# Eco-evolutionary Guided Pathomics Analysis to Predict DCIS Upstaging

Yujie Xiao[1], Manal Elmasry[2,3], Ji Dong K. Bai[2], Andrew Chen[2], Yuzhu Chen[2], Brooke Jackson[4],

Joseph O. Johnson[4], Robert J. Gillies[4], Prateek Prasanna[5], Chao Chen[5*], Mehdi Damaghi[1,2,3*]

1- Department of Applied Mathematics and Statistics, Stony Brook University, NY, USA
2- Department of Pathology, Stony Brook Medicine, Stony Brook University, NY, USA
3- Department of Pathology, Faculty of Medicine, Mansoura University, Mansoura, Egypt
4- Moffitt Cancer Center, Tampa, Fl, USA
5-Department of Biomedical Informatics, Stony Brook Medicine, Stony Brook University, NY, USA

**Running title:** Eco-Evolutionary Designed Biomarker

* Corresponding authors:  Chao Chen: chao.chen.1@stonybrook.edu,

Mehdi Damaghi: mehdi.damaghi@stonybrookmedicine.edu
Stony Brook Cancer Center
MART building, 9th floor, 9M0802
Lauterbur Drive
Stony Brook, NY 11794-7263
Phone: (631) 216-2920

**Declaration of interests**

The authors declare no competing interest.

## Abstract

Cancers evolve in a dynamic ecosystem. Thus, characterizing cancer's ecological dynamics is crucial to understanding cancer evolution and can lead to discovering novel biomarkers to predict disease progression. Ductal carcinoma in situ (DCIS) is an early-stage breast cancer characterized by abnormal epithelial cell growth confined within the milk ducts. In this study, we show that ecological habitat analysis of hypoxia and acidosis biomarkers can significantly improve prediction of DCIS upstaging. First, we developed a novel eco-evolutionary designed approach to define habitats in the tumor intra-ductal microenvironment based on oxygen diffusion distance. Then, we identified cancer cells with metabolic phenotypes attributed to their habitat conditions, such as the expression of CA9 indicating hypoxia responding phenotype, and LAMP2b indicating the acid adaptation. Traditionally these markers have shown limited predictive capabilities for DCIS upstaging, if any. However, when analyzed from an ecological perspective, their power to differentiate between pure DCIS and upstaged DCIS increased significantly. Second, using eco-evolutionary guided computational and digital pathology techniques, we discovered distinct niches with spatial patterns of these biomarkers and used the distribution of such niches to predict patient upstaging. The niches patterns were characterized by pattern analysis of both cellular and spatial features. With a 5-fold validation on the biopsy cohort, we trained a random forest classifier to achieve the area under curve (AUC) of 0.74. Our results affirm the importance of using eco-evolutionary-designed approaches in biomarkers discovery studies in the era of digital pathology by demonstrating the role of tumor ecological habitats and niches.

**Keywords:**

Tumor ecology and evolution, DCIS, Eco-evolutionary biomarkers, Metabolic phenotypes, Habitat analysis, Niche analysis, Pathomics, Machine learning, Digital pathology

## Introduction:

In recent years, the understanding that cancer is a dynamic ecological and evolutionary process has become deeply entrenched (1,2,3). To date, several evolutionary approaches have been adapted and applied in cancer biology, such as diversity measures to predict disease progression; however, tumor ecosystem and ecological habitat and niche studies are still overlooked (3,4). Within the human body and much like organisms in the natural world, cancer cells follow evolutionary principles, utilizing resources and establishing habitats and niches within tissues (5,6). This ecological perspective of cancer is crucial for discovering the natural selection driving cancer evolution. Recognizing the parallels between organismal ecology and the tumor microenvironment opens up untapped opportunities to incorporate ecological measures, improving our understanding of both tumor dynamics and selective pressures shaping tumors' evolutionary landscapes. Such insights may potentially lead to improved cancer prognosis, progression prediction, risk stratification, and therapeutic strategies. If tumor evolutionary state and/or its evolutionary trajectories could be reliably achieved using a single biopsy tissue, clinical translation would be comparatively more manageable. Nevertheless, studies have yet to determine whether measures of tumor evolvability derived from a single biopsy sample are adequate, or if the inclusion of multiple samples significantly enhances predictions of clinical outcomes (7).

Breast cancer incidence in the US has been increasing over the past decade at a rate of 0.5% per year(8). With increased mammographic screening, there has been a substantial increase in detecting the early non-invasive forms of breast cancer, such as ductal carcinoma in situ (DCIS)(2,9). About one-third of breast cancers detected by mammography are DCIS (10). As the most common pre-cancer state, DCIS can progress to invasive disease in a linear evolution pattern, or can be part of other clonal evolutionary dynamics such as branching, punctuated, or neutral evolution (2,9,11). Since DCIS and IDC (invasive ductal carcinoma) are indistinguishable by (epi-)genetic mutations, gene expression, or protein biomarkers, and because it is not possible to predict whether DCIS will remain indolent or progress to more aggressive disease, almost all early tumors are treated with aggressive interventions(2,12–14). To avoid such over treatment, more research is needed to fully understand evolution from pre-cancer to indolent DCIS or progress to IDC(9).

DCIS is a heterogeneous group of neoplastic lesions confined to the mammary ducts. The confinement of proliferating neoplastic cells inside the duct and growth of pre-cancer cells toward the center of the duct, which is far from vasculature, causes limitations in oxygen and nutrients. This intraductal oxygen microenvironment is also influenced by complex ecosystems surrounding the duct, such as vascular activity(15), stiffness of extracellular matrix (ECM) (16), and metabolites (6,17,18,19) (**Figure 1A**) . Local microinvasion is the main difference between DCIS and IDC and might also be the first evolutionary step of progressing in the case of linear evolution(11). Microinvasion consists of cohorts of cancer cells that breach the basement membrane into the surrounding ECM. Recently, genomic analysis of matched DCIS and IDC samples has revealed that in 75% of cases, the invasive recurrence was found to be clonally related to the initial DCIS. This implies that tumor cells derived from DCIS could evolve in a linear or branching fashion with 18% new transformations and/or clonogenesis (11). These new findings emphasize the extraordinary heterogeneity in genotype and phenotypic plasticity in breast cancer that must be studied in the light of evolution and ecological studies. Thus, we designed our study to capture the phenotypic heterogeneity of cancer cells in their selective microenvironments. We hypothesize that non-genetic ecological factors, such as intra-ductal microenvironmental conditions, may

2

89 be responsible for transitioning from a DCIS to IDC phenotype, in the case of linear and branching
90 evolution, or may select clones with pre-existing IDC phenotypes in the case of the other evolutionary
91 trajectories, including punctuated and neutral evolution(6,11,18,20).
92 To validate this hypothesis, we propose a novel method to study DCIS evolution, by capturing and
93 characterizing "tumor habitats" and "cell niches" and their interactions in the tumor ecosystem. Natural
94 selection requires phenotypic diversity within a population undergoing microenvironmental selection
95 forces (21). Cells that adapt in response to natural selection may present similar phenotypes,
96 corresponding to the microenvironment exerting the selection.  We started by defining the habitats based
97 on availability of oxygen into: a) oxygenated habitat and b) hypoxic habitat. Following previous
98 theory(18,22), these habitats are defined by distance from the duct boundary. However, a uniform
99 distance threshold hardly captures the true oxidate/hypoxic states of cells. Therefore, we further proposed
100 to fine-tune these habitats using protein expression indicative of phenotypes resulting from cancer cell
101 adaptation to variation in oxygen availability. Therefore, we defined *intraductal DCIS niches* inside
102 habitats as clusters of cells with similar phenotypic behavior responding to hypoxia. Through analysis
103 via these niches, we can identify more aggressive phenotypes leading to microinvasion and DCIS
104 upstaging to IDC or possible direct evolution to IDC without going through DCIS sub-stages.
105 Our biomarkers are designed based on prior biological knowledge. Oxygen availability determines the
106 source of energy production as of either mitochondrial respiration or glycolysis. Hypoxic cells switch to
107 glycolysis, causing lactic acid production that can lead to acidosis when lactic acid is not cleared from
108 the tumor space. Peri-luminal cells will experience hypoxia if they are far (>0.125 - 0.160 mm) from a
109 blood supply.  These cancer cells inhabit a microenvironment of hypoxia, acidosis, and severe nutrient
110 deprivation (18,22). These environmental properties exert a strong selection pressure upon the cancer
111 cells, which in turn feeds back to the microenvironment, creating a dynamically changing tumor
112 ecosystem containing several habitats. We have shown that cancer cells within breast ducts subjected to
113 chronic hypoxia and acidosis evolve mechanisms of adaptations to survive in this harsh
114 microenvironment (17,18,20). We have also shown that cells adapted to hypoxic and/or acidic niches
115 have developed specific metabolic vulnerabilities that can be targeted to push them back to a more
116 physiologically normal state(17). Both these studies strengthen the acid-induced evolution model of
117 breast cancer and our proposed evolutionary designed biomarkers including CA9 and LAMP2b in this
118 research(6,17,20,23,24). Here we examined the role of these biomarkers within an eco-evolutionary
119 concept as a predictor of DCIS upstaging for the first time. We used these markers as representative of
120 the cancer cell metabolic states to define niches inside habitats that can select for more aggressive
121 phenotypes, leading to microinvasion and DCIS upstaging to IDC or possible direct evolution to IDC
122 without going through DCIS sub-stages.
123 To perform our analysis, we curated a retrospective cohort of DCIS patients, with specimens collected
124 from Biopsy (Bx) samples before surgery and after Excision (Ex).  All the patients had histologically
125 confirmed DCIS on core biopsy, followed by diagnosis confirmed on surgical excision specimens with
126 either DCIS or IDC (**Figure 1B**). Our niche-based prediction model is trained and tested on the Bx
127 samples. These best fits future clinical applications that machine learning model can be subsequently
128 applied to predict upstaging at Bx for future patients. We then stained 3 sequentially sectioned slides for
129 hematoxylin and eosin (HE), CA9 and LAMP2b. We manually annotated ducts bigger than 400 μms in
130 diameter. The 200 μms in radius annotation ensures each duct has both oxygenated and hypoxic habitats
131 to build a balanced cohort for analysis. We developed a novel algorithm to detect intra-ductal DCIS cell

132   niches based on biomarker expression similarity. Then, we studied the spatial organization of CA9- and
133   LAMP2b-positive cells as the eco-evolution markers of cancer cells in hypoxic and acidic habitats at
134   three different scales: whole slide, duct, and oxygen habitats (normoxic and hypoxic). We also applied
135   multiple spatial functions and spatial entropies were used to define niche and micro-niches describing
136   the spatial patterns of the cell groups. After a systematic and comprehensive analysis, we observed that
137   the spatial features at the finest habitat level possess the most predictive power where the micro-niches
138   were defined by the expression of CA9 and LAMP2b in hypoxic habitats. By characterizing these niches
139   and micro-niches with spatial and pathomics features, we then developed a risk scoring system by
140   integrating principles of ecological-evolutionary dynamics with pathological imaging and molecular
141   features of early-stage breast tumors (**Figure 1C**).  We show that quantitative analyses of immuno-
142   histological images combined with the tumor's eco-evolution dynamics and underlying molecular
143   pathophysiology can significantly improve predicting if the neoplasm has already evolved to invasive
144   disease and cancer. We developed a machine learning model fine-tuning the tumor habitats into micro-
145   niches using specific molecular signatures of resident cancer cells to provide informed decision support.
146   In summary, we show that specific habitats containing micro-niches of cells with similar phenotypes
147   responding to hypoxia and acidosis, or adaptation to long term exposure of these conditions, are
148   responsible for DCIS progression, and hence would be correlated to upstaging. To test this hypothesis,
149   we applied machine learning techniques to calculate the niches inside the tumor to define spatial and
150   temporal distribution of habitats in solid tumors of DCIS patients with pure DCIS and upstaged disease.
151   By deploying eco-evolutionary principles and machine learning techniques, our work proposes a novel
152   consilient approach - as opposed to the traditional single biomarker studies - to stratify DCIS patients

## Materials and Methods

### Method overview

155   Our evolutionary analysis pipeline takes 3 consecutive slides of each patient sample, detects intra-ductal
156   cell niches, characterizes these niches with their spatial and morphological features, and then predicts
157   whether the patient will be pure DCIS or upstaged based on the distribution of these niches. In particular,
158   the pipeline has 4 modules. First, we annotate and align ducts from different whole slide images (WSIs)
159   of the same patient sample. This ensures cells of different slides are aligned and we can characterize their
160   interactions. In the second module, we detect and map all eco-evo positive cells (i.e., cells activated with
161   the selected stains) into the same duct and detect different clusters of cells as niches. In the third module,
162   we characterize these niches with comprehensive spatial statistical features, as well as their
163   morphological features as observed in HE. Finally, we categorize these niches into different subclasses
164   through deep learning-based dimension reduction and clustering based on their features. We use the
165   distribution of different niche subclasses to characterize different samples/patients. We demonstrate the
166   discriminative power of this niche-based characterization in predicting whether a patient will be pure
167   DCIS or upstaged in the future. **Figure 1C** illustrates the overview of our pipeline.

### Data preparation and usage

169   The data used in this study is the biopsy samples collected after mammography and before surgery. 84
170   samples including 68 pure DCIS and 16 progressed to IDC were analyzed. This study complied with the
171   Health Insurance Portability and Accountability Act and was approved by the institutional review board,
172   with a waiver of the requirement for informed consent. Women with a core biopsy diagnosis of DCIS
173   between 2012 and 2022 who consented to at Moffitt Cancer Center Total Cancer Care protocol were

174 included in this analysis. Cases were excluded if surgical excision was performed more than 6 months
175 after the core biopsy, if there was concurrent ipsilateral invasive breast cancer or metastatic malignancy,
176 or if neoadjuvant chemotherapy (for a concurrent contralateral breast malignancy) or chemotherapy for
177 a non-breast primary malignancy was administered between the dates of the DCIS core biopsy and
178 surgery. Additional exclusions included a personal history of invasive breast cancer or DCIS within 12
179 months preceding the core biopsy or a concurrent diagnosis of Paget disease in the ipsilateral breast.
180 After applying these inclusion and exclusion criteria, 84 cases of biopsy-proven DCIS were identified,
181 of which 16 were upgraded at surgery and 68 remained non-upgraded.
182 Pure DCIS and upstaged patients were matched across clinical features, including age, race, ethnicity,
183 grade, ER status, and PR status, to minimize their influence on the analysis (**Figure S1**). To validate the
184 comparability of these groups, we conducted a Wilcoxon rank-sum test for the continuous variable (age)
185 and chi-square tests for the categorical variables (race, ethnicity, grade, ER status, and PR status). None
186 of these tests showed significant differences between the two groups, with all p-values larger than 0.1,
187 indicating that the groups were well-matched.
188 For each sample, we obtained 3 whole slide images, including 1 HE and 2 IHC slides. We conducted 5-
189 fold stratified cross validation, where 4 folds are used for niche clustering and for the training of the pure
190 DCIS/upstaged classifier and 1-fold is used for validation. This fits the clinical application we are aiming
191 for; we would like our model to estimate the risk based on biopsy samples, which are much less invasive
192 and can be used for patient stratifications before surgery and hopefully decrease over treatment. Further
193 details on HE and IHC acquisition are provided below.
194 **Sample selection, immunohistochemistry and HE staining.** Patients' tumor blocks were selected by
195 pathologists using the archived HE stained slides. The blocks were sequentially sectioned 4 µms and de-
196 identified for research use. 3 slides were stained with primary antibodies of 1:100 dilution of anti-LAMP2
197 (#ab18529, Abcam), and 1 ug/ml concentration of anti-CA9 (#AF2188, R&D), and HE staining using
198 standard hematoxylin and eosin protocol. Positive and negative controls were used. Normal placenta was
199 used as a positive control for LAMP2b and clear cell renal cell carcinoma was used as a positive control
200 for CA9. For the negative control, an adjacent section of the same tissue was stained without application
201 of primary antibody and any stain pattern observed was considered as non-specific binding of the
202 secondary. Primary immunohistochemical analysis was conducted using digitally scanning slides. The
203 scoring method used by the pathologist reviewer to determine (a) the degree of positivity scored the
204 positivity of each sample ranged from 0 to 3 and was derived from the product of staining intensity (0–
205 3+). A zero score was considered negative, score 1 was weak positive, score 2 was moderate positive,
206 and score 3 was strong positive. (b) The percentage of positive tumors stained (on a scale of 0-3). Whole
207 slide imaging (WSI) of IHC and HE slides were obtained by scanning at 20X magnification (of 0.5022
208 micrometer per pixel) using Aperio AT2 from Leica Biosystems. Images were transferred to cloud
209 storage and locally to be uploaded in QuPath software for analysis. QuPath software was used to detect
210 the positive pixels for each IHC marker (CA9 and LAMP2b) and to segment the HE images into hypoxic
211 and normoxic tumor habitats based on their distance from the basement membrane. The 'Positive Cell
212 Detection' function from Qupath was used to automatically classify the positivity of CA9 and LAMP2b
213 markers and validated by the study pathologist.

214 **MODULE 1: Duct annotation and alignment**
215 **Manual annotation of ducts in the Bx cohort.** We annotate and align ducts within all input slides (1
216 HE + 2 IHCs per sample). After annotating ducts, we align the ducts from the three modalities via co-

5

217 registration. This alignment enables us to map cells into the same spatial domain and analyze their
218 interaction. Details are provided below. QuPath was used as the interface to annotate ducts by the
219 pathologist (Dr. Bai) and the trained students and reviewed by D. Damaghi. We annotate ducts from
220 WSIs of all three modalities. To ensure best characterization, we only identify ducts of >400 μms
221 diameter, with visible myoepithelial layer and basement membrane. Following this, based on distance,
222 each duct was annotated with four layers: adjacent stroma, oxidative/normoxia, hypoxic/hypoxia, and
223 necrosis. Adjacent stroma was defined as the stroma up to 125 μms outside a given duct. Within the duct,
224 necrosis was defined as any area containing dead cells, as identified by a lack of nuclei. Oxidative layer
225 was defined as the area containing cells inside the duct within 125 μms of the basement membrane.
226 Hypoxia was defined as the area containing cells inside the duct further than 125 μms from the basement
227 membrane. The annotations were done for all 84 samples in the Bx cohort, and then were exported as
228 standard GeoJSON files.
229 **Co-registration.** To characterize the interactions of different modalities from single-plexed slides, an
230 alignment strategy was utilized. We register both CA9 and LAMP2b IHC slides towards the HE slides.
231 A direct co-registering at the whole slide level with manual landmarks does not give us satisfactory
232 alignment at each duct, due to the variable deformations across slides. We further co-register the slides
233 in a duct-by-duct fashion. Using initially registered whole slides, and spatial proximity, we identify the
234 corresponding ducts at the HE and 2 IHC slides. Next, we register both the CA9 duct and LAMP2b duct
235 into the corresponding HE ducts. We use Virtual Alignment of pathology Image Series (VALIS), which
236 provides a fully automated pipeline to register whole slide images (WSI) using rigid and/or non-rigid
237 transformations (34). For each sample, we chose non-rigid registration and registered the ducts from CA9
238 and LAMP2b towards the reference HE ducts. The co-registration procedure and the qualitative results
239 are shown in **Figure S4 and S5**. The co-registration provides a mapping of any cells detected in CA9 or
240 LAMP2b towards a shared spatial domain, enabling the analysis of their interactions.

## MODULE 2: Cell and niche detection

242 **Cell detection.** With the duct annotations in place, we automatically detect cells from the 2 IHCs and
243 determine if they are positive in CA9 or LAMP2b based on their intensities. As we are only interested in
244 intra-ductal cell niches, we only detect cells within each duct. For each IHC duct, we detect cells using
245 Qupath watershed cell detection algorithm (25). Based on the intensity level, we categorize the cells into
246 4 groups: 'Negative', '1+', '2+', and '3+'. The detection of cells within a duct is done by starDist (25,35)
247 extension in Qupath on HE slide.
248 **Graph construction for niche detection.** After annotating all of the positive cells (i.e., CA9 or LAMP2b
249 positive cells), they were mapped on HE slides, enabling us to detect niches on HE slides. Since there is
250 a large amount of positive cells within each duct, with diverse spatial context and morphological features,
251 we construct a graph with these cells by connecting cells whose distances are smaller than a certain
252 threshold and detect connected components of the graph as representatives of cells living in "niches".
253 Multiple thresholds have been experimented and an optimum value is selected based on performance.
254 Each positive cell niche is supposed to have a similar eco-evo phenotype and be spatially coherent.
255 Therefore, we overlay both CA9 positive and LAMP2b positive cells into the same domain as an
256 approximation of the local eco-evo cell distribution (**Figure S6**). This gives us the opportunity to measure
257 their interaction via spatial statistical functions as defined later. Based on the same principle, we use cell
258 morphological features extracted from HE within the region of each niche to characterize the niche.

## MODULE 3: niche characterization and feature extraction

Once niches are detected. We extracted both spatial and morphological features to characterize them. To describe the spatial interaction patterns, we utilized various spatial functions as features. We also extract cell features consisting of morphology features and texture features that are commonly adopted in HE image analysis.

**Cellular features.** For cellular features we measured both morphological and texture features. The morphological features include area, eccentricities, circularity, elongation, extent, major axis length, minor axis length, solidity and curvature. The texture features include angular second moment (ASM) of co-occurrence matrix, contrast, correlation, entropy, homogeneity and intensity. All features were calculated following the implementations in the sc-MTOP(36) package.

Although we do not have exact cell-to-cell correspondence between the cells within a niche and cells detected in HE, we still can aggregate morphological and texture features within the proxy of the cells part of a niche to characterize the niche. For each niche, we identify the concave hull region enclosing its eco-evo positive cells within a duct on HE slide. Next, we aggregate cell features across all HE-detected cells within the corresponding region. For each cell feature dimension, we calculated its mean, standard deviation, maximum, minimum, kurtosis and skewness.

**Spatial features.** We extract various spatial statistical functions (37) to characterize residingcells and their interactions inside habitats to define niches. These functions are listed below:

G Function: The G function, denoted as G(r), is the cumulative distribution function of nearest-neighbor distance. The G function provides insights into the clustering or dispersion behavior of the point pattern.

$$G(r) = P\{d(u, X\backslash u) \leq r | u \in X\}, d(\bullet) \text{ is the minimum distance}$$

F Function: The F function, known as the empty space function, is the cumulative distribution function of the empty-space distance. The F function is commonly used to assess the regularity or inhibition patterns in point patterns.

$$F(r) = P\{d(u, X) \leq r\}, d(\bullet) \text{ is the minimum distance}$$

K Function: Ripley's K function, denoted as K(r), is a measure of second-order intensity or spatial interaction. It assesses whether points tend to be more clustered or dispersed within a certain distance r compared to a CSR process. It considers both the distance and intensity of points to capture the clustering behavior of the point pattern.

$$K(r) = \frac{|W|}{n(n-1)} \sum_{i=1}^{n} \sum_{j=1, j \neq i}^{n} \mathbf{1}\{d_{ij} \leq r\} e_{ij}(r), e_{ij}(\bullet) \text{ is the edge correction weight}$$

L Function: L function is a variance stabled version of K function.

$$L(r) = \frac{\sqrt{K(r)}}{r}$$

We calculated G, F, and L functions in both univariate and multivariate fashions. For each of the functions, the distances between source cell and the target cells are considered. Univariate spatial functions sample source cells and target cells from the same type of cells while multivariate counterparts' sample from different types of cells. Univariate G, F, and L are calculated for the single-marker cell subsets, and multivariate G_cross and L_cross for different subsets such as CA9-LAMP2b. 'Gest'

7

299 function and 'Fest' function from 'spatstat' R package were used with Kaplan-Meier estimator(38), and
300 'Lest' function was used with isotropic correction(39,40).

## MODULE 4: Diagnostic risk estimation with pattern proportion

302 In the last module, we train a classifier using these niches to predict whether a patient will be "upstaged"
303 or "pure DCIS". This establishes the diagnostic power of these niches. A direct aggregation of niche
304 information within each sample/patient is not sufficient. Tumor microenvironment is heterogeneous, and
305 niches demonstrate diverse spatial and morphological behavior. To account for the diversity, we will
306 focus on how different niches are distributed across a sample. We show that the distributions of different
307 niches essentially characterize the tumor ecology in a much more refined manner compared with previous
308 distance-based definitions of hypoxia/oxidative layers.

309 One technical challenge is that the niche features computed in the previous module are high dimensional
310 and the niche features are diversely distributed. We propose to first find a simplified distributional
311 description of the niches, and then use the simplified description for prediction. First, we cluster the
312 niches into different sub-classes based on their features. The clustering is carried out using K-means
313 clustering with a tunable parameter k. Once the niche sub-classes are determined. We use their
314 distribution on each sample to predict its upstaged/pure DCIS status. The prediction power of the
315 classifier sheds light on the diagnostic power of the niches and their spatial and cellular features. Five-
316 fold cross-validation was employed, with one fold designated as the test set in each run. This approach
317 prevents\s data leakage and helps mitigate overfitting.

318 To understand the contribution of each feature to the prediction model, we employed SHAP (SHapley
319 Additive exPlanations) analysis. SHAP is a unified approach to interpreting machine learning models by
320 assigning each feature an importance value for a particular prediction. In our study, SHAP values were
321 computed for the features representing the proportions of different patterns within the niches. By
322 calculating the SHAP values, we could determine the impact of each feature on the model's output,
323 thereby identifying the most influential patterns that contribute to predicting DCIS upstaging. This step
324 is crucial for ensuring the transparency and interpretability of the machine learning model.

325 Furthermore, we select features that are highly relevant to the sub-classes using different approaches
326 including covariance, mutual information scoring and maximum relevance minimum redundancy
327 (mRMR)(41) and choose the features identified by both approaches. **Figure 4C** shows the gradient map
328 of each of these features on niches in the latent space.

329 **Niche distribution for diagnosis.** After assigning each duct to its sub-class, we aggregate across all
330 niches of each sample and use its sub-class distribution to characterize this sample. Assuming k niche
331 sub-classes, each sample has a k dimensional histogram to describe its niche sub-class distribution. We
332 call this the niche distributional (Nbd-Dist) feature. We trained a classifier to predict whether a sample
333 is pure DCIS or upstage. Repeating the iteration 10 times and comparing the mean area under curve
334 (AUC) on the test set. The classifier types experimented include lightGBM, soft vector machine (SVM),
335 logistic regression and random forest, and the random forest classifier yields the best performance.

336

### Data Availability

338 The data generated in this study are available within the article and its supplementary data files. All the
339 staining and annotations are also deposited in the physical sciences in oncology network (PSON).

8

## Results:

### Sample curation and cohort building

We built a retrospective cohort from 84 patients with histologically confirmed DCIS on core biopsy, followed by surgical excision, with available FFPE blocks at both Bx and Ex from Moffitt Cancer Center Biobank. The cohort has two arms: the first one is pure DCIS including the patient diagnosed with DCIS at both Bx and Ex. The second arm includes the upstaged group with DCIS at Bx and IDC at Ex (**Figure 1B**). HE stained slides of DCIS biopsy cores were retrieved from both the biobank core at the Moffitt Cancer Center tissue core and reviewed by a study pathologist (49). Then the selected blocks were pulled and sequentially cut for HE staining, CA9, and LAMP2b IHC staining. The HE and subsequent 2 IHC slides are digitally scanned using the Aperio XT® high-throughput slide scanner and housed on the web-based Aperio server/Spectrum database package. Upstage status was pulled from the electronic medical record and confirmed by our study pathologist from the Ex tissues (**Figure 1C**). All images were then segmented and annotated using Qupath supervised by study pathologist (25,49).

### Annotation and eco-evolutionarily mapping of habitats at the individual duct level

We have shown previously that peri-luminal cells that are far (>0.125 - 0.160 mm) from a blood supply inhabit a microenvironment of hypoxia and acidosis (18,20,26). Thus, we created two simple annotation zones on HE slides based on oxygen diffusion distance representing oxygen defined habitats: i) hypoxic zone or habitat that is above 125 μms from the duct boundary, basement membrane, and ii) normoxic habitat that is the outer regions adjacent to the basement membrane (**Figure 2A**). We used the basement membrane as our zero point of reference. We also annotated necrotic zones inside the hypoxic habitats that also represent the anoxic habitat falling perfectly above 0.160 mm distance from basement membrane. Since adjacent stroma is also of interest to our group and others, we annotated adjacent stroma for each duct with binary scoring of 1 for having adjacent stroma or 0 for lacking it (**Supplementary Table 1**). To ensure a balanced representation of hypoxic and normoxic habitats, we excluded small ducts by establishing a duct size threshold of minimum 400 μms in diameter (or 200 μms radius) for manual annotation (**Figure S2**). After annotating all the ducts bigger than 200 μms of radius on HE slides, we expanded our annotations to other 2 consecutive IHC slides stained with CA9 and LAMP2b antibodies (**Figure 2B**). Subsequently, our pathologist, Dr. Bai, manually scored each duct for hypoxic and normoxic habitats based on CA9 and LAMP2b positivity using a scoring scale of 0–3 (**Supplementary Table 1**). Following this, positive cells in IHC slides were counted using Qupath (25), habitats were categorized into different classes based on the count of positive cells. The distribution of these habitat categories was compared between pure DCIS and upstaged groups (**Figure 2C, 2D, and S3**). Using the Wilcoxon test, it was shown that there existed significant differences between pure DCIS and Upstaged group when habitats considered at the duct level. The tests were carried out for both hypoxic and oxidative layers for both CA9 (**Figure 2C and 2D**), and for LAMP2b (**Figure S3**) as well as architecture, grade, lymphocytes, microcalcifications, and necrosis (**Supplementary Table 1**). As shown in Figure 2D, CA9 scoring within hypoxic habitats provides a much clearer distinction between pure DCIS and upstaged groups compared to the normoxic zone. Interestingly, CA9 did not show significant differences between the groups when analyzed at the whole duct or whole-slide level, as is traditionally done (**Figure S2B**). However, focusing on hypoxic or oxidative habitats revealed that CA9-positive cells are distributed differently between the two patient groups. This analysis underscores the value of examining fine-scale

382  habitats within ducts. The improved performance of habitat-level scoring compared to whole-duct
383  scoring highlights the necessity and significance of exploring the cellular composition and interactions
384  within these microhabitats.
385

**Mapping Metabolic Niches Within Habitats to Enhance Spatial Machine Learning Models**

387  Previous analyses of hypoxic and normoxic habitats in breast cancer ducts were limited to scoring each
388  biomarker individually, focusing solely on the count of positive cells within each habitat. To broaden the
389  scope and incorporate interactions and relationships between these two eco-evolutionary marker-positive
390  cells, a co-registration step was essential. This step enabled the creation of a virtual multiplex IHC
391  (mIHC) by mapping cells onto a unified 2D reference space. HE slides were selected as the reference,
392  and all IHC slides were registered onto this common framework. (**Figure S4**). Note that since our analysis
393  is carried out duct-by-duct, it is not necessary to register the whole slide. Instead, for each duct, we
394  register its IHC stainings to its HE staining. This ensures all the downstream analyses could be performed
395  on the same HE slide coordinates system, providing consistency and precision in the spatial data
396  integration. Then we used these mIHC images to define niches of cells that are positive for CA9,
397  LAMP2b, or both. We hypothesized that niches within habitats characterized by both markers together
398  would provide greater biological insight than analyzing each marker individually, given the established
399  correlation between hypoxia and acid phenotypes. Then, we focus on the cell features such as nuclear
400  morphology and texture and cell spatial features inside these niches to explore their predictive power on
401  DCIS upstaging. As illustrated in **Figure 3**, we first map each positive cells to the reference HE slide
402  using the co-registration described above. Then, by treating each positive cell as a node and connecting
403  the cells within a distance threshold, we construct a cell-proximity graph out of mIHC positive cells
404  whereby each connected component of this graph represents a continuous region or niche that is hypoxic,
405  acidic, or both. The threshold is a tunable parameter that is optimized by the classifying power of the
406  downstream analysis. And depending on the selection of the eco-evo markers, there can be CA9 positive
407  niches, LAMP2b positive niches, or both CA9 and LAMP2b positive niches. We then developed a pattern
408  differential analysis pipeline, which comprises two stages: First, the samples are clustered based on the
409  features and classified into one of the clusters or patterns. Then for each patient, we calculate the
410  proportion of each pattern, forming a distribution profile of the patterns.
411  By using these proportion features, we train a classifier aiming to predict the upstaging status. From this
412  pipeline, we are able to predict the clinical outcome of a patient based on his/her spatially defined pattern
413  distributions (**Figure 1C**). Then, to test the hypothesis that finer regions with biological meanings could
414  provide better predictive power, we conduct a multi scale analysis performing a series of experiments
415  using the same set of features and with the same pattern differential analysis pipeline at 3 different scales:
416  duct, habitat, and niche (**Figure 1C**). At the habitat level, normoxic and hypoxic zones are analyzed
417  independently. At the niche level, analyses are further refined to separately examine CA9-positive cells,
418  LAMP2b-positive cells, and cells co-positive for both CA9 and LAMP2b.
419  For all the experiments, the biopsy dataset underwent 5-fold stratified cross-validation, where in each
420  round, 4 folds served as the training dataset and 1-fold as the test dataset, with the goal of predicting the
421  patients' clinical outcome at the time of biopsy. Upon comparing the mean accuracy score and the mean
422  AUC score of all the classifiers, the niche level classifier yielded the best predictive results particularly
423  under both metrics (**Table 1**). This result confirms that niche-based analysis outperforms our primary
424  habitat analysis. The higher accuracy of the niche measurements may be implying the phenotype-based

425 niche measurement is better than inferring habitat from oxygen diffusion rate measure based on the
426 distance of the cells from basement membrane. Also, it is worth mentioning that oxygen habitat analysis
427 is a rough estimate in our analysis since we do not know the exact location of the vasculature and their
428 activity.

| | Duct | Habitat | | Niche | | |
|---|---|---|---|---|---|---|
| | | Normoxia | Hypoxia | CA9 | LAMP2b | CA9 & LAMP2b |
| **Accuracy** | 0.78 ± 0.06 | 0.86 ± 0.03 | 0.83 ± 0.06 | 0.82 ± 0.06 | 0.90 ± 0.03 | 0.90 ± 0.03 |
| **AUC** | 0.61 ± 0.08 | 0.67 ± 0.03 | 0.66 ± 0.10 | 0.64 ± 0.10 | 0.72 ± 0.07 | 0.74 ± 0.13 |

429 **Table 1**. **Performance scores of multi scale classifiers.** While habitat-level analysis enhanced
430 performance, the niche-level classifier produced the most accurate predictive results.
431

432 **Post analysis to reveal contributing features and prototype visualization on mIHC.**
433 After identifying the best-performing classifier based on the AUC metric we employed SHAP (48)
434 (Shapley Additive exPlanations) analysis to interpret the model by calculating SHAP values for each
435 feature, specifically on the proportions of distinct patterns (**Figure 4b**). The pattern with the maximum
436 SHAP value, identified as the most impactful, underwent further differential analysis to uncover features
437 that significantly differentiated this pattern from others. This differential analysis employed methods
438 including correlation analysis, mutual information (MI), and maximum relevance minimum redundancy
439 (MRMR), which together identified Area_min, Perimeter_min, AreaBbox_min, and $F\_0 <= r < 10$ as the
440 top distinguishing features for Pattern 5 (**Figure 4c**). A prototype for Pattern 5, selected based on its
441 alignment with the mean values of these features, was visualized to illustrate its characteristics (**Figure
442 4d**). Using a multi-scale analytical approach, we integrated spatial interactions of CA9-positive and
443 LAMP2b-positive cells into the machine learning pipeline to distinguish between pure DCIS and
444 progressed DCIS. Niche-level analysis yielded the highest accuracy and AUC, emphasizing the
445 importance of fine-scale regions in predicting clinical outcomes. The use of SHAP analysis and
446 differential analysis provided an interpretable framework to highlight influential patterns and features,
447 such as Area_min and Perimeter_min, offering insights into the tumor microenvironment. This approach
448 not only advanced our understanding of key spatial and morphological features but also demonstrated
449 significant potential for precise diagnostic tools in clinical applications.
450

451

## Discussion:
453 Ductal carcinoma in situ is the most prevalent type of precancer that can range from indolent to
454 aggressive. DCIS lesions are highly heterogeneous in their intra- and inter- ductal physical
455 microenvironments, genetics, and molecular expression patterns. They can be described as complete
456 ecosystems containing habitats and niches including normal epithelial cells, pre-cancer cells, stromal
457 cells, vasculature, structural proteins, signaling proteins and physical factors such as pH and oxygen
458 concentration (18). These habitats and niches of micro-domains can contain unique mixtures of cells with

459    physical and biochemical characteristics, with differential evolutionary potential and trajectories (27).
460    The niches with similar mixtures of cells usually are also similar in their physiology and phenotypes
461    mainly due to living in similar habitats. Our hypothesis is that knowledge of these niches and their
462    habitats can potentially provide patient benefit by stratifying their tumor progress and therapeutic choices.
463    However, tools and techniques are lacking to distinguish them. Proper tools and techniques can identify
464    and define habitats and niches to map (pre-)cancer ecosystems to discriminate between the different types
465    of DCIS to design the right treatment for breast cancer patients.

466    In this study, we argue that the overdiagnosis and overtreatment of DCIS stem from conventional
467    frameworks that focus primarily on genetic signatures while neglecting the phenotypic heterogeneity
468    within tumor ecosystems. Thus, we interpreted complex eco-evolutionary data of cancer cells within their
469    niche using machine learning and pathomics, all framed within an innovative ecological and evolutionary
470    dynamic model. Oxygen habitats are identified based on varying levels of perfusion and oxygenation,
471    which are believed to play a crucial role in driving ecological diversity by changing cancer cells
472    metabolism, creating new habitats, and enhancing tumor heterogeneity, ultimately leading to diverse
473    evolutionary trajectories. (28, 29). Solid tumors often exhibit an impaired vascular system, leading to
474    habitats within tumors that vary in hypoxia, nutrient deficiency, and acidity. These habitats can
475    significantly influence the spatial selection of cellular phenotypes in distinct subregions. Inhabiting
476    hypoxia, acidosis, and severe nutrient deprived habitats, face (pre-)cancer cells to strong selective
477    pressures leading to divergence to novel phenotypes in population. These new phenotypes can
478    reciprocally influence the microenvironment reshaping due to their new metabolic phenotypes resulting
479    in a dynamically changing tumor ecosystem with multiple habitats. Therefore, the phenotype of the cells
480    residing in these habitats can also be leveraged to define the habitats with a certain degree of accuracy.
481    Previous research from our group and others demonstrated that cancer cells within breast ducts, exposed
482    to chronic hypoxia and acidosis, develop adaptive mechanisms for survival in this challenging
483    microenvironment including expression of CA9 or LAMP2b at the cell surface (18,20,30). However,
484    none of these findings were used in a relevant translational study for biomarker discovery. In this study,
485    we explore these biomarkers within an eco-evolutionary framework for the first time, using them as
486    indicators of the metabolic state of cancer cells residing in a niche as part of oxygen habitats that may
487    favor the selection of more aggressive phenotypes to predict the upstaging of DCIS. While a longitudinal
488    study would indeed be a better study design for direct observation of evolutionary changes over time, our
489    current cross-sectional approach enables us to capture a snapshot of the tumor microenvironment at two
490    near time points, providing valuable insight into the conditions that distinguish DCIS from IDC. We
491    recognize the assumption that the synchronous IDC microenvironment may contribute to the progression
492    from DCIS to IDC. However, our study design allows us to test whether specific microenvironmental
493    factors and related habitats and niche correlate with the presence of IDC, which can provide strong
494    hypotheses for future longitudinal investigations. A future prospective or retrospective longitudinal
495    (multiple long time points) study would indeed help distinguish whether these microenvironmental
496    changes in tumor ecosystem locally belonged to habitats or niches can drive progression from DCIS to
497    IDC or if IDC-induced those changes in the tumor ecosystem contribute to the synchronous DCIS
498    phenotype.

499    In our curated retrospective cohort of 84 DCIS patients with histologically confirmed DCIS on core
500    biopsy, we manually annotated 916 single ducts and more than 3000 habitats on all three slides and scored
501    them at habitat levels. This unique detailed eco-evolutionary annotation can be used for future similar

502  eco-evolutionary designed studies including stroma habitats. Our risk scoring system integrating
503  principles of ecological-evolutionary dynamics with pathological imaging and molecular features of
504  early-stage breast tumors showed improvement on prediction power of biomarkers alone and in
505  combination.

506  We employed a 5-fold stratified cross-validation approach to ensure robust internal validation of our
507  model. While this method helps mitigate overfitting and provides reliable performance estimates, we
508  acknowledge the absence of an independent validation set, which is crucial for assessing the model's
509  generalizability. The unique design of our cohort, which integrates specific ecological and
510  microenvironmental factors, limits the availability of comparable external datasets for validation. As
511  such, there is no current dataset with similar characteristics for cross-validation. We recognize this as a
512  key limitation and emphasize that future studies should aim to validate the model on independent cohorts
513  when such datasets become available. Furthermore, although our model achieved an AUC of 0.74, this
514  performance is not yet sufficient for clinical translation. Additional efforts to refine the model and test it
515  in larger, independent cohorts will be essential before its use in clinical practice can be considered.
516  Interestingly, a recent approach using multiplex IF on DCIS cohort reached the same AUC(2). While
517  both our study and the Risom et al. paper aim to leverage spatial relationships to predict DCIS
518  progression, we would like to emphasize that the two approaches are fundamentally different in terms of
519  the markers used. Risom et al. focused on a broad panel of markers, including those related to the stroma,
520  immune cells, and tumor cells, which provide a comprehensive view of the tumor microenvironment. In
521  contrast, our approach centers on eco-evolutionary markers derived from adaptation of cancer cells to
522  physical microenvironment, specifically CA9 and LAMP2b, which are associated with hypoxia and
523  tumor acidity and their spatial distribution, respectively. These differences reflect divergent hypotheses
524  about the key drivers of DCIS progression. The fact that both studies report a similar AUC of 0.74, with
525  the distinct marker sets and biological processes, suggests that our findings offer complementary insights
526  into DCIS progression and combination of approaches might increase the accuracy.

527  Our study demonstrates the utility of eco-evolutionary principles in understanding DCIS progression. In
528  our study, we proposed that specific tumor microenvironmental conditions, such as hypoxia and acidosis,
529  are associated with phenotypic changes that may indicate DCIS progression. However, although we have
530  shown previously that these microenvironments can cause aggressive phenotypes, we acknowledge that
531  our findings here do not conclusively demonstrate that these environmental factors are causative agents
532  in the transition from DCIS to IDC. Instead, our data suggest that these conditions could serve as
533  biomarkers for identifying lesions that are more likely to be upstaged. However, the ability to define
534  more refined cell phenotypes within each region of interest (ROI) could further enhance our analysis. If
535  we can identify and characterize more detailed phenotypes, it would allow us to extract additional features
536  that describe the spatial interactions of these phenotypes. This, in turn, could potentially improve the
537  classifier's performance and make the results more interpretable. By capturing the intricate interactions
538  between various cell types and their microenvironments, we could gain deeper insights into the ecological
539  dynamics driving DCIS progression and improve predictive models for patient outcomes.

540  In recent years, there has been a growing trend towards adopting a "watchful waiting" approach for certain
541  cases of DCIS, rather than immediate surgical excision(31,32). This strategy aims to reduce overtreatment
542  by closely monitoring DCIS lesions that may not progress to invasive cancer. In this context, our upstaging
543  predictions become particularly relevant. Identifying microenvironmental and phenotypic factors that
544  indicate a higher likelihood of progression to IDC could help clinicians make more informed decisions

13

545 about when to intervene and when to adopt a more conservative, observational approach. The ability to
546 predict which DCIS cases are at higher risk of progressing to invasive disease would provide critical
547 information for optimizing patient management, minimizing unnecessary treatments, and reducing the
548 psychological and physical burdens associated with overtreatment(33). Further validation of these
549 predictive models could therefore have important implications for guiding treatment strategies in the
550 context of DCIS.

**Lead contact**

552 Further information and any related requests should be directed to and will be fulfilled by the lead contact
553 Mehdi Damaghi (Mehdi.Damaghi@stonybrookmedicine.edu).
554

563

**Author contributions**

565 M.D. conceptualized and designed the research; Y.X., M.A., J.D.B., A.C., Y.C., M.D., performed the
566 experiment and analysis; J.D.B. reviewed all the slides and scored them as the project pathologist; P.P.,
567 C.C., and M.D. contributed to results interpretation; and M.D. wrote the paper. All authors revised the
568 paper.

569

**Figure Captions**

571

572 **Figure 1. Ecological and evolutionary designed biomarkers of DCIS upstaging. A)** Model of
573 microenvironment-driven evolution of breast cancer from normal breast tissue to DCIS and IDC: Our schematic
574 is overlaid on HE staining of breast cancer specimens at different stages of DCIS and IDC. Different patients may
575 experience various types of evolutionary trajectory following different evolutionary models, including linear and
576 branched progression from DCIS to IDC shown here. Note that these events are not sequential or stepwise. **B)** The
577 patient cohort was curated from retrospective DCIS samples, with two sample collections at biopsy and excision.
578 The main criterion was the diagnosis of DCIS at the biopsy stage. **C)** Eco-evolutionary designed- machine learning
579 assisted pipeline to define cancer cell niches inside oxygen habitats in DCIS. *i)* Data preprocessing steps including
580 duct annotation, cell detection and classification for HE and IHC slides, followed by co-registration to map IHC-
581 identified cells onto the HE slides. *ii)*The analysis is carried out at multiple scales, namely duct, habitat and niche,
582 from the largest to smallest. At each scale the nucleus morphology texture feature and spatial features are extracted.
583 *iii)* The pattern differential analysis approach where the patterns are firstly identified and then the proportions of
584 such patterns are used as features to predict the upstaging status of a patient.
585

14

**Figure 2**. **Eco-evolutionarily designed biomarker discovery to predict upstaging in DCIS. A)** Illustration of normoxic, hypoxic and necrotic habitats in a duct. **B)** Illustration of annotation and scoring on 2 IHCs and how cells are scored in each habitat. **C)** and **D)** Dot plots of counts of CA9 expression in each habitat per duct. Cells are scored 0 for 'negative' or '1+','2+','3+' for positive cells based on their intensity. Scoring was performed and analyzed separately for normoxic (oxidative) habitat (C) or hypoxic habitat (D). In the dot plot, each dot is a single duct. The color of dots reflects their score as follows: Blue = 0, yellow ='1+', orange ='2+', and red = '3+'. The number of dots reflects how many ducts were detected in each patient's biopsy with size bigger than 400 μms in diameter. The distribution in hypoxic habitat is significantly different between pure DCIS and upstaged groups in hypoxic habitats and not in oxygenated habitat. Data was analyzed using the Wilcoxon signed-rank test. The same graph is created for LAMP2b (supplementary fig. 2).

**Figure 3. Niches are defined inside habitats from the hypoxia and acidosis markers expression. A)** One sample duct from CA9 slide. Top: The original IHC slide. Middle: Cell detection and intensity-based classification using Qupath overlaid on the slide. Bottom: the graph constructed from the CA9 positive cells and the connected components of the graph (Niches) highlighted in different colors. **B)** The HE staining of the same duct as A. Top: The original HE slide. Middle: Duct annotation overlaid on the HE slide. Bottom: Co-registered CA9-positive niches mapped and overlaid on HE slides as mIHC to be able to extract HE features from CA9 positive niches. Note the orientation of HE and CA9 slide was opposite, and our co-registration technique successfully created a mIHC of the ducts with similar coordinates. The same approach was used for LAMP2b and the combination.

**Figure 4. Post Analysis reveals the top contributing patterns and features. A)** UMAP of the features of the niches, different colors represent different clusters(patterns) **B)** Top: The impact of each pattern on the classifying result, blue and red colors represent impact on pure DCIS and progressed predictions respectively, the proportion of pattern 5 has the greatest impact for both categories. Bottom: Using correlation, MI, and MRMR to obtain the most contributing features in the pattern 5 clustering phase, identifying a common feature set that includes 4 features: Area_min, Perimeter_min, AreaBbox_min, and F_0<=r<10. **C)** UMAP showing the value of the 4 identified features for different samples, and it can be seen that samples in the pattern 5 tend to have higher values in Area_min, Perimeter_min, AreaBbox_min and low values for F_0<=r<10. **D)** A niche belonging to pattern 5, it contains no small size cells and exhibits a relatively dispersed distribution.

**References:**

1.  Greaves M, Maley CC. Clonal evolution in cancer. Nature. 2012;481:306–13.

2.  Risom T, Glass DR, Averbukh I, Liu CC, Baranski A, Kagel A, et al. Transition to invasive breast cancer is associated with progressive changes in the structure and composition of tumor stroma. Cell. 2022;185:299–310.e18.

3.  Maley CC, Aktipis A, Graham TA, Sottoriva A, Boddy AM, Janiszewska M, et al. Classifying the evolutionary and ecological features of neoplasms. Nat Rev Cancer. 2017;17:605–19.

627    4.    Boutry J, Tissot S, Ujvari B, Capp J-P, Giraudeau M, Nedelcu AM, et al. The evolution
628        and ecology of benign tumors. Biochim Biophys Acta Rev Cancer. 2022;1877:188643.

629    5.    Amend SR, Pienta K. Abstract 2884: Tumor-driven eutrophication of the tumor ecosystem
630        selects for cancer cell clones that overcome evolutionary inertia leading to increased
631        metastatic capacity. Cancer Res. American Association for Cancer Research;
632        2015;75:2884–2884.

633    6.    Damaghi M, Mori H, Byrne S, Xu L, Chen T, Johnson J, et al. Collagen production and
634        niche engineering: A novel strategy for cancer cells to survive acidosis in DCIS and
635        evolve. Evol Appl. 2020;13:2689–703.

636    7.    Lipinski KA, Barber LJ, Davies MN, Ashenden M, Sottoriva A, Gerlinger M. Cancer
637        Evolution and the Limits of Predictability in Precision Cancer Medicine. Trends Cancer
638        Res. 2016;2:49–63.

639    8.    Giaquinto AN, Sung H, Miller KD, Kramer JL, Newman LA, Minihan A, et al. Breast
640        Cancer Statistics, 2022. CA Cancer J Clin. 2022;72:524–41.

641    9.    Strand SH, Rivero-Gutiérrez B, Houlahan KE, Seoane JA, King LM, Risom T, et al.
642        Molecular classification and biomarkers of clinical outcome in breast ductal carcinoma in
643        situ: Analysis of TBCRC 038 and RAHBT cohorts. Cancer Cell. 2022;40:1521–36.e7.

644    10.    Lehman CD, Arao RF, Sprague BL, Lee JM, Buist DSM, Kerlikowske K, et al. National
645        Performance Benchmarks for Modern Screening Digital Mammography: Update from the
646        Breast Cancer Surveillance Consortium. Radiology. 2017;283:49–58.

647    11.    Lips EH, Kumar T, Megalios A, Visser LL, Sheinman M, Fortunato A, et al. Genomic
648        analysis defines clonal relationships of ductal carcinoma in situ and recurrent invasive
649        breast cancer. Nat Genet. 2022;54:850–60.

650    12.    Sarhadi S, Salehzadeh-Yazdi A, Damaghi M, Zarghami N, Wolkenhauer O, Hosseini H.
651        Omics Integration Analyses Reveal the Early Evolution of Malignancy in Breast Cancer.
652        Cancers [Internet]. 2020;12. Available from: http://dx.doi.org/10.3390/cancers12061460

653    13.    Heselmeyer-Haddad K, Berroa Garcia LY, Bradley A, Ortiz-Melendez C, Lee W-J,
654        Christensen R, et al. Single-cell genetic analysis of ductal carcinoma in situ and invasive
655        breast cancer reveals enormous tumor heterogeneity yet conserved genomic imbalances
656        and gain of MYC during progression. Am J Pathol. 2012;181:1807–22.

657    14.    Hanna WM, Parra-Herran C, Lu F-I, Slodkowska E, Rakovitch E, Nofech-Mozes S. Ductal
658        carcinoma in situ of the breast: an update for the pathologist in the era of individualized
659        risk assessment and tailored therapies. Mod Pathol. 2019;32:896–915.

660    15.    Carmeliet P, Jain RK. Principles and mechanisms of vessel normalization for cancer and
661        other angiogenic diseases. Nat Rev Drug Discov. 2011;10:417–27.

662    16.    Wu B, Liu D-A, Guan L, Myint PK, Chin L, Dang H, et al. Stiff matrix induces exosome
663        secretion to promote tumour growth. Nat Cell Biol. 2023;25:415–24.

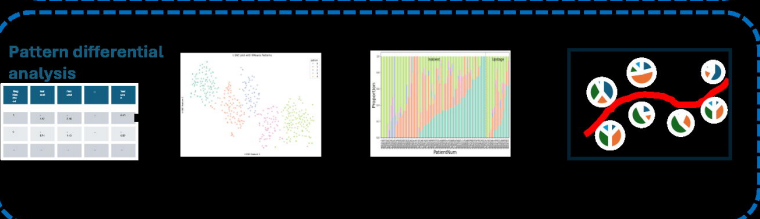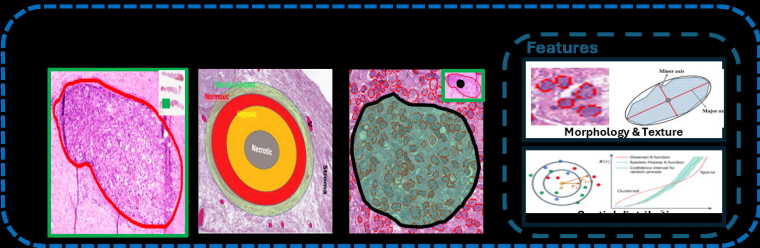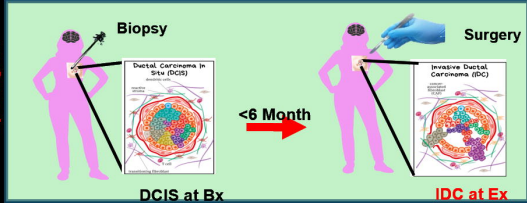664    17.    Persi E, Duran-Frigola M, Damaghi M, Roush WR, Aloy P, Cleveland JL, et al. Systems
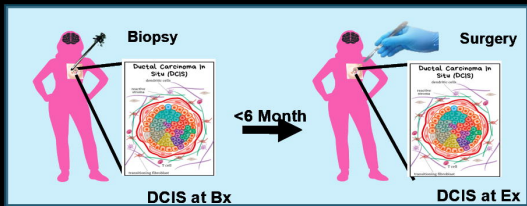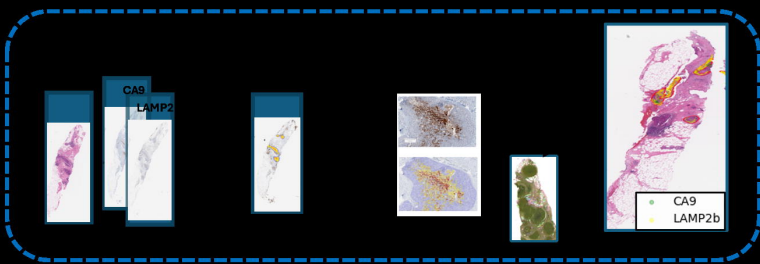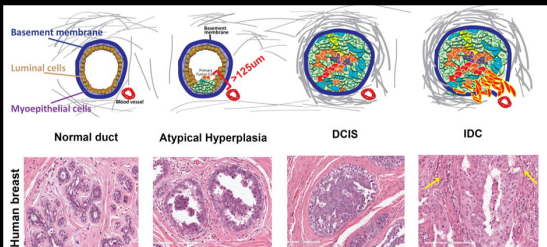
16

665       analysis of intracellular pH vulnerabilities for cancer therapy. Nat Commun. 2018;9:2997.

666  18.  Damaghi M, West J, Robertson-Tessi M, Xu L, Ferrall-Fairbanks MC, Stewart PA, et al.
667       The harsh microenvironment in early breast cancer selects for a Warburg phenotype. Proc
668       Natl Acad Sci U S A [Internet]. 2021;118. Available from:
669       http://dx.doi.org/10.1073/pnas.2011342118

670  19.  Lobo RC, Hubbard NE, Damonte P, Mori H, Pénzváltó Z, Pham C, et al. Glucose Uptake
671       and Intracellular pH in a Mouse Model of Ductal Carcinoma In situ (DCIS) Suggests
672       Metabolic Heterogeneity. Front Cell Dev Biol. 2016;4:93.

673  20.  Damaghi M, Tafreshi NK, Lloyd MC, Sprung R, Estrella V, Wojtkowiak JW, et al.
674       Chronic acidosis in the tumour microenvironment selects for overexpression of LAMP2 in
675       the plasma membrane. Nat Commun. 2015;6:8752.

676  21.  Ordway B, Swietach P, Gillies RJ, Damaghi M. Causes and Consequences of Variable
677       Tumor Cell Metabolism on Heritable Modifications and Tumor Evolution. Front Oncol.
678       2020;10:373.

679  22.  Gillies RJ, Verduzco D, Gatenby RA. Evolutionary dynamics of carcinogenesis and why
680       targeted therapy does not work. Nat Rev Cancer. 2012;12:487–93.

681  23.  Ibrahim-Hashim A, Robertson-Tessi M, Enriquez-Navas PM, Damaghi M, Balagurunathan
682       Y, Wojtkowiak JW, et al. Defining Cancer Subpopulations by Adaptive Strategies Rather
683       Than Molecular Properties Provides Novel Insights into Intratumoral Evolution. Cancer
684       Res. 2017;77:2242–54.

685  24.  Damaghi M, Gillies R. Phenotypic changes of acid-adapted cancer cells push them toward
686       aggressiveness in their evolution in the tumor microenvironment. Cell Cycle.
687       2017;16:1739–43.

688  25.  Bankhead P, Loughrey MB, Fernández JA, Dombrowski Y, McArt DG, Dunne PD, et al.
689       QuPath: Open source software for digital pathology image analysis. Sci Rep.
690       2017;7:16878.

691  26.  Freischel AR, Damaghi M, Cunningham JJ, Ibrahim-Hashim A, Gillies RJ, Gatenby RA, et
692       al. Frequency-dependent interactions determine outcome of competition between two
693       breast cancer cell lines. Sci Rep. 2021;11:4908.

694  27.  Jardim-Perassi BV, Huang S, Dominguez-Viqueira W, Poleszczuk J, Budzevich MM,
695       Abdalah MA, et al. Multiparametric MRI and Coregistered Histology Identify Tumor
696       Habitats in Breast Cancer Mouse Models. Cancer Res. 2019;79:3952–64.

697  28.  Sobhani F, Muralidhar S, Hamidinekoo A, Hall AH, King LM, Marks JR, et al. Spatial
698       interplay of tissue hypoxia and T-cell regulation in ductal carcinoma in situ. NPJ Breast
699       Cancer. 2022;8:105.

700  29.  Compton ZT, Mallo D, Maley CC. Stronger together: Cancer clones cooperate to alleviate
701       growth barriers in critical cancer progression transitions. Cancer Res. American
702       Association for Cancer Research (AACR); 2023;83:4013–4.

703 30. Chafe SC, McDonald PC, Saberi S, Nemirovsky O, Venkateswaran G, Burugu S, et al.
704     Targeting hypoxia-induced carbonic anhydrase IX enhances immune-checkpoint blockade
705     locally and systemically. Cancer Immunol Res. American Association for Cancer
706     Research; 2019;7:1064–78.

707 31. Ryser MD, Worni M, Turner EL, Marks JR, Durrett R, Hwang ES. Outcomes of active
708     surveillance for ductal carcinoma in situ: A computational risk analysis. J Natl Cancer Inst.
709     Oxford University Press (OUP); 2016;108:djv372.

710 32. Glencer AC, Miller PN, Greenwood H, Maldonado Rodas CK, Freimanis R, Basu A, et al.
711     Identifying good candidates for active surveillance of ductal carcinoma in situ: Insights
712     from a large neoadjuvant endocrine therapy cohort. Cancer Res Commun. American
713     Association for Cancer Research (AACR); 2022;2:1579–89.

714 33. Fortunato A, Mallo D, Cisneros L, King LM, Khan A, Curtis C, et al. Evolutionary
715     Measures Show that Recurrence of DCIS is Distinct from Progression to Breast Cancer.
716     medRxiv [Internet]. 2024; Available from: http://dx.doi.org/10.1101/2024.08.15.24311949

717 34. Gatenbee CD, Baker A-M, Prabhakaran S, Swinyard O, Slebos RJC, Mandal G, et al.
718     Virtual alignment of pathology image series for multi-gigapixel whole slide images. Nat
719     Commun. 2023;14:4502.

720 35. Schmidt U, Weigert M, Broaddus C, Myers G. Cell Detection with Star-Convex Polygons.
721     Medical Image Computing and Computer Assisted Intervention – MICCAI 2018. Springer
722     International Publishing; 2018. page 265–73.

723 36. Zhao S, Chen D-P, Fu T, Yang J-C, Ma D, Zhu X-Z, et al. Single-cell morphological and
724     topological atlas reveals the ecosystem diversity of human breast cancer. Nat Commun.
725     2023;14:6796.

726 37. Baddeley A, Rubak E, Turner R. Spatial Point Patterns: Methodology and Applications
727     with R. CRC Press; 2015.

728 38. Baddeley A, Gill R. Kaplan-Meier estimators of interpoint distance distributions for spatial
729     point processes. IEEE Trans Inf Theory [Internet]. 1993; Available from:
730     https://ir.cwi.nl/pub/5247/05247D.pdf

731 39. Ohser J. On estimators for the reduced second moment measure of point processes. Series
732     Statistics. Taylor & Francis; 1983;14:63–71.

733 40. Kendall WS. Stochastic Geometry: Likelihood and Computation. Routledge; 2019.

734 41. Peng H, Long F, Ding C. Feature selection based on mutual information: criteria of max-
735     dependency, max-relevance, and min-redundancy. IEEE Trans Pattern Anal Mach Intell.
736     2005;27:1226–38.

737 48. Lundberg, S. M. & Lee, S.-I. A unified approach to interpreting model predictions. *Proc.*
738 *Adv. Neural Inf. Process. Syst.* 30 (2017).

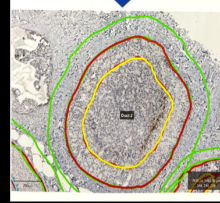739 49. Mayfield, J. D., Ataya, D., Abdalah, M., Stringfield, O., Bui, M. M., Raghunand, N., Niell,
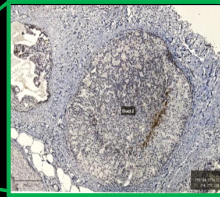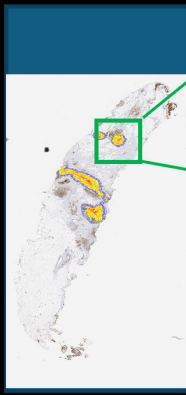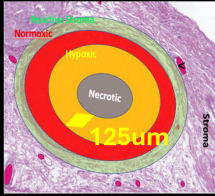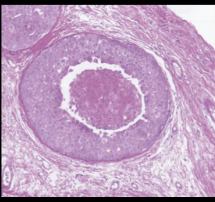
740    B., & El Naqa, I. Presurgical upgrade prediction of DCIS to invasive ductal carcinoma using
741    time-dependent deep learning models with DCE MRI. Radiology: Artificial Intelligence 6,
742    e230057 (2024).

743

Normal duct

Atypical Hyperplasia

DCIS

IDC

Basement membrane

Luminal cells

Myoepithelial cells

Blood vessel

>125um

Human breast

Biopsy

Surgery

Ductal Carcinoma In Situ (DCIS)

<6 Month

DCIS at Bx

DCIS at Ex

Upstaged

Biopsy

Surgery

Ductal Carcinoma In Situ (DCIS)

Invasive Ductal Carcinoma (IDC)

<6 Month

DCIS at Bx

IDC at Ex

CA9

LAMP2

CA9

LAMP2b

Features

Morphology & Texture

Pattern differential analysis

**c)**

## Oxidative

Pure DCIS | Progressed

Number of ducts

Wilcoxon test between indolent and upstaged
p-value: 0.06452

Patient

**d)**

## Hypoxia

Pure DCIS | Progressed

Number of ducts

Wilcoxon test between indolent and upstaged
p-value: 0.002253

Patient

CA9