cells, TC: T-cells, TREG: regulatory T-cells), proliferative state (p : Ki67+, np : Ki67-, t : total) and regional location (Epi: Epithelium, Str: Stroma) and **b** H-matrix which classifies normal and DCIS regions into three immune states according to Meta-Markers. **c** Fraction of stromal and epithelial regions from samples in each immune state with high, low or no T-cells (T), B-cells (B) and regulatory T-cells (Treg) densities. **d** Immune-state comparison in 20 samples (rows) with matching normal (left column) and DCIS (right column) regions. **e** Expression of immune checkpoint receptors genes, *PDCD1* and *CTLA4* in each immune state. **f** GSEA normalized enrichment score (NES) for a Reactome gene set across immune states. **g** Distribution of the expression of the MHC-I complex scored by immunohistochemical staining in DCIS and normal adjacent regions. The median scores of the adjacent and DCIS region in each immune state are connected with a dotted line. Error bars in the box and whiskers plot represent 1.5 fold the interquartile range above (resp. below) the first (resp. third) quartile of the distribution.

in regions in Excluded state for all three cell types (Fig. 4c, Supplementary Fig. 6 and Supplementary Table 7b). Furthermore, regions in Active states had the highest epithelial T-cell density (120 cells/mm²) while regions in Suppressed state had the highest T-regs and B-cells epithelial densities (10.4 and 8.1 cells/mm² respectively). A larger fraction of the normal regions were found in Active (7/21) or Suppressed state (9/21) rather than in Excluded state (4/21) and premalignant regions in Excluded state were more likely to be high grade (7/15 vs 2/17 $p = 0.049$). Interestingly, the immune states of normal and premalignant regions were concordant in 12/19 matched cases and discordant in seven whose lesions were specifically in the Excluded state (Fig. 4d). This suggests that the eExcluded state may be acquired in response to premalignant growth, while other states may be intrinsic to various breast microenvironments. Furthermore, premalignant regions in Suppressed state were more likely identified in cases

younger than 55 (5/8 vs 4/24 OR = 7.6 $p = 0.02$), consistent with the younger age of DCIS patients with infiltrating PD-L1 + lymphocytes³⁸. We did not observe any associations between immune states and intrinsic subtype, ER or Her2 status, tumor size, breast density, adipose fraction, or DCIS architecture suggesting that they may be independent from traditional histopathological progression risk factors.

In order to identify functional differences between immune states, we evaluated the differential activity of Hallmark and Reactome processes among the 29 DCIS regions with available gene expression information (Supplementary Fig. 7). Compared to Active and Suppressed states, the Excluded state was associated with upregulation of Type 1 and 2 Interferon response, PD1 signaling, and proliferation-related processes as well as the repression of Calreticulin-Calnexin cycle (Supplementary Fig. 7). Noting that the epithelium of DCIS in Excluded state were not

completely depleted of infiltrating lymphocytes, the upregulated processes were consistent with the higher expression of *PCDC1* or *CTLA4* genes in DCIS in Excluded state (Fig. 4e), albeit not significant, and suggesting a likely continuum of increasing immuno-suppression from Suppressed to Excluded states. More interestingly, the repression of the Calreticulin-Calnexin cycle was confirmed via single-sample enrichment analysis and showed a progressive repression from Active, to Suppressed, to Excluded states ($p = 0.022$, ANOVA, Fig. 4f). This suggests that the export of glycoproteins—including components of MHC1 complex—via the endoplasmic reticulum, impacts immune-surveillance. To verify this hypothesis, we measured the in situ expression of MHC1 complex in 15 samples (Supplementary Table 8 and Supplementary Fig. 8) and compared its levels in adjacent normal and DCIS in each immune state. While the level of MHC1 expression in DCIS region were not significantly different between Excluded and non-Excluded samples, the change between normal and DCIS was different, with the non-Excluded samples displaying increased expression between normal and DCIS, while the Excluded samples remained constant ($p = 0.0009$, Mann–Whitney, $N = 15$, Fig. 4g). This therefore suggests that the Excluded immune state may be mediated by both intrinsic expression level of MHC1 and ability to increase it in DCIS.

DISCUSSION

There is a compelling requirement for a DCIS atlas that delivers a relatively unbiased, multi-modal perspective of pre-invasive breast cancer. Here, we report the multi-modal profiling of a diverse set of pure DCIS. This comprehensive atlas both confirms previous molecular findings and provides a higher resolution histological and spatial context to interpret them. However, with only three known recurrences, the significance of our observations for progression prognosis could not be formally established. Our findings provide a landscape of representative pure DCIS identified in absence of invasive lesions. While some lesions were small, others were quite extended ($N = 14 > 4$ cm), which should capture factors that may be associated with robust containment. The cohort therefore spans a variety of clinical, histological, phenotypic, and genotypic features. Such variety and contrast are critical to ensure this atlas' utility in designing larger studies, or perhaps providing more cautionary interpretation of observations from cohorts enriched for specific risk factors.

At the heart of our study's innovation was the ability to generate molecular profiles from limited amounts of dissected archival tissue specimen. Similar approaches are used to study clonal expansion in normal tissues^{28,29}, but generally not performed in parallel for RNA and DNA. Importantly some limitations remain and not all assays were successful. The large variability in success rate was not easy to predict. Likely the age of the specimen, its size, fixation conditions and storage conditions all contribute to success variability which cannot be controlled in a retrospective investigation. Additional limitations are analytical, such as the absence of a matched source of normal DNA from every sample which can result in residual germline variants, perhaps inflating the overall mutation rate observed. The use of adjacent normal tissue can also be problematic and there is ample evidence that they also accumulate somatic mutations³⁹. In our study, we clearly identified known breast cancer driver mutations in samples from ADH or other benign alterations. Overall, while some samples are unlikely to ever contain sufficient material for profiling or dissection of adjacent normal, as methodologies evolve and advance, the success rate and data quality will improve to make molecular pre-malignant profiling more accessible and as routine as is the case in invasive cancer.

Our report contributes to two major advances for understanding pre-malignant lesions. First, we characterized most samples across four important modalities all within a maximum

of 50 μ m sequentially sectioned tissue. Such advances were enabled by pre-analytical improvements allowing us to reduce the tissue requirement, to include small lesions, and to precisely match regions of interest across each modality: histology, epithelial gene expression, DNA mutations, and immune landscape. As a result, we could isolate regions with different histological features that may coexist within a specimen and more confidently establish their association with expression subtypes, clonal heterogeneity, or immune state. For example, the integration of histology and expression subtypes showed a clear correlation between cribriform architecture and the Luminal A subtype. By integrating histology, expression subtype, and immune state we showed that some immune states are found in normal areas and that there is no clear association between immune state and expression subtype. Hence, the depth and interpretability of the analysis are considerably increased by integrating all modalities at the regional level. This has been clearly the case in large cancer studies such as the TCGA, or, more recently through the integrated analysis of histological and somatic features in normal, aging tissues^{28,29,40}. While most studies do not typically include immunohistochemical or other multiplexed spatial analysis, other important advancements in this field in the past year include spatial proteomics used to evaluate the structure of the myoepithelium in DCIS, and spatial transcriptomics used to identify the transcriptional effect of driver mutations in DCIS, representing the emerging frontier of pre-malignant tissue characterization^{10,41}. It is therefore likely that additional spatial profiling compatible with FFPE specimens will bring additional prognostic and mechanistic insights in future DCIS studies.

The other important contribution of our study is the sub-histological analysis to compare regions of interest from the same sample and infer phylogenetic relationships between them. While we determined that the majority of the DCIS samples were classified as Normal-like and Luminal A subtype, typically considered less-aggressive subtypes in breast cancer and reflective of the known precursor stage that DCIS represents, we showed evidence for intrinsic heterogeneity in the PAM50 probabilities, either from the distribution of probabilities within a region or from physically separated regions. This is not entirely surprising as bulk expression subtypes are the result of averaging heterogeneity, similar to glioblastoma subtypes⁴² or IBC subtypes⁴³ from single-cell analysis. Such heterogeneity, especially in DCIS, had been proposed before on the basis of marker staining⁴⁴ and our results confirm that it may be rather common. Similarly to the frequency of heterogeneity between region subtypes, we identified evidence of genetic heterogeneity in 7/12 cases, including the presence of private putative driver mutations. This fraction may be an underestimate given the close proximity of many selected pairs. However, the majority of putative genetic drivers, copy number hallmarks and even WGD were clonal, shared by all regions investigated, including a few normal regions. This observation supports evolutionary models derived from invasive cancer, including multi-sample studies, that suggest that most driver mutations occur early followed by a phase of clonal expansion. Similar observations were also made in early multi-regional studies in DCIS^{44–46} and studies comparing synchronous DCIS-IBC cases using single-cell sequencing¹³, providing further evidence that breast cancer genetic evolution starts in the pre-invasive stage and possibly in normal regions. It is likely that driver alterations may even be present in adjacent histologically normal tissue as observed in field effects studies in normal ducts^{39,47}. Such effects support an important contribution of host factors to the initial genetic injury. Hence, unlike previous attempts which were focused on histopathological features, including grade, surgical margins^{48,49}, future DCIS prognostic models will likely need to be derived from lifetime cancer risk models like GAIL⁵⁰ or BOADICEA⁵¹ and incorporate host-specific factors, such as

polygenic risk scores and reproductive factors, that likely contribute to the DCIS initiation and trajectory.

The immune microenvironment of DCIS has been previously investigated, using both quantification of tumor infiltrating lymphocytes (TIL) and more specific immunohistochemical approaches and revealed clear quantitative and qualitative variation in lymphocyte infiltration, including higher TIL number and more immunosuppressive features in high-risk lesions¹⁹. Importantly, previous studies in pure DCIS did not quantify stromal and epithelial TILs separately^{12,19}. This distinction may be hard to make in IBC, where both compartments interact at the invasive front and pathologist subjectivity can have a major impact^{52,53}. However, this separation can be more clearly established in the analysis of DCIS and was critical in the identification of the Excluded immune state in our atlas. While the Active and Suppressed states have been observed before and could readily be identified in our data, the identification of the Excluded state required the use of an analytical method (NMF) to account for the strong correlations that can exist between TILs type and compartments. The inclusion of adjacent normal areas was also important to interpret the significance of the immune states, as the Excluded state appeared more likely in reaction to the DCIS growth and increased grade. The Excluded state exhibited features of immune evasion and could represent a more advanced level of immuno-suppression than the Suppressed state, with the consequence of a topological exclusion from the duct. The downregulation of components of the Calreticulin-Calnexin cycle in the epithelium in Excluded state could impact MHC-I export or maturation, as suggested by the lack of MHC-I expression induction in DCIS of the eExcluded state, hence providing an evasion mechanism, and contrasting with evasion mediated by MHC-I genetic loss observed in IBC^{54,55}. It would be interesting to determine whether the immune states identified can explain the variability of response to local injection of anti-PD1 antibody in DCIS patients and whether any of the states would elicit, or prevent, the desired ductal infiltration by T-cells⁵⁶.

As illustrated by our study and recent advances in the profiling of normal tissues^{28,29}, histopathology, and molecular pathology are becoming more integrated fields, generating deeper and broader datasets at increased cellular and spatial resolution, from the most challenging human samples. Future studies of early transformation and pre-cancer biology such as the one presented here will likely benefit the most from such approaches which capture heterogeneity at scale and can help reconcile analog (optical) and digital (genomics and multiplex) observations. As a result, such multi-dimensional integration may help identify common factors mediating epithelial transformation and progression across multiple glands and organs.

METHODS

Sample collection and preparation

FFPE blocks were obtained from UCSD or UVM Pathology Departments after surgical biopsy, excision or mastectomy. The study was reviewed and approved by each institutional review board and they granted a waiver of consent. Eligibility criteria were: (1) adult female, (2) pure DCIS diagnosis (without evidence of invasive disease), and (3) with available pathology blocks. Few cases also had bilateral disease or were matching index and recurrent lesion (ipsilateral or contralateral—Table 1 and Supplementary Table 1). Importantly there was no attempt to enrich high-risk cases or investigate specifically the role of certain candidate risk factors. Factors such as age, grade, race, ER, or Her2 status were not part of the selection criteria and the cohort was designed to reflect patients seen in a regular DCIS clinic. All specimen blocks were de-identified and sectioned sequentially for the following purpose: Hematoxylin-Eosin (H&E) staining ($N = 1$; 4 μM glass slide), Laser Capture Microdissection (LCM; $N = 3$; 7 μM glass slide coated with polyethylene naphthalate—ThermoFisher #LCM0522), multiplex or regular immunohistochemistry ($N \geq 3$ 4 μM glass slide) and a final H&E staining ($N = 1$; 4 μM glass slide). The H&E slides were

scanned at high resolution and reviewed and annotated by the study pathologist. The LCM slides were stored at -20°C in an airtight container with desiccant until ready for dissection (1 day to 3 months). H&E sections were diagnosed according to the standard of care criteria (AJCC TNM 8th ed./CAP Breast DCIS Reporting Protocol v4.3). DCIS features recorded included lesion grade: Grade I (low), Grade II (intermediate), or Grade III (high), and associated histology: e.g., papillary, cribriform, solid, comedo necrosis. DCIS lesion, normal glands (and in some cases hyperplasia) were delineated on H&E images to assist LCM. DCIS laterality and size, patient age, and menopausal status, and lesion mode of detection were obtained from the original pathology reports or from the Vermont Breast Cancer Surveillance System (UVM specimen) or local cancer registry and chart review (UCSD specimen). Hormone receptor and Her2 statuses (where available) were gathered from the patient reports and/or by de novo IHC staining. The LCM sections were thawed and stained with eosin, sections were kept in xylene and dissected within 2 h of staining. LCM was performed using the ArcturusXM Laser Capture Microdissection System (ThermoFisher). Matching regions from six adjacent sections were collected on CapSure Macro Cap (for DNA, $N = 3$ slides) or HS caps (for RNA, $N = 3$ slides), region size, and unambiguous match permitting. Post-dissection, all caps were covered and stored at -20°C with desiccant.

DNA extraction and QC

The membrane and adhering tissue were peeled off the caps using a razor blade and the peeled membrane was incubated in proteinase K digestion reaction overnight for 16 h at 56°C to maximize DNA yield after cell lysis. The DNA was extracted using the QIAamp DNA Micro Kit (Qiagen) and the elution was done in 20 μL . The extracted DNA was quantified by fluorometry (HS dsDNA kit Qbit—ThermoFisher).

RNA sequencing and analysis

Library Preparation. RNA sequencing was performed using SMART-3Seq, a 3' tagging strategy specifically designed for degraded RNA directly from FFPE LCM specimen²⁷. LCM dissected SMART-3Seq libraries were prepared using the standard protocol for FFPE tissue on Arcturus HS LCM Cap and the individual library SPRI purification option. All FFPE LCM dissected libraries were amplified using 19 PCR cycles during indexing to minimize over-amplification of high abundance mRNAs in each library. Libraries were individually analyzed for size distribution on an Agilent 2200 TapeStation with High Sensitivity D1000 reagent kits to verify average library size of 190 bp and stored at -20°C until sequencing. When all libraries were ready for sequencing, 1 μL of each library was then used to create two library pools used for sequencing and quantified by Qubit 2.0 Fluorometer HS DNA assay. Library pools were sequenced with a 1% PhiX spike-in control library and sequenced on an Illumina HiSeq 4000, a run type of single read 75 (SR75), and dual index sequencing.

Transcriptome analysis. Read count data was obtained using a dedicated analysis workflow <https://github.com/danielanach/SMART-3SEQ-smk>. Briefly, sequencing reads were trimmed using cutadapt 1.18, UMIs were processed using the umi_homopolymer.py script in the SMART-3SEQ tools (<https://github.com/jwfoley/3SEQtools>), aligned using STAR 2.6.1a, deduplicated using the dedup.py script from <https://github.com/jwfoley/umi-dedup> and read counts were calculated using featureCounts 1.6.3^{57,58}. Count data were then merged and filtered to remove samples with <55,000 counts and genes with <10 read counts across all samples. Filtered count data was then loaded into Seurat version 3.2.3 and processed using the *SCTransform()* function version 0.3.2 to regress out the high mitochondrial content variability across the samples⁵⁹. Batch correction was then performed using ComBat to remove variation attributable to the sequencing center (UCSD vs UVM)⁶⁰. PAM50 subtype probabilities were calculated from the *SCTransform* and batch normalized data using the geneFu package⁶¹. Gene set enrichment analysis (GSEA) was performed as in⁶² and single-sample GSEA as in⁶³. Gene sets from the REACTOME and Hallmark collections in MSigDB were used to compare the Excluded to the non-Excluded groups, a permutation test was performed to assess the significance of the GSEA results^{64,65}. ANOVA was used to compare the ssGSEA results between the three miHC groups. FDR of <0.1 and p -values of <0.05 were considered significant.

Whole exome sequencing and primary analysis

Library preparation. DNA was sheared down to 200 base pairs (bp) using Adaptive Focused Acoustics on the Covaris E220 (Covaris Inc) following

manufacturer recommendations with 10 μ L Low EDTA TE buffer supplemented with 5 μ L of truSHEAR buffer using a microTUBE-15. Libraries were prepared using the Accel-NGS 2S PCR-Free DNA Library Kit (Swift Biosciences). Ligated and purified libraries were amplified using KAPA HiFi HotStart Real-time PCR 2X Master Mix (KAPA Biosystems). Samples were amplified with 5 μ L of KAPA P5 and KAPA P7 primers. The reactions were denatured for 45 seconds (sec) at 98 °C and amplified 13–15 cycles for 15 sec at 98 °C, for 30 sec at 65 °C, and for 30 sec at 72 °C, followed by final extension for 1 min at 72 °C. Samples were amplified until they reached Fluorescent Standard 3, cycles being dependent on input DNA quantity and quality. PCR reactions were then purified using 1x AMPure XP bead clean-up and eluted into 20 μ L of nuclease-free water. The resulting libraries were analyzed using the Agilent 4200 TapeStation (D1000 ScreenTape) and quantified by fluorescence (Qubit dsDNA HS assay).

Capture and sequencing. Samples were paired and combined (12 μ L total) to yield a capture “pond” of at least 350 ng, and supplemented with 5 μ L of SureSelect XT HS and XT Low Input Blocker Mix. The hybridization and capture were performed using the Human All Exon V7 panel (S31285117) paired with the Agilent SureSelect XT HS Target Enrichment Kit following the manufacturer’s recommendations. Post-capture amplification was performed on the beads in a 25 μ L reaction: 12.5 μ L of nuclease-free water, 10 μ L 5x Herculase II Reaction Buffer, 1 μ L Herculase II Fusion DNA Polymerase, 0.5 μ L 100 millimolar (mM) dNTP Mix and 1 μ L SureSelect Post-Capture Primer Mix. The reaction was denatured for 30 sec at 98 °C, then amplified for 12 cycles of 98 °C for 30 sec, 60 °C for 30 sec and 72 °C for 1 min, followed by an extension at 72 °C for 5 min and a final hold at 4 °C. Libraries were purified with a 1x AMPure XP bead clean-up and eluted into 20 μ L nuclease-free water in preparation for sequencing. The resulting libraries were analyzed using the Agilent 4200 TapeStation (D1000 ScreenTape) and quantified by fluorescence (Qbit–ThermoFisher). All libraries were sequenced using the HiSeq 4000 sequencer (Illumina) for 100 cycles in Paired-End mode. Libraries with distinct indexes were pooled in equimolar amounts. The sequencing and capture pools were later deconvoluted using program bcl2fastq [19].

Sequencing reads processing and coverage quality control. Sequencing data were analyzed using bcbio-nextgen (v1.1.6) as a workflow manager [20]. Adapter sequences were trimmed using Atropos (v1.1.22), the trimmed reads were subsequently aligned with BWA-MEM (v0.7.17) to reference genome hg19, then PCR duplicates were removed using biobambam2 (v2.0.87)^{66–68}. Additional BAM file manipulation and collection of QC metrics were performed with Picard (v2.20.4) and Samtools (v1.9)⁶⁹. The summary statistics of the sequencing and coverage results are presented in Supplementary Table 9.

Identification of somatic mutation and copy number alterations

Variant calling. Single nucleotide variants (SNVs) and short insertions and deletions (indels) were called with VarDictJava (v1.6.0), and Mutect2 (v2.2)^{70,71}. Variants were required to fall within a 10 bp boundary of targeted regions that overlapped with RefSeq genes (v 109.20190905). A pool of normal DNA was created using whole exome sequencing data of blood of 18 unrelated individuals and was used to eliminate artifacts and common germline variants. Only variants called by both algorithms were considered. These variants were then subjected to an initial filtering step with default bcbio-nextgen tumor-only variant calling filters and the following parameters were used: position covered by at least five reads, mapping quality > 45, mean position in read > 15, number of average read mismatches < 2.5, microsatellite length < 5, tumor log odds threshold > 10, Fisher strand bias Phred-scaled probability < 10 and VAF > 0.1⁷². Functional effects were predicted using SnpEff (v4.3.1)⁷³. All samples were re-evaluated for the presence of COSMIC (v91) database mutations which have been previously observed in at least 15 patients and fall within 137 known breast cancer driver genes (Supplementary Table 10)⁷.

Germline variant filtering. In absence of matched normal tissue for DCIS samples, somatic mutations were prioritized computationally using the approach from the bcbio-nextgen tumor-only configuration then additionally subjected to more stringent filtering⁷². Briefly, common variants (MAF > 10^{−3} or > 9 individuals) present in population databases - 1000 genomes (v2.8), ExAC (v0.3), or gnomAD exome (v2.1)—were removed unless in a tier 1 gene from the cancer gene consensus and present in either COSMIC (v91) or ClinVar (20190513)^{7,74–77}. Variants were removed as

likely germline if found at a variant allelic fraction (VAF) greater or equal to 0.9 in non-LOH genomic segments—as determined by CNA analysis (below). Lastly, variants were also removed as potential germline (or artifact) if found in > 2 patients in the pool of normal (described above).

Single-sample CNA calling. CNVkit⁷⁸ was used for calling somatic copy number alterations (CNA) to measure both overall CNA burden, arm and gene level CNA, and identify LOH as previously described in²⁶. Allele-specific copy number calling algorithm, ASCAT, was used on a select number of samples for which there was sufficient coverage and the algorithm converged on a solution, in order to identify whole-genome doubling events as well as confirm CNA identified by CNVkit⁷⁹. Default parameters were used with ASCAT with the exception of a segmentation penalty of 100 and a gamma of 1.

Multi-region CNA segmentation. To generate harmonized segmentation breakpoints between regions belonging to the same sample, multi-region segmentation was performed with the R CopyNumber (v1.26.0) package⁸⁰. Outliers in CNVkit bin-level log₂ copy ratios were detected and modified using Median Absolute Deviation Winsorization with the winsorize() function, segments were then called using the *mutipcf()* function with a gamma of 40.

Mutational signatures. Mutational signatures were called on merged region samples using a single-sample variation of SigProfiler with default parameters to decompose into known single-base substitutions (SBS) reported in COSMIC^{81,82}.

Analysis of the clonal evolution and genetic heterogeneity

Measurement of genetic divergence. Divergence was measured on each pair of related regions, a and b, using Eq. (1):

$$Divergence_{a,b} = \sum_{k=0}^n |copy\ ratio\ a_k - copy\ ratio\ b_k| * \left(\frac{bins_k}{total\ bins} \right) \quad (1)$$

Where k is the copy number segment, n is the total segments, and bins is the number of bins covered by a segment from the CNVkit input file. The $\frac{bins_k}{total\ bins}$ term was used as a weighted correction factor for the number of bins contributing to a segment. For samples with > 2 regions, the maximum divergence between any two regions was used to represent the sample.

CNA-based phylogenetic reconstruction. Construction of phylogenetic trees was performed similarly to the methodology outlined in⁸³. Briefly, for each sample the log₂ copy ratios from multi-sample copy number segments with at least 12 probes (see above), were translated into a matrix containing −1 for loss (log₂ copy ratio < −0.6), 0 for neutral (−0.4 ≤ log₂ copy ratio ≤ 0.3) and undetermined for anything else. This matrix was then used to generate Maximum Parsimony trees using phangorn using default parameters⁸⁴.

Mutation-based phylogenetic reconstruction. To allow the analysis of clonal relationships between regions of the same sample, the coverage depth of each allele at any remaining mutated position in any region was extracted using Mutect2 joint variant caller on the sets of aligned reads from each region. In order to call a mutation either absent or present in a region, we used a Bayesian inference model specifically designed for multi-region variant calling⁸⁵. Treeomics (v1.7.10) was run with the default parameters except for $e = 0.02$. The tree solution which matched the CNA-based reconstruction was then integrated into a single tree for Fig. 3 and S4.

Multiplex immunohistochemistry

Staining. Tissue sections were prepared from formalin-fixed paraffin embedded tissue blocks and cut to 4 μ m serial sections and mounted on Superfrost Plus (VWR). The procedure for multiplex immunohistochemistry (mIHC) was followed by a manufacturer’s protocol for Opal7-color automation IHC kit (Akoya Bioscience), and the staining was performed with Autostainer DISCOVERY ULTRA (Ventana). Antibodies used in mIHC are anti-CD3 (clone 2GV6, Ventana), anti-CD20 (clone L26, Ventana), anti-Ki67 (clone 30-9, Ventana), anti-FOXP3 (clone SP97, Spring), anti-pan cytokeratin (CK; clone AE1/AE3, DAKO), anti-CD117 (clone c-kit, DAKO). The molecular markers of immune panel (CD3, CD20, Ki67, CKs, FOXP3, and CD117) were visualized with Opal520, Opal540, Opal570, Opal620,

Opal650, and Opal690, respectively. DAPI counterstaining was performed with Discovery QD DAPI (Roche). ProLong Diamond Antifade Mounting (ThermoScientific) was used for mounting the coverslip. Detailed staining conditions and autostainer's protocols are reported in our recent report⁸⁶.

Visualization and analysis. Tissue samples stained with mIHC were scanned with multispectral imaging microscopy (Vectra 3, Akoya Bioscience). Scanned multispectral images were unmixed on inForm software (ver.2.4.0, Akoya Bioscience) to acquire the fluorescence signal from each marker⁸⁶. Imaging analysis was performed on inForm software by identifying tumor (CK+ area) and stroma (CK- area proximal to the epithelium), each nucleated cell, and its cell type. Alternatively, QuPath software 2.3.1⁸⁷ was also used to perform similar imaging analysis on unmixed images converted to multi-layered TIFF format by inForm software⁸⁶. The images of the regions of the same type (DCIS or normal) from the same case, were typically stitched together and stored and shared as one single larger multi-layered TIFF image (data availability below). Scanned image areas were aggregated into up to three histological regions per sample: main pre-invasive lesion, alternate pre-invasive lesion, normal epithelium. In each region, the stromal and epithelial densities of each cell type and state were calculated, including when cells were not present (density = 0). Regularized marker densities into distribution deciles were then used to classify samples using non-negative matrix factorization (Supplementary Table 7c). The immune states were assigned and named after the hierarchical clustering of the H-matrix (meta-marker values).

MHC1 immunostaining and analysis

Four-micron sections were baked at 60 ° for 1 h, followed by deparaffinization through three successive changes of xylene. Tissue was then rehydrated in decreasing grades of alcohol, with two changes of 100%, 95%, and then 70% EtOH, followed by diH2O. Antigen retrieval was performed using Antigen Unmasking Solution Citrate Based pH6, H-3300 (Vector) at 95 °C for 30 min. Staining was performed using the IntelliPATH Automated IHC stainer (Biocare). Endogenous peroxidase was blocked using Bloxall blocking solution, SP-6000 (Vector) for 10 min, followed by two washes in TBST. Afterward, tissue was blocked with a 3% Donkey Serum for 10 min, followed by blocking with Anti-HLA Class I ABC Primary Mouse Antibody, ab70328m (Abcam) at 1:1000 for 1 h and subsequently washed twice with TBST. Tissues were then blocked with Anti-Mouse HRP UltraPolymer IgG, 2MH-100 (Cell IDx) for 30 min, and washed twice with TBST. The reaction was then developed with 3,3'-Diaminobenzidine Chromogen, 95041-478 (VWR) for 5 min, and then stopped with two washes in diH2O. Counterstaining was performed with Mayer's Hematoxylin Solution, 51275 (Sigma) for 5 min. Lastly tissues were washed twice in TBST, and once in diH2O, dehydrated in increasing grades of EtOH, then cleared and mounted with xylene based mountant. MHC1 expression was scored from 0 to 3 separately for DCIS and normal epithelium throughout the entire section, away from possible biopsied areas. The scores were established as follows: 0: no staining or weak staining in <50% of cells; 1: weak staining in >50% of cells; 2: intermediate staining in >50% of cells; 3: strong staining in >50% of cells.

Whole slide image digital analysis

High resolution whole slide images of H&E stains were loaded into a QuPath (v2.3) project⁸⁷. One analysis area was defined for each specimen, avoiding the location of biopsies as well as dust or marked areas. The analysis areas were segmented into superpixels (sigma = 5 µm, spacing = 50 µm, maxiterations = 10, regularization = 0.25) and each superpixel was annotated with both Hematoxylin and Eosin Intensity features (size = 2 µm, tile size = 25 µm). The mean, median, min, max, and standard deviation values were then smoothed (Haralick distance = 1, Haralick bins = 32). Multiple training areas were annotated from each of the following classes: adipose, stroma, inflammation, epithelium (normal and atypical), void, necrosis, blood vessels. Multiple areas across 2–4 samples were used to train a Random Tree classifier. The classifier was then applied to all superpixels included in the analysis area. The accuracy of the classifier was assessed both visually and with multiple test areas for each class. Superpixels of the same class were merged into single annotations and the resulting areas were recorded. Separate classifiers were used for images from different institutions, to mitigate possible variation staining, scanning, or image format. The fraction of adipose area was compared to breast density using Mann-Whitney test comparing dense & heterogeneously

dense breast to other lower densities, or comparing solid DCIS to non-solid DCIS lesions.

Reporting summary

Further information on research design is available in the Nature Research Reporting Summary linked to this article.

DATA AVAILABILITY

The raw RNA and DNA sequencing data have been deposited in dbGAP phs002225. High resolution whole slide images of the H&E stains and corresponding annotations can be viewed on the JPL LabCAS portal (digital object identifiers included in Supplementary Table 1). Images corresponding to the stitched field of views of the region of interest in the multiplex immunohistochemistry are made available as multi-layered tiff files on the JPL LabCAS portal <https://doi.org/10.48577/rmy-pj94> (UVM) and <https://doi.org/10.48577/3gns-rn74> (UCSD).

Received: 26 June 2021; Accepted: 6 December 2021;

Published online: 13 January 2022

REFERENCES

- Siegel, R. L., Miller, K. D. & Jemal, A. Cancer statistics, 2020. *CA Cancer J. Clin.* **70**, 7–30 (2020).
- Bleyer, A. & Welch, H. G. Effect of three decades of screening mammography on breast-cancer incidence. *N. Engl. J. Med.* **367**, 1998–2005 (2012).
- Independent, U. K. Panel on breast cancer screening. The benefits and harms of breast cancer screening: an independent review. *Lancet* **380**, 1778–1786 (2012).
- Sprague, B. L. et al. Time-varying risks of second events following a DCIS diagnosis in the population-based Vermont DCIS cohort. *Breast Cancer Res. Treat.* **174**, 227–235 (2019).
- Gorringe, K. L. & Fox, S. B. Ductal carcinoma in situ biology, biomarkers, and diagnosis. *Front. Oncol.* **7**, 248 (2017).
- Pang, J.-M. B. et al. Breast ductal carcinoma in situ carry mutational driver events representative of invasive breast cancer. *Mod. Pathol.* **30**, 952–963 (2017).
- Tate, J. G. et al. COSMIC: the catalogue of somatic mutations in cancer. *Nucleic Acids Res.* **47**, D941–D947 (2019).
- Parker, J. S. et al. Supervised risk predictor of breast cancer based on intrinsic subtypes. *J. Clin. Oncol.* **27**, 1160–1167 (2009).
- Lin, C.-Y. et al. Genomic landscape of ductal carcinoma in situ and association with progression. *Breast Cancer Res. Treat.* **178**, 307–316 (2019).
- Nagasawa, S. et al. Genomic profiling reveals heterogeneous populations of ductal carcinoma in situ of the breast. *Commun. Biol.* **4**, 438 (2021).
- Pareja, F. et al. Whole-Exome sequencing analysis of the progression from non-low-grade ductal carcinoma in situ to invasive ductal carcinoma. *Clin. Cancer Res.* **26**, 3682–3693 (2020).
- Abba, M. C. et al. A molecular portrait of high-grade ductal carcinoma in situ. *Cancer Res.* **75**, 3980–3990 (2015).
- Casasent, A. K. et al. Multiclonal invasion in breast tumors identified by topographic single. *Cell Sequencing. Cell* **172**, 205–217.e12 (2018).
- Gerdes, M. J. et al. Single-cell heterogeneity in ductal carcinoma in situ of breast. *Mod. Pathol.* **31**, 406–417 (2018).
- Pruneri, G. et al. The prevalence and clinical relevance of tumor-infiltrating lymphocytes (TILs) in ductal carcinoma in situ of the breast. *Ann. Oncol.* **28**, 321–328 (2017).
- Campbell, M. J. et al. Characterizing the immune microenvironment in high-risk ductal carcinoma in situ of the breast. *Breast Cancer Res. Treat.* **161**, 17–28 (2017).
- Trinh, A. et al. Genomic alterations during the in situ to invasive ductal breast carcinoma transition shaped by the immune system. *Mol. Cancer Res.* **19**, 623–635 (2021).
- Lesurf, R. et al. Molecular features of subtype-specific progression from ductal carcinoma in situ to invasive breast cancer. *Cell Rep.* **16**, 1166–1179 (2016).
- Gil Del Alcazar, C. R. et al. Immune escape in breast cancer during in situ to invasive carcinoma transition. *Cancer Discov.* **7**, 1098–1115 (2017).
- Allen, M. D., Marshall, J. F. & Jones, J. L. $\alpha\text{v}\beta 6$ Expression in myoepithelial cells: a novel marker for predicting DCIS progression with therapeutic potential. *Cancer Res.* **74**, 5942–5947 (2014).
- Delort, L. et al. The adipose microenvironment dysregulates the mammary myoepithelial cells and could participate to the progression of breast cancer. *Front Cell Dev. Biol.* **8**, 571948 (2020).
- Allinen, M. et al. Molecular characterization of the tumor microenvironment in breast cancer. *Cancer Cell* **6**, 17–32 (2004).

23. Hu, M. et al. Regulation of in situ to invasive breast carcinoma transition. *Cancer Cell* **13**, 394–406 (2008).
24. Unsworth, A., Anderson, R. & Britt, K. Stromal fibroblasts and the immune microenvironment: partners in mammary gland biology and pathology? *J. Mammary Gland Biol. Neoplasia* **19**, 169–182 (2014).
25. Sinha, V. C. & Piwnica-Worms, H. Intratumoral heterogeneity in ductal carcinoma in situ: chaos and consequence. *J. Mammary Gland Biol. Neoplasia* **23**, 191–205 (2018).
26. Nachmanson, D. et al. Mutational profiling of micro-dissected pre-malignant lesions from archived specimens. *BMC Med. Genom.* **13**, 173 (2020).
27. Foley, J. W. et al. Gene expression profiling of single cells from archival tissue with laser-capture microdissection and Smart-3SEQ. *Genome Res.* **29**, 1816–1825 (2019).
28. Martincorena, I. et al. Tumor evolution. High burden and pervasive positive selection of somatic mutations in normal human skin. *Science* **348**, 880–886 (2015).
29. Martincorena, I. et al. Somatic mutant clones colonize the human esophagus with age. *Science* **362**, 911–917 (2018).
30. Cancer Genome Atlas Network. Comprehensive molecular portraits of human breast tumours. *Nature* **490**, 61–70 (2012).
31. Bielski, C. M. et al. Genome doubling shapes the evolution and prognosis of advanced cancers. *Nat. Genet.* **50**, 1189–1195 (2018).
32. D'Antonio, M., Tamayo, P., Mesirov, J. P. & Frazer, K. A. Kataegis expression signature in breast cancer is associated with late onset, better prognosis, and higher HER2 levels. *Cell Rep.* **16**, 672–683 (2016).
33. Afzaljavan, F., Sadr, A. S., Savas, S. & Pasdar, A. GATA3 somatic mutations are associated with clinicopathological features and expression profile in TCGA breast cancer patients. *Sci. Rep.* **11**, 1679 (2021).
34. Emmanuel, N. et al. Mutant GATA3 actively promotes the growth of normal and malignant mammary cells. *Anticancer Res.* **38**, 4435–4441 (2018).
35. Kader, T. et al. Atypical ductal hyperplasia is a multipotent precursor of breast carcinoma. *J. Pathol.* **248**, 326–338 (2019).
36. Kader, T. et al. The genetic architecture of breast papillary lesions as a predictor of progression to carcinoma. *NPJ Breast Cancer* **6**, 9 (2020).
37. Cai, Y. et al. Loss of chromosome 8p governs tumor progression and drug response by altering lipid metabolism. *Cancer Cell* **29**, 751–766 (2016).
38. Thompson, E. et al. The immune microenvironment of breast ductal carcinoma in situ. *Mod. Pathol.* **29**, 249–258 (2016).
39. Danforth, D. N. Jr Genomic changes in normal breast tissue in women at normal risk or at high risk for breast cancer. *Breast Cancer* **10**, 109–146 (2016).
40. Cancer Genome Atlas Research Network. et al. The cancer genome atlas pan-cancer analysis project. *Nat. Genet.* **45**, 1113–1120 (2013).
41. Risom, T. et al. Transition to invasive breast cancer is associated with progressive changes in the structure and composition of tumor stroma. Preprint at *bioRxiv* <https://doi.org/10.1101/2021.01.05.425362> (2021).
42. Patel, A. P. et al. Single-cell RNA-seq highlights intratumoral heterogeneity in primary glioblastoma. *Science* **344**, 1396–1401 (2014).
43. Roden, D. L. et al. Single cell transcriptomics reveals molecular subtype and functional heterogeneity in models of breast cancer. Preprint at *bioRxiv* <https://doi.org/10.1101/282079> (2018).
44. Allred, D. C. et al. Ductal carcinoma in situ and the emergence of diversity during breast cancer evolution. *Clin. Cancer Res.* **14**, 370–378 (2008).
45. Sun, R., Hu, Z. & Curtis, C. Big bang tumor growth and clonal evolution. *Cold Spring Harb. Perspect. Med.* <https://doi.org/10.1101/cshperspect.a028381> (2018).
46. Polyak, K. Is breast tumor progression really linear? *Clin. Cancer Res.: Off. J. Am. Assoc. Cancer Res.* **14**, 339–341 (2008).
47. Zeng, Z. et al. Somatic genetic aberrations in benign breast disease and the risk of subsequent breast cancer. *NPJ Breast Cancer* **6**, 24 (2020).
48. Silverstein, M. J. The University of Southern California/Van Nuys prognostic index for ductal carcinoma in situ of the breast. *Am. J. Surg.* **186**, 337–343 (2003).
49. Mannu, G. S. et al. Invasive breast cancer and breast cancer mortality after ductal carcinoma in situ in women attending for breast screening in England, 1988–2014: population based observational cohort study. *BMJ* **369**, m1570 (2020).
50. Gail, M. H. et al. Projecting individualized probabilities of developing breast cancer for white females who are being examined annually. *J. Natl Cancer Inst.* **81**, 1879–1886 (1989).
51. Lee, A. et al. BOADICEA: a comprehensive breast cancer risk prediction model incorporating genetic and nongenetic risk factors. *Genet. Med.* **21**, 1708–1718 (2019).
52. Kos, Z. et al. Pitfalls in assessing stromal tumor infiltrating lymphocytes (sTILs) in breast cancer. *NPJ Breast Cancer* **6**, 17 (2020).
53. Hendry, S. et al. Assessing tumor-infiltrating lymphocytes in solid tumors: a practical review for pathologists and proposal for a standardized method from the international immunooncology biomarkers working group: part 1: assessing the host immune response, TILs in invasive breast carcinoma and ductal carcinoma in situ, metastatic tumor deposits and areas for further research. *Adv. Anat. Pathol.* **24**, 235–251 (2017).
54. Cornel, A. M., Mimpfen, I. L. & Nierkens, S. MHC class I downregulation in cancer: underlying mechanisms and potential targets for cancer immunotherapy. *Cancers* **12**, 1760 (2020).
55. Garrido, M. A. et al. HLA class I alterations in breast carcinoma are associated with a high frequency of the loss of heterozygosity at chromosomes 6 and 15. *Immunogenetics* **70**, 647–659 (2018).
56. Campbell, M. J. et al. Abstract 961: Intralesional injection of anti-PD-1 (pembrolizumab) results in increased T cell infiltrate in high risk DCIS. *Cancer Res.* **78**, 961–961 (2018).
57. Dobin, A. et al. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15–21 (2013).
58. Liao, Y., Smyth, G. K. & Shi, W. featureCounts: an efficient general purpose program for assigning sequence reads to genomic features. *Bioinformatics* **30**, 923–930 (2014).
59. Stuart, T. et al. Comprehensive Integration of Single-Cell Data. *Cell* **177**, 1888–1902.e21 (2019).
60. Hafemeister, C. & Satija, R. Normalization and variance stabilization of single-cell RNA-seq data using regularized negative binomial regression. *Genome Biol.* **20**, 296 (2019).
61. Gendoo, D. M. A. et al. Genefu: an R/Bioconductor package for computation of gene expression-based signatures in breast cancer. *Bioinformatics* **32**, 1097–1099 (2016).
62. Subramanian, A. et al. Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proc. Natl Acad. Sci. USA* **102**, 15545–15550 (2005).
63. Barbie, D. A. et al. Systematic RNA interference reveals that oncogenic KRAS-driven cancers require TBK1. *Nature* **462**, 108–112 (2009).
64. Fabregat, A. et al. The Reactome pathway knowledgebase. *Nucleic Acids Res.* **46**, D649–D655 (2018).
65. Liberzon, A. et al. Molecular signatures database (MSigDB) 3.0. *Bioinformatics* **27**, 1739–1740 (2011).
66. Didion, J. P., Martin, M. & Collins, F. S. Atropos: specific, sensitive, and speedy trimming of sequencing reads. *PeerJ* **5**, e3720 (2017).
67. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. Preprint at <https://arxiv.org/abs/1303.3997> (2013).
68. Tischler, G. & Leonard, S. biobambam: tools for read pair collation based algorithms on BAM files. *Source Code Biol. Med.* **9**, 13 (2014).
69. Li, H. et al. The sequence alignment/map format and SAMtools. *Bioinformatics* **25**, 2078–2079 (2009).
70. Lai, Z. et al. VarDict: a novel and versatile variant caller for next-generation sequencing in cancer research. *Nucleic Acids Res.* **44**, e108 (2016).
71. Van der Auwera, G. A. & O'Connor, B. D. *Genomics in the Cloud: Using Docker, GATK, and WDL in Terra.* (O'Reilly Media, Inc., 2020).
72. Guimera, R. V. bcbio-nextgen: Automated, distributed next-gen sequencing pipeline. *EMBnet. J.* **17**, 30 (2011).
73. Cingolani, P. et al. A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w1118; iso-2; iso-3. *Fly* **6**, 80–92 (2012).
74. 1000 Genomes Project Consortium. et al. An integrated map of genetic variation from 1,092 human genomes. *Nature* **491**, 56–65 (2012).
75. Lek, M. et al. Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536**, 285–291 (2016).
76. Landrum, M. J. et al. ClinVar: improving access to variant interpretations and supporting evidence. *Nucleic Acids Res.* **46**, D1062–D1067 (2018).
77. Karczewski, K. J. et al. The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature* **581**, 434–443 (2020).
78. Talevich, E., Shain, A. H., Botton, T. & Bastian, B. C. CNVkit: genome-wide copy number detection and visualization from targeted DNA sequencing. *PLoS Comput. Biol.* **12**, e1004873 (2016).
79. Van Loo, P. et al. Allele-specific copy number analysis of tumors. *Proc. Natl Acad. Sci. USA* **107**, 16910–16915 (2010).
80. Nilsen, G. et al. Copynumber: Efficient algorithms for single- and multi-track copy number segmentation. *BMC Genom.* **13**, 591 (2012).
81. Islam, S. M. A. et al. Uncovering novel mutational signatures by de novo extraction with SigProfilerExtractor. Preprint at *bioRxiv* <https://doi.org/10.1101/2020.12.13.422570> (2021).
82. Alexandrov, L. B. et al. The repertoire of mutational signatures in human cancer. *Nature* **578**, 94–101 (2020).
83. Kim, S. et al. Evaluating tumor evolution via genomic profiling of individual tumor spheroids in a malignant ascites from a patient with ovarian cancer using a laser-aided cell isolation technique. *Sci. Rep.* <https://doi.org/10.1101/282277> (2018).
84. Schliep, K. P. phangorn: phylogenetic analysis in R. *Bioinformatics* **27**, 592–593 (2011).
85. Reiter, J. G. et al. Reconstructing metastatic seeding patterns of human cancers. *Nat. Commun.* **8**, 14114 (2017).
86. Mori, H. et al. Characterizing the tumor immune microenvironment with Tyramide-based multiplex immunofluorescence. *J. Mammary Gland Biol. Neoplasia* **25**, 417–432 (2020).

87. Bankhead, P. et al. QuPath: open source software for digital pathology image analysis. *Sci. Rep.* **7**, 16878 (2017).

ACKNOWLEDGEMENTS

We are grateful to Drs. Michael Campbell, Christina Yau, and Kathleen Curtius for helpful conversations, Drs. Alfredo Molinolo, Oluwale Fadare, Sharneela Kaushal, Valeria Estrada, and Mrs. Kimberly McIntyre for their support and assistance in the tissue collection, preparation, and dissection, Mrs. Eliza Jeong, Marcy Andersen, and Nicole Lee for their assistance retrieving clinical information. We thank the technical assistance of the Vermont Integrative Genomics Resource Massively Parallel Sequencing Facility with the combined support of the University of Vermont Cancer Center, Lake Champlain Cancer Research Organization, UVM College of Agriculture and Life Sciences, and the UVM Larner College of Medicine. We acknowledge the work of all working groups from the NCI Consortium for the Molecular Characterization of Screen Detected Lesions (MCL), in particular Christopher Amos, Daniel Crichton, Heather Kincaid, Kirsten Anton, Luca Cinquini, and David Liu for their assistance in the data management and sharing. Figure 1 was created with BioRender.com and printed with permission. This work is supported by funding from the National Institute of Health (U01CA196406, U01CA196406-03S1, U01CA196383, R01DE026644, T32GM008806, T15LM011271), the National Cancer Institute (P30CA023100), and the California Tobacco Related Disease Research Program pre-doctoral fellowship to DN (28DT-0011). The funding bodies had no role in the design of the study; collection, analysis, and interpretation of data; or in the writing of the manuscript.

AUTHOR CONTRIBUTIONS

J.L.S., G.S.S., G.L.H., L.J.E., A.D.B., and O.H. designed the study, F.H., A.D.B., D.L.W., B.L.S., and T.O.K. selected and annotated the specimen. J.S., K.J., H.Y., H.M., M.F.E., and J.G. generated the data, O.H., D.N., A.O., and H.M. analyzed the data, O.H. and D.N. wrote the manuscript. All authors reviewed and approved the manuscript.

COMPETING INTERESTS

The authors declare no competing interests.

ADDITIONAL INFORMATION

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41523-021-00365-y>.

Correspondence and requests for materials should be addressed to Olivier Harismendy.

Reprints and permission information is available at <http://www.nature.com/reprints>

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2022