

1 **The supplementary file of “Decoding Tumor**  
2 **Phenotype: A Radiologist-Inspired Deep Learning**  
3 **System for Breast Cancer Recurrence Prediction.”**

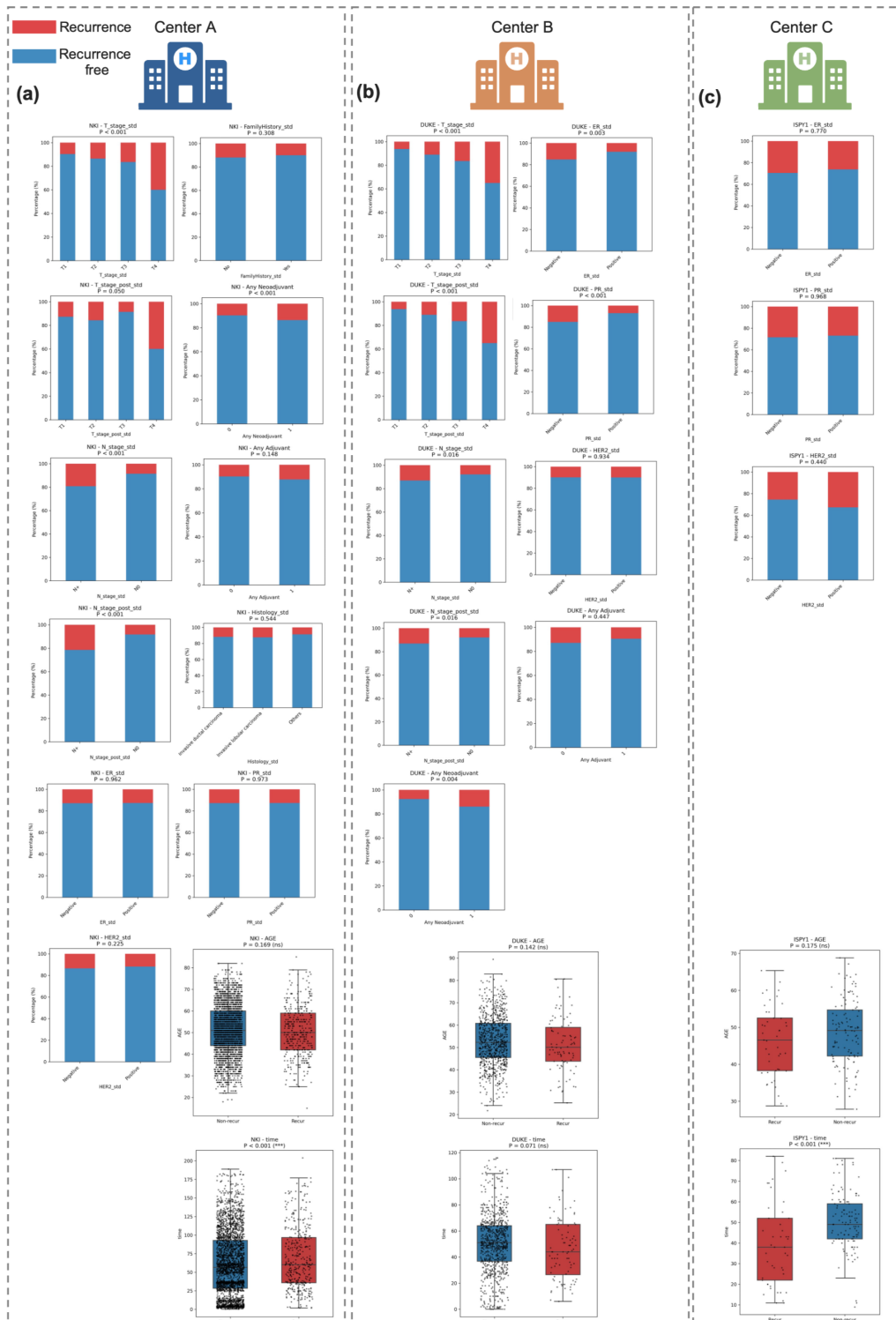
4

## 5 **Distribution of clinicopathological characteristics across** 6 **cohorts**

7 Supplementary Fig. S1 presents the distribution of clinicopathological characteristics  
8 across the multi-center cohorts used in this study, including the NKI development cohort,  
9 the Duke external validation cohort, and the I-SPY1 external validation cohort. Stacked bar  
10 charts illustrate the proportions of recurrence and recurrence-free cases across categorical  
11 clinical variables, while box plots summarize the distributions of continuous variables such  
12 as patient age and time-to-event. Recurrence cases are indicated in red and recurrence-  
13 free cases in blue.

14 Statistical comparisons between recurrence and recurrence-free groups were performed  
15 using the Chi-square test for categorical variables and the Mann–Whitney U test for  
16 continuous variables. Several established clinical risk factors, including tumor stage (T  
17 stage) and nodal status (N stage), showed statistically significant associations with  
18 recurrence in univariable analyses in certain cohorts. However, a notable proportion of  
19 patients within clinically high-risk categories (e.g., higher T stage or node-positive disease)  
20 remained recurrence-free, suggesting that conventional clinicopathological variables alone  
21 provide limited discrimination for recurrence risk at the individual patient level.

22 In addition, the distributions of both recurrence prevalence and clinical characteristics  
23 varied across the three cohorts. For example, the I-SPY1 cohort exhibited a relatively  
24 higher proportion of recurrence events compared with the other cohorts. These differences  
25 highlight the presence of inter-cohort heterogeneity and potential distributional shifts  
26 between institutions, underscoring the importance of evaluating model robustness and  
27 generalizability across independent datasets.



28

29 **Fig. S1 | Clinicopathological characteristics and outcome distributions across multi-**  
 30 **center cohorts.** Stacked bar charts and box plots displaying the distribution of key clinical  
 31 variables stratified by recurrence status (Red: Recurrence; Blue: Recurrence-free) for  
 32 the NKI development cohort (a), Duke external validation cohort (b), and I-SPY1 external  
 33 validation cohort (c). P-values were calculated using the Chi-square test for categorical

34 variables and the Mann-Whitney U test for continuous variables. While univariable  
35 analyses confirm that standard risk factors (e.g., T-stage, N-stage) are statistically  
36 associated with recurrence ( $P < 0.05$ ), the substantial proportion of recurrence-free  
37 survivors within high-risk groups (e.g., T3/N+) illustrates the limited distinctiveness of  
38 clinical variables alone. Furthermore, the varying prevalence of recurrence across centers  
39 (e.g., higher event rate in I-SPY1) highlights the distributional shift challenges addressed  
40 by the AI model.

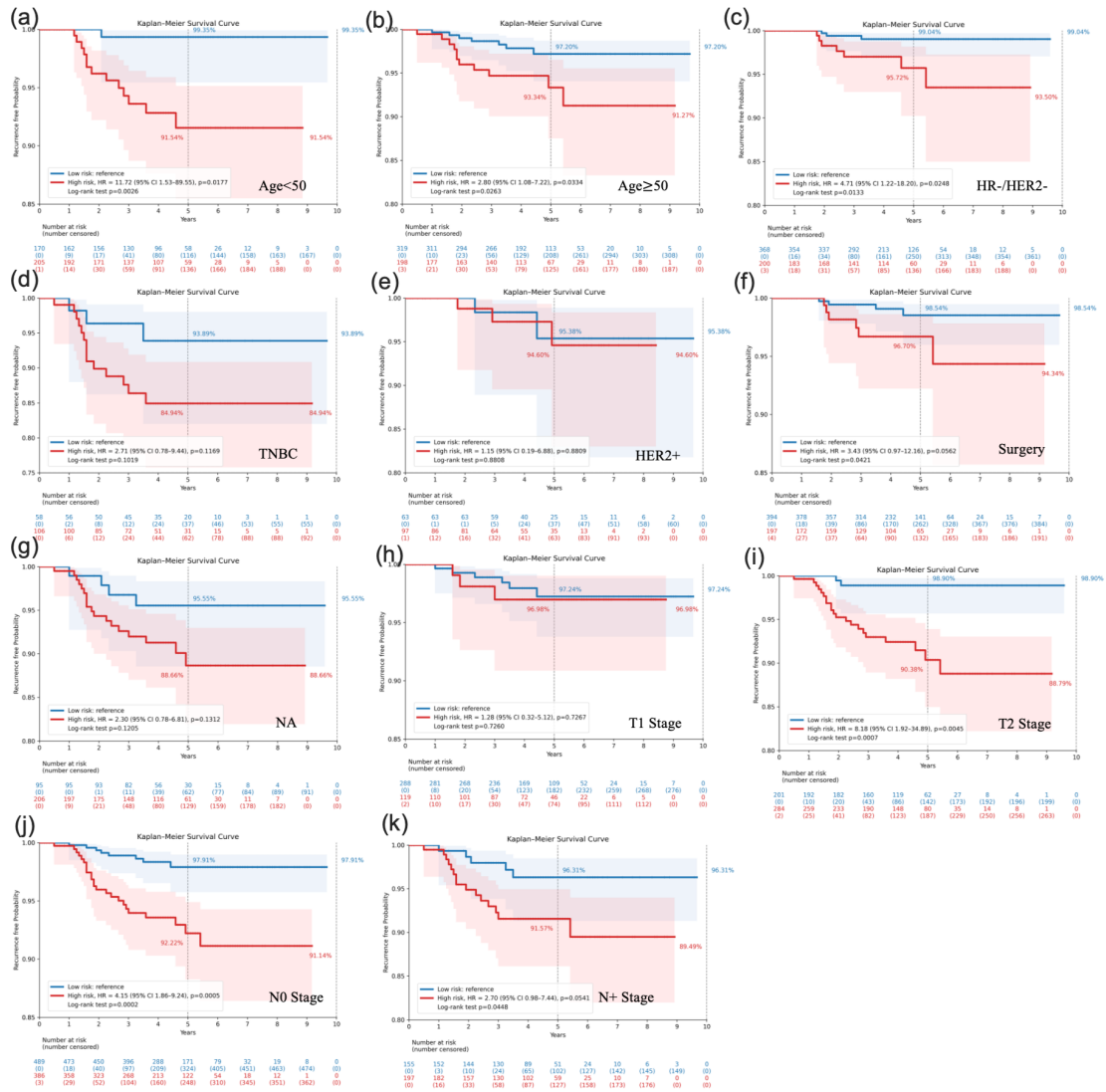
41

42 **Subgroup survival analyses performed in the Duke external**  
43 **validation cohort**

44 Supplementary Fig. S2 presents subgroup survival analyses performed in the Duke  
45 external validation cohort. Patients were stratified into model-predicted low-risk and high-  
46 risk groups based on the median risk score derived from the training cohort. Kaplan–Meier  
47 survival curves were then generated within clinically relevant subgroups, including age  
48 (<50 vs ≥50 years), molecular subtype (HR-/HER2-, TNBC, HER2+), treatment strategy  
49 (primary surgery vs neoadjuvant therapy), and pathological staging categories (T1 vs T2  
50 stage and N0 vs N+ stage).

51 Across most clinical strata, the model-predicted risk groups remained significantly  
52 associated with recurrence-free survival, demonstrating consistent risk stratification  
53 performance in the independent Duke dataset. Notably, even within conventionally defined  
54 clinical categories, substantial heterogeneity in recurrence outcomes was observed  
55 between the predicted low- and high-risk groups. This finding suggests that the imaging-  
56 based model captures prognostic information beyond standard clinicopathological factors.

## KM CURVE-Subgroup Analysis in DUKE Dataset



57

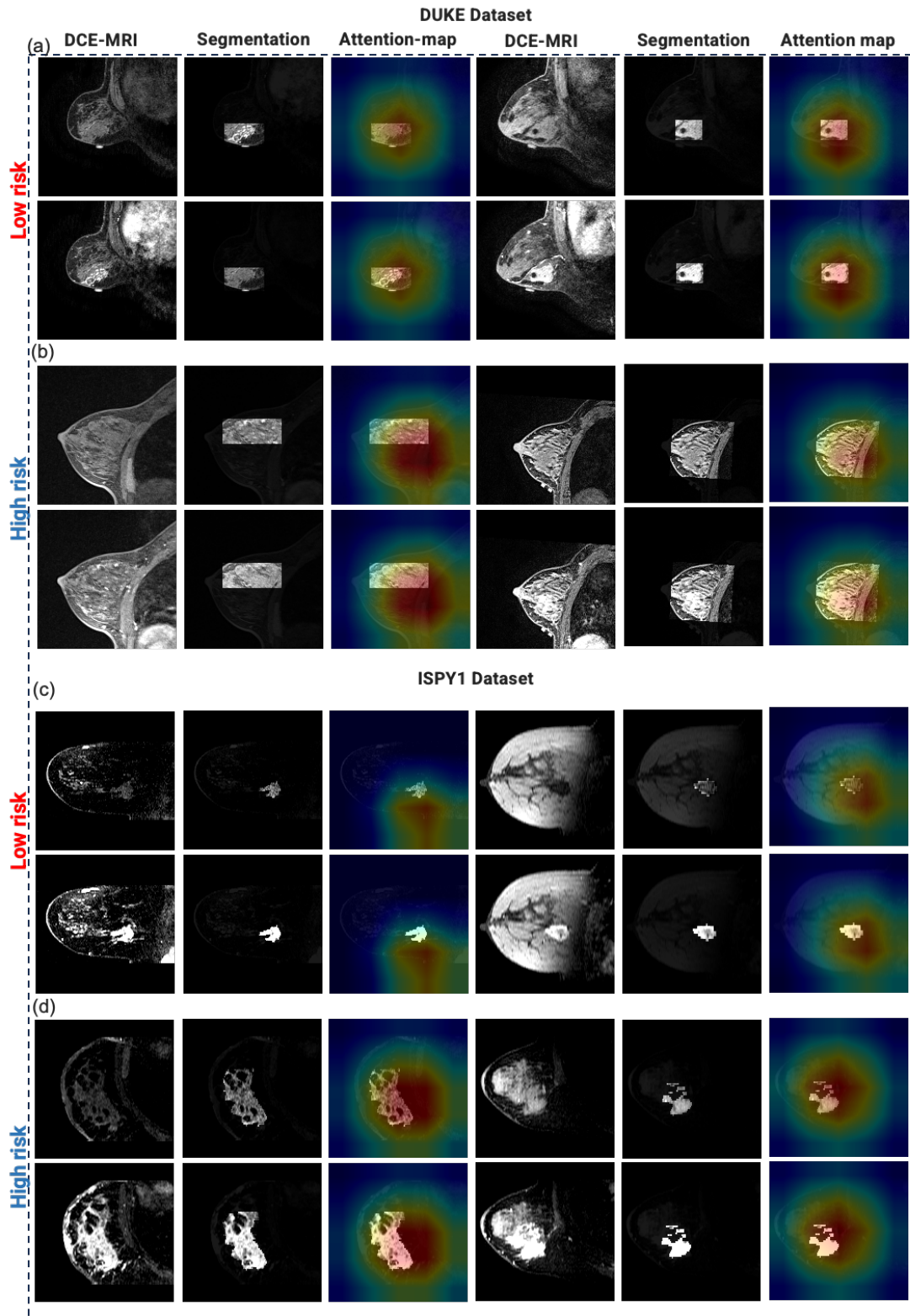
58 **Fig. S2 | Stratified Kaplan-Meier analysis demonstrating model robustness across**  
 59 **clinically relevant subgroups in the NKI cohort.** Recurrence-free survival curves  
 60 stratified by the AI Risk Score (Red: High-risk; Blue: Low-risk) within distinct patient  
 61 subgroups. a-j, Subgroup analyses performed on the NKI internal test set, categorized by:  
 62 Age: < 50 years and ≥ 50 years. Molecular Subtype: HR+/HER2-, TNBC, and  
 63 HER2+. Treatment: Surgery and Neoadjuvant Therapy (NA). Tumor Stage: T1 and T2, T3,  
 64 T4. Nodal Status: N0 and N+.

65

66 **Visual-cognitive alignment and interpretability of the**  
67 **multimodal framework from the external validation cohorts**

68 Supplementary Fig. S3 shows representative visualization examples from the external  
69 validation cohorts (Duke and I-SPY1). For each case, the original DCE-MRI image, lesion  
70 segmentation, and corresponding model attention map are displayed. Low-risk and high-  
71 risk examples are presented separately.

72 In both external datasets, the model attention maps predominantly localize to tumor  
73 regions defined by the lesion masks. The spatial attention patterns observed in these  
74 cohorts are consistent with those observed in the development dataset, suggesting that  
75 the model focuses on similar lesion-related imaging features across different institutions.



76

77 **Fig. S3 | Cross-center visual explainability and phenotypic characterization in**  
 78 **external validation cohorts.** Representative examples of model interpretability applied to  
 79 the Duke (Top panel) and I-SPY1 (Bottom panel) external validation datasets.

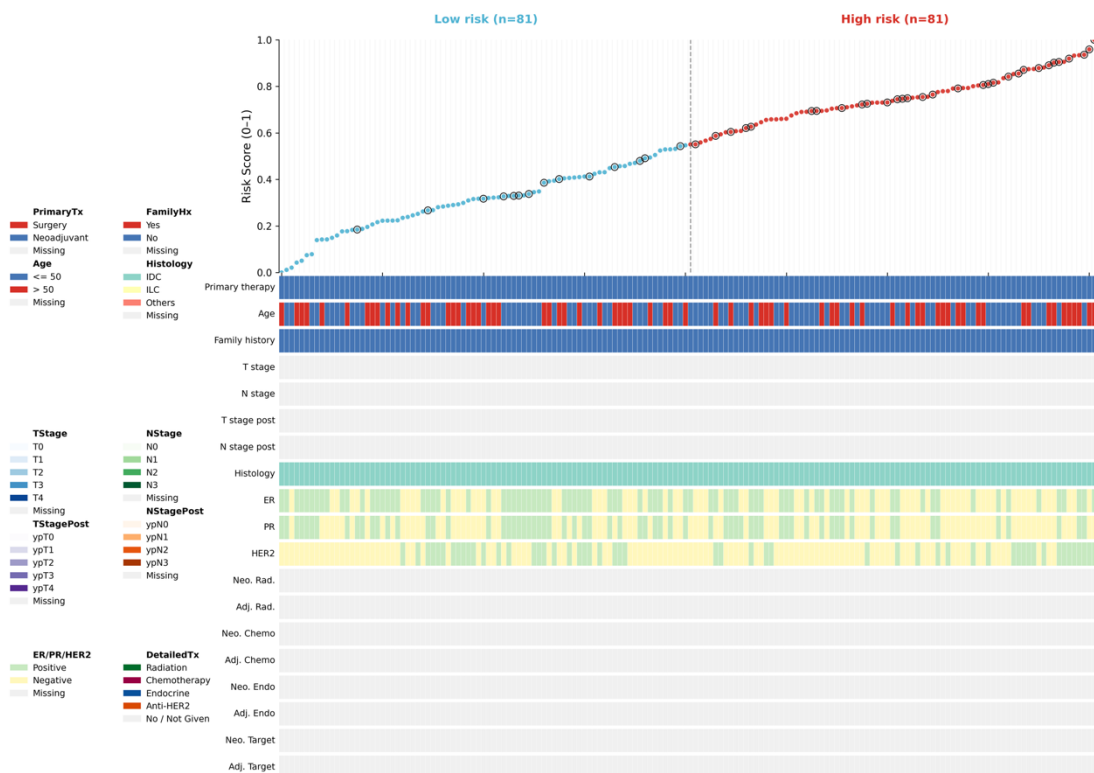
80

81 **Distribution of predicted risk scores and associated**  
 82 **clinicopathological characteristics within I-SPY1 cohort**

83 Supplementary Fig. S4 illustrates the distribution of AI-derived risk scores in the  
 84 independent I-SPY1 cohort, which served as an external test dataset. Patients are ranked  
 85 according to the predicted risk score generated by the model. Recurrence events are  
 86 indicated by black circles. The vertical dashed line represents the median risk score used  
 87 to divide the cohort into low-risk and high-risk groups.

88 Despite the relatively small cohort size and the higher overall recurrence prevalence  
 89 characteristic of the I-SPY1 neoadjuvant trial population, recurrence events are  
 90 predominantly concentrated toward the higher end of the predicted risk spectrum. The  
 91 lower panel displays the distribution of available clinicopathological variables across the  
 92 ranked patients, including treatment strategy, age group, histological subtype, hormone  
 93 receptor status, HER2 status, and systemic therapies. Many of these variables are either  
 94 relatively homogeneous or partially missing within this cohort.

95 The observed enrichment of recurrence events within the model-defined high-risk region  
 96 suggests that the imaging-derived risk score retains prognostic relevance in this  
 97 independent dataset, despite differences in cohort composition and clinical variable  
 98 availability compared with the development cohort.



99  
 100 **Fig. S4 | Risk landscape in the high-risk I-SPY1 neoadjuvant trial cohort. Patients are**

101 ranked by AI Risk Score. Despite the cohort's high baseline risk (reflected in the higher  
102 overall density of recurrence events, black circles) and relatively small sample size, the  
103 model successfully enriches events towards the high-risk spectrum (right). The ability of  
104 the model to achieve risk stratification in this homogenous, data-scarce setting—where  
105 standard clinical discriminators are either uniform or unavailable—confirms that the  
106 prognostic signal is primarily driven by intrinsic imaging phenotypes captured by the Visual  
107 Encoder, validating the robustness of the "Zero-Shot Modality Dropout" protocol.

108

## 109 **Detailed patient inclusion and exclusion criteria across cohorts**

110 This study included one internal cohort for model development and validation, and two  
111 independent external cohorts for testing. The detailed patient inclusion and exclusion  
112 criteria for each cohort are presented in Supplementary Figure S5.

### 113 **Development Cohort (Center A)**

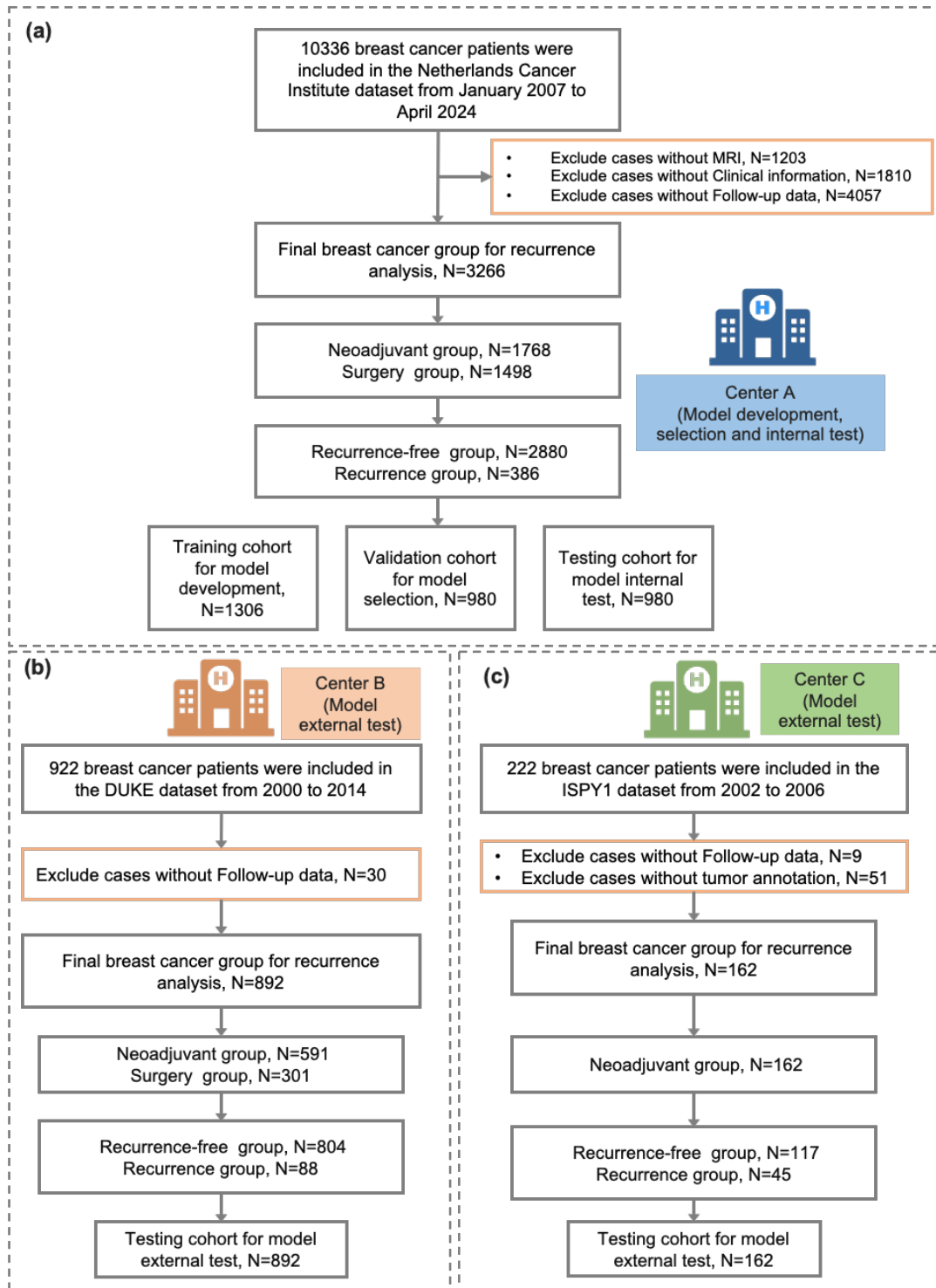
114 The primary cohort for model development, selection, and internal testing was  
115 retrospectively established from the Netherlands Cancer Institute (NKI) database. We  
116 initially identified 10,336 patients diagnosed with breast cancer between January 2007 and  
117 April 2024. Patients were excluded based on the following criteria: absence of pre-  
118 treatment Magnetic Resonance Imaging (MRI) (n=1,203), missing clinical information  
119 required for the model (n=1,801), or inadequate follow-up data for recurrence analysis  
120 (n=4,057). After applying these criteria, a final cohort of 3,266 patients remained for the  
121 analysis. This cohort was partitioned into a training set (n=1,306), a validation set used for  
122 hyperparameter tuning and model selection (n=980), and an internal test set for  
123 performance evaluation (n=980).

### 124 **External Test Cohort 1 (Center B)**

125 The first external test cohort was derived from the publicly available DUKE dataset,  
126 comprising patients treated between 2000 and 2014. From an initial set of 922 patients, 30  
127 were excluded due to the lack of follow-up data. The final cohort for external testing  
128 consisted of 892 patients.

### 129 **External Test Cohort 2 (Center C)**

130 The second external test cohort was sourced from the publicly available ISPY1 dataset,  
131 which also included patients from 2000 to 2014. The initial cohort consisted of 222 patients.  
132 We excluded 9 patients without follow-up data and 51 patients who lacked tumor  
133 annotations necessary for feature extraction. This resulted in a final cohort of 162 patients  
134 for the second external test.



135

136 **Fig. S5 | CONSORT flow diagrams illustrating patient inclusion, exclusion, and**  
 137 **cohort stratification.** a, Derivation of the primary development cohort (Center A: NKI). A  
 138 total of 10,336 breast cancer patients diagnosed between January 2007 and April 2024  
 139 were initially screened. Following stringent quality control, 7,061 patients were excluded  
 140 due to missing pre-treatment MRI (N=1,203), unavailable clinical information (N=1,801), or  
 141 insufficient follow-up data (N=4,057). The final NKI cohort comprised 3,266 patients,  
 142 including both neoadjuvant (N=1,768) and upfront surgery (N=1,498) groups. This dataset

143 was randomly split into independent training (N=1,306), validation (N=980), and internal  
144 testing (N=980) sets for model development and hyperparameter tuning. b, Selection of  
145 the Duke external validation cohort (Center B). Out of 922 patients screened from the Duke  
146 dataset (2000–2014), 30 were excluded for missing follow-up. The final cohort consisted  
147 of 892 patients (591 neoadjuvant, 301 surgery), serving as an independent test set to  
148 evaluate cross-center generalizability. c, Selection of the I-SPY1 external validation cohort  
149 (Center C). From the I-SPY1 clinical trial dataset (2002–2006), 60 patients were excluded  
150 due to missing follow-up or lack of tumor annotations. The final cohort comprised 162 high-  
151 risk patients, all of whom underwent neoadjuvant chemotherapy, providing a rigorous  
152 testbed for the model's performance in a specialized trial setting.

153

154

155 **MRI acquisition parameters of the in-house and external**  
156 **datasets**

157 DCE-MRI data were acquired at three centers with center-specific imaging protocols (Table  
158 S1). The in-house NKI cohort was acquired on 1.5-T and 3-T Philips scanners, with a  
159 repetition time (TR) of 3.66–6.33 ms, an echo time (TE) of 1.72–2.30 ms, and a flip angle  
160 of 10–12°. The in-plane spatial resolution was  $\leq 1.0 \times 1.0$  mm, with a slice thickness of  
161 1.20–1.80 mm. The public Duke cohort was acquired on 1.5-T and 3-T GE and Siemens  
162 scanners, with TR ranging from 4.1 to 6.0 ms, TE from 1.2 to 2.5 ms, and flip angle from 7  
163 to 12°. The in-plane resolution ranged from  $0.6 \times 0.6$  mm to  $1.0 \times 1.0$  mm, and slice  
164 thickness ranged from 1.1 to 2.5 mm. The public ISPY-1 cohort included scans acquired  
165 on 1.5-T GE, Siemens, and Philips systems, with TR  $\leq 20$  ms, TE of approximately 4.5 ms,  
166 and flip angle  $\leq 45^\circ$ . The in-plane resolution was  $\leq 1.0 \times 1.0$  mm and the slice thickness was  
167  $\leq 2.5$  mm.

168

169 **Table S1. MRI acquisition parameters of the in-house and external datasets.**  
 170 **Note:** DCE-MRI = dynamic contrast-enhanced magnetic resonance imaging; TR =  
 171 repetition time; TE = echo time. Values are presented as ranges, approximate values, or  
 172 upper limits according to the acquisition protocols reported by each center/dataset. NKI  
 173 represents the in-house dataset, whereas Duke and ISPY-1 represent publicly available  
 174 external datasets.

<b>Parameter</b>	<b>NKI (in-house)</b>	<b>Duke (public, external)</b>	<b>ISPY-1 (public, external)</b>
Vendor	Philips	GE, Siemens	GE, Siemens, Philips
Magnetic Field Strength	1.5 T, 3 T	1.5 T, 3 T	1.5 T
Sequence	DCE-MRI	DCE-MRI	DCE-MRI
TR (ms)	3.66 - 6.33	4.1 - 6.0	≤ 20
TE (ms)	1.72 - 2.30	1.2 - 2.5	~4.5
Flip angle (°)	10 - 12	7 - 12	≤ 45
In-plane resolution (mm)	≤ 1.0 × 1.0	0.6 × 0.6 - 1.0 × 1.0	≤ 1.0 × 1.0
Slice thickness (mm)	1.20 - 1.80	1.1 - 2.5	≤ 2.5

175