

# Machine learning radiomics for predicting recurrence risk in patients with early-stage invasive breast cancer

**Herui Yao** (✉ [yaohherui@mail.sysu.edu.cn](mailto:yaohherui@mail.sysu.edu.cn))

Sun Yat-sen Memorial Hospital, Sun Yat-sen University

**Yunfang Yu**

Sun Yat-sen Memorial Hospital, Sun Yat-sen University <https://orcid.org/0000-0003-2579-6220>

**Wei Ren**

Sun Yat-sen Memorial Hospital, Sun Yat-sen University

**Zifan He**

Sun Yat-sen Memorial Hospital, Sun Yat-sen University

**Yongjian Chen**

The Third Affiliated Hospital of Sun Yat-sen University

**Yujie Tan**

Sun Yat-sen Memorial Hospital, Sun Yat-sen University

**Jingwen Liu**

Sun Yat-sen Memorial Hospital, Sun Yat-sen University

**Anlin Li**

The First Clinical Medical College, Guangdong Medical University

**Nian Lu**

Sun Yat-sen University Cancer Center, State Key Laboratory of Oncology in South China, Collaborative Innovation Center for Cancer Medicine

**Jie Ouyang**

Tungwah Hospital, Sun Yat-sen University

**Yaping Yang**

Sun Yat-sen Memorial Hospital, Sun Yat-sen University

**Kai Chen**

Sun Yat-sen Memorial Hospital, Sun Yat-sen University

**Chenchen Li**

Sun Yat-sen Memorial Hospital, Sun Yat-sen University

**Mudi Ma**

Sun Yat-sen Memorial Hospital, Sun Yat-sen University

**Xiaohong Li**

Shunde Hospital, Southern Medical University

**Rong Zhang**

Shunde Hospital, Southern Medical University

**Zhuo Wu**

Sun Yat-sen Memorial Hospital, Sun Yat-sen University

**Fengxi Su Su**

Sun Yat-sen Memorial Hospital, Sun Yat-sen University

**Qiugen Hu**

Shunde Hospital, Southern Medical University

**Chuan-Miao Xie**

Sun Yat-sen University Cancer Center

**Erwei Song**

Sun Yat-sen Memorial Hospital <https://orcid.org/0000-0002-5400-9049>

---

**Article**

**Keywords:** Machine learning radiomics, recurrence risk, early-stage invasive breast cancer

**Posted Date:** October 7th, 2020

**DOI:** <https://doi.org/10.21203/rs.3.rs-81589/v1>

**License:**  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

# Abstract

There are no satisfying approaches to identify high- and low-risk recurrence patients with early-stage breast cancer in current clinical practice. Patients might be overtreated or undertreated due to the inaccurate prediction of recurrence risk. Herein, machine learning magnetic resonance imaging radiomic-based signature that integrates the intratumoral and peritumoral radiomic signatures, and clinicopathological characteristics was developed to classify high- and low-risk recurrence patients and predict recurrence within multicentre cohorts. The radiomic-clinical signature could also discriminate high- from low-risk recurrence patients among different breast cancer molecular subtype, and HR+/Her2-, T1N0M0 stage patients. Furthermore, it was observed that the neoadjuvant chemotherapy improved survival in high-risk Luminal subtype patients compared with the adjuvant chemotherapy. The survival-associated radiomic features also showed the correlation with the immune microenvironment. The radiomic-clinical signature presented the feasibility of predicting recurrence risk and assisting clinical decision-making in early-stage invasive breast cancer patients.

## Introduction

Breast cancer is the first leading cause of cancer death among women globally and approximately 10–15% of patients experience a recurrence in the first 5 years from diagnosis<sup>1,2</sup>. The St Gallen<sup>3</sup> consensus proposed clinicopathological risk categories to identify patients with high or low likelihood of recurrence, and suggested the use of endocrine monotherapy without adjuvant chemotherapy for clinical low-risk group. Nevertheless, there remained part of misclassified patients who might be overtreated or undertreated. Nowadays, 70-gene expression profile<sup>4</sup> and 21-gene recurrence score assay<sup>5</sup>, are currently recommended in clinical practice to predict recurrence risk and the benefit of adjuvant chemotherapy<sup>6</sup>, but the cost of these assays remains prohibitive and is an appropriate option only for Luminal subtype patients. Therefore, a more widely applicable and accurate method to pinpoint patients who are at high or low risk of recurrence is expected.

Several studies indicated that the radiomic features were significantly associated with tumor microenvironment and prognosis<sup>7</sup>, and a previous study has established a radiomic signature based on 294 invasive breast cancer patients to predict disease-free survival, but the model was difficult to be applied to clinical practice due to a small and single-centre dataset they used, no validation in different molecular subtype and lack of high-level evidences<sup>8</sup>. Recently, some studies showed that peritumoral radiomic features were also be predictive of prognosis rather than just focus on the tumor region<sup>9,10</sup>. This multicentre study aimed to construct a magnetic resonance imaging (MRI) radiomic-clinical signature that integrates intratumoral and peritumoral radiomic signatures, and clinicopathological characteristics for predicting the high and low recurrence risk in patients with early-stage invasive breast cancer.

## Results

### Patient characteristics

This study eligible 1,084 patients from four academic institutions in China (Supplementary Table 1), The study workflow was shown in Fig. 1. The table 1 showed the clinicopathological characteristics of patients in the training cohort (n=799), the prospective-retrospective validation cohort (n=105), and the external validation cohort (n=180). Adjuvant chemotherapy was administered to 709 (89%) of 799 patients in the training cohort, 57 (54%) of 105 patients in the prospective-retrospective validation cohort, and 156 (87%) of 180 patients in the external validation cohort. 105 patients underwent neoadjuvant chemotherapy from the prospective-retrospective validation cohort. Median follow-up was 22.8 months (IQR 15.5–35.4) for patients in the training cohort, 24.2 months (IQR 14.3–34.9) for those in the prospective-retrospective validation cohort, and 23.0 months (IQR 9.7–48.8) for those in the external validation cohort. The detailed information regarding the patient recruitment was described in Supplementary Fig. 1.

In the univariate analysis, which was presented in Supplementary Table 2, six differentially expressed clinical characteristics were found to be associated with RFS in the training cohort, including number of tumor ( $P < .001$ ), pathological N stage ( $P < .001$ ), histological grade ( $P < .001$ ), pathological tumor–node–metastasis (pTNM) stage ( $P < .001$ ), Ki-67 status ( $P < .001$ ), and Progesterone receptor status ( $P = .009$ ).

### **Intratumoral and peritumoral signatures for predicting recurrence risk**

The key radiomic features were selected by the Random forest algorithm to construct T1+C, T2WI, or DWI-ADC sequence signature by Cox regression, the detailed results were summarized in Supplementary Tables 3.

The intratumoral radiomic signature incorporated T1+C, T2WI, and DWI-ADC single sequence signature was conducted, which could assign patients into high- and low-risk groups. Patients with low-risk had better RFS in the training cohort (HR 0.06, 95% CI 0.02-0.15;  $P < .001$ ), the prospective-retrospective validation cohort ( $P = .039$ ), and the external validation cohort ( $P < .001$ ) (Supplementary Fig. 2a-c). In addition, the efficacy of the intratumoral radiomic signature showed AUCs of 0.86, 0.88, and 0.92 for 1-, 2-, 3-year RFS prediction in the training cohort, 0.86, 0.88, and 0.86 in the prospective-retrospective validation cohort, and 0.92, 0.94, and 0.91 in the external validation cohort, respectively (Supplementary Fig. 2d-f).

Simultaneously, the peritumoral radiomic signature was constructed and also presented the ability of discriminating high- from low-risk patients in the training cohort (HR 0.04, 95% CI 0.02-0.10;  $P < .001$ ), the prospective-retrospective validation cohort ( $P < .001$ ), and the external validation cohort ( $P = .043$ ) (Supplementary Fig. 3a-c). The peritumoral radiomic signature predicted AUCs of the 1-, 2-, and 3-year RFS of 0.96, 0.92, and 0.96 for the training cohort, 0.83, 0.86, and 0.86 for the prospective-retrospective validation cohort, and 0.87, 0.87 and 0.85 for the external validation cohort, respectively (Supplementary Fig. 3d-f).

### **Combined intratumoral and peritumoral signatures for predicting recurrence risk**

A radiomic signature combined both intratumoral radiomic signature and peritumoral radiomic signature was developed. The intratumoral-peritumoral radiomic signature categorized patients into high- and low-risk groups, which were significantly different in terms of RFS in the training cohort (HR 0.03, 95% CI 0.01–0.07;  $P < .001$ ), the prospective-retrospective validation cohort ( $P = .039$ ), and the external validation cohort ( $P < .001$ ) (Fig. 2a-c). Moreover, the intratumoral-peritumoral radiomic signature showed improved prediction in AUCs of the 1-, 2-, and 3-year RFS of 0.97, 0.95, and 0.98 in the training cohort, 0.86, 0.89, and 0.87 in the prospective-retrospective validation cohort, and 0.93, 0.94, and 0.91 in the external validation cohort, respectively, which presented a better predictive value than utilizing intratumoral or peritumoral radiomic signature alone (Fig. 2d-f).

In addition, the intratumoral-peritumoral radiomic signature was employed to classify high- and low-risk recurrence patients with the consideration of molecular subtype. Encouragingly, the radiomic signature could identify high- from low-risk patients in the subgroups of Luminal A ( $P < .001$ ), Luminal B ( $P < .001$ ), human epidermal growth factor receptor 2 (Her-2) positive ( $P = .007$ ), and triple-negative breast cancer (TNBC) ( $P < .001$ ) patients (Supplementary Fig. 4).

### **Radiomic-clinical signature for predicting recurrence risk**

To develop a more precisely and clinically applicable method that could predict an individual's recurrence, we took the clinicopathologic characteristics that associated with RFS in the univariate analysis into consideration. Multivariable analysis indicated that intratumoral-peritumoral radiomic signature, number of tumors, histological grade, pTNM stage, and Ki-67 status were independent factors of RFS (Supplementary Table 4), and these factors were used to construct the radiomic-clinical signature.

According to the radiomic-clinical signature, an optimal cutoff value (281) was generated to classify patients into high- and low-risk groups in the training cohort. A radiomic-clinical signature-print was showed to illustrate the association of these factors with the recurrence risk. The intratumoral-peritumoral radiomic signature presented the largest proportion in both high-risk (83%) and low-risk (45%) recurrence groups, followed by the histological grade (high-risk, 68%; low-risk, 44%) (Fig. 3).

This radiomic-clinical signature assigned 47 (10.4%) of 452 patients to the high-risk group, and there were significant differences in RFS between high-risk and low-risk groups (HR 0.03, 95% CI 0.01-0.08,  $P < .001$ ). In the prospective-retrospective validation cohort, 40 (59.7%) of 67 patients were separated into high-risk group, which had shorter RFS ( $P = .030$ ). In the external validation cohort, 34 (26.6%) of 128 patients with high-risk had shorter RFS ( $P < .001$ ) (Fig. 4a-c). Besides, the radiomic-clinical signature showed better performance of RFS prediction, which achieved the higher 1-, 2-, 3-year AUCs (0.97, 0.96, and 0.98) in the training cohort, the prospective-retrospective validation cohort (AUCs of 0.88, 0.93 and 0.93), and the external validation cohort (AUCs of 0.94, 0.96 and 0.94), respectively (Fig. 4d-f).

In addition, the radiomic-clinical signature demonstrated the capability of precisely predicting recurrence risk and could be used for identifying high- and low-risk patients among different molecular subtype ( $P < .001$  for Luminal A,  $P < .001$  for Luminal B,  $P = .007$  for Her2-positive,  $P < .001$  for TNBC) (Fig. 5a-d). For

Luminal subtype patients in the high-risk group who received the neoadjuvant chemotherapy showed significantly prolonged RFS ( $P = .048$ ) compared with patients who received the adjuvant chemotherapy, whereas there was no added benefit of the neoadjuvant chemotherapy for patients in the low-risk group (Supplementary Fig. 5). Moreover, among Luminal subtype (T1N0M0 stage, HR-positive and Her2-negative status) patients, the radiomic-clinical signature could recognize high and low-risk patients ( $P < .001$ ; Fig. 5e), including the subgroups analysis of patients who received adjuvant chemotherapy ( $P < .001$ ; Fig. 5f).

### **Radiomic features associated with tumor immune microenvironment and genomics**

The key radiomic features from intratumoral T1+C and T2WI sequences of The Cancer Genome Atlas (TCGA) and The Cancer Imaging Archive (TCIA) were found to be correlated linearly with the immune cells (Fig. 6a). The activated natural killer cells were observed to have a positive correlation with the most radiomic features. The M0 macrophages, T cells regulatory Tregs and T cells follicular helper also presented a strong correlation. In addition, we had previously identified 29 lncRNAs which were associated with survival and immune response<sup>11</sup>. In this study, most of the lncRNAs were indicated to be remarkably correlated with the radiomic features (Fig. 6a), including the NKILA, which had been proved to play an important role in immune microenvironment in a previous study<sup>12</sup>. These results illustrated that the radiomic features could provide important information about tumor immune microenvironment.

Different classes of the radiomic features were identified using the unsupervised consensus clustering analysis in patients from TCGA and TCIA. A total of 536 differentially expressed lncRNAs and 835 differentially expressed genes were identified to be associated with radiomic features. Then the unsupervised consensus clustering analysis was performed with these lncRNAs and genes in 1,082 breast cancer patients. Two main radiomic-based lncRNA subtypes were identified to be associated with significant difference in overall survival (HR 0.71, 95% CI 0.51–0.97;  $P = .031$ ) (Fig. 6b). Next, the Gene Ontology (GO) and Kyoto Encyclopedia of Genes and Genomes (KEGG) analysis were conducted to evaluate the enrichment of the radiomic-based genes. The GO enrichment analysis indicated that the radiomic-based genes were enriched in various physiological metabolic processes, such as affection of oxidoreductase activity, lipid metabolism and potassium channel complex, details were illustrated in Fig. 6c. The KEGG pathway enrichment analysis found these genes were involved in the vitamin digestion and absorption and peroxisome proliferator-activated receptor signaling pathway.

## **Discussion**

In this multicentre cohort study, the intratumoral-peritumoral radiomic signature based on machine learning algorithm discriminated high- from low-risk recurrence patients and performed well in predicting RFS. The radiomic-clinical signature comprised the intratumoral-peritumoral radiomic signature and clinicopathological characteristics was found to be significantly associated with RFS and presented higher predictive value in RFS. In addition, the radiomic-clinical signature successfully classified high and low recurrence risk among different breast cancer molecular subtype patients, and HR+/Her2- (T1N0M0

stage) patients. The key radiomic features were also found to be associated with immune microenvironment. Therefore, this study developed and validated a prognostic, radiomic-clinical signature for individualized prediction of high and low recurrence risk, which provided an effective tool for prediction of survival and clinical decision-making in patients with early-stage invasive breast cancer.

While previous studies<sup>13,14</sup> showed the potential of MRI-based radiomics for predicting recurrence in breast cancer, their clinical value were limited because of the small sample size from single-centre and the radiomic features only extracted from tumor region. Our study built a radiomic-clinical signature with integrating the intratumoral-peritumoral radiomic signature and clinicopathological characteristics based on a more than 1000-patient size from multicentre and independent external validation cohort. Our results showed that the radiomic-clinical signature played an important role in predicting recurrence, and the signature-print indicated that radiomic features were more associated with recurrence risk than clinicopathological characteristics. Thus we proposed to combine radiomic features with clinicopathological characteristics, which could better predict recurrence and utilize in clinical practice.

In the past few decades, the high- and low-risk recurrence was mainly evaluated by the clinicopathological characteristics of the patients. A study retrospectively reviewed 1,500 patients with node-negative breast cancer and found that using the 2007 St Gallen risk categories resulted in different outcomes<sup>3,15</sup>. The St Gallen divided patients with node negative, pathological tumor size  $\leq 2$  cm, histological grade 1, absence of extensive peritumoral vascular invasion, HR+, HER2- and age  $\geq 35$  years into low-risk group. In our study, using St Gallen categories assigned only 13 (11%) of 118 patients into low-risk group among patients with Luminal subtype (HR+/HER2-, and T1N0M0 stage), and was unable to further recognize high and low recurrence risk (Supplementary Fig. 6). Additionally, the intermediate- and high-risk groups defined according to St Gallen categories included many patients who had a good outcome, which indicated that St Gallen criteria contained misclassified patients who might be potentially overtreated with adjuvant chemotherapy. However, our radiomic-clinical signature could assign 110 (93%) of 118 HR+/Her2-, T1N0M0 stage patients into low-risk group, which could minimize the probability of being overtreated.

Nowadays, multigene profiles were constructed for risk stratification and therapy strategies guidance in breast cancer<sup>16,17</sup>. The randomized trial TAILORx<sup>5,18</sup> enrolled patients with HR+, HER2-, and axillary node-negative breast cancer, and demonstrated that patients with low range 21-gene recurrence score could avoid adjuvant chemotherapy. Another randomized trial MINDACT<sup>4</sup> allowed enrollment of patients with up to three positive axillary nodes, and showed the ability of identifying patients with high clinical risk who can avoid chemotherapy by testing 70-gene signature. In this study, the radiomic-clinical signature displayed the competence of discriminating high- from low-risk patients in different breast cancer molecular subtype, which indicated that the combination of multigene profiles and the radiomic-clinical signature might distinct high and low recurrence risk more precisely and reduce the rate of undertreatment and overtreatment in future clinical practice.

In current clinical practice, patients with large tumor size, positive axillary nodes would consider to receive neoadjuvant chemotherapy. The results of randomized trials NSBPA-18 and NSABP-27<sup>19</sup> manifested that though neoadjuvant chemotherapy was equivalent to adjuvant chemotherapy, patients who achieved a pathologic complete responses after neoadjuvant chemotherapy had significantly prolonged survival and lower risk of recurrence compared with patients who did not. Although pathologic complete responses have been shown to be predictive of benefit from neoadjuvant chemotherapy, it could only be evaluated after surgery, therefore a preoperative approach is urgently needed to distinguish patients who could benefit from neoadjuvant chemotherapy.

In this study, the radiomic-clinical signature could predict RFS and identify high- and low-risk recurrence in the prospective-retrospective validation cohort of which all of the patients underwent neoadjuvant chemotherapy. It is worth noting that the neoadjuvant chemotherapy improved RFS in high-risk Luminal subtype patients compared with the adjuvant chemotherapy. However, larger sample sizes and longer follow-up time were needed to further validate radiomic-clinical signature's potential of recognizing patients who could obtain more benefit from neoadjuvant chemotherapy.

Several limitations still existed in this study. The heterogeneity MRI scans from multiple centres was inevitable, and the median follow-up was about 24 months, the signatures could not be applied for predicting overall survival. Previous studies have shown the association between radiomic features and immune response<sup>8,20</sup>. In this study, we analyzed the correlation of radiomic features with immune cells and lncRNAs, and significant linear correlation was presented, which indicated that the radiomic features can provide tumor immune microenvironment information. We also evaluated the radiomic-based genes and related pathways. However, due to the retrospective approach taken in this study and the lack of available data of gene expression, we were unable to further analyze the more association between radiomic features and tumor microenvironment, especially the mechanisms of using the radiomic features to predict recurrence need to be further explored. It may be beneficial to comprise radiomic signatures with genetic signatures such as genomics and transcriptomics, which had better prediction ability of recurrence and clinical application value.

In conclusion, this study presented a radiomic-clinical signature that incorporated MRI intratumoral-peritumoral radiomic signature and clinicopathological characteristics, which could identify high- and low-risk recurrence among different molecular subtype and be conveniently used for individualized prediction of RFS in patients with early-stage invasive breast cancer.

## Methods

### Study design and patients

This study was conducted in accordance with the STROBE guideline checklist<sup>21</sup>. A total of 1,161 early-stage invasive breast cancer patients were retrospective recruited from four institutions in China, of which 1,084 patients (a total of 178,847 images) passed quality control. A total of 799 patients recruited from

Sun Yat-sen Memorial Hospital of Sun Yat-sen University, a national hospital (Guangzhou, China) and Sun Yat-sen University Cancer center, a national hospital (Guangzhou, China) between March 23, 2011, and August 26, 2019 were assigned into a training cohort. The patients consisted of 105 patients from the prospective phase III clinical trials [NCT01503905] collected from the Sun Yat-sen Memorial Hospital of Sun Yat-sen University (Guangzhou, China) between July 20, 2015, and April 22, 2019 were assigned into the prospective-retrospective validation cohort. A total of 180 patients collected from the Shunde Hospital of Southern Medical University (Foshan, China) and Tungwah Hospital of Sun Yat-sen University (Dongguan, China) between March 09, 2012, and September 21, 2019 were used as an independent external validation cohort.

The primary outcome was recurrence-free survival (RFS), RFS was calculated from the surgery date to the date of most recent medical review or diagnosis of recurrence. The inclusion criteria were female patients aged at least 18 years with histologically confirmed as stage I–III invasive breast cancer<sup>22</sup> and patients underwent breast tumor and axillary MRI scans before surgery and axillary lymph node dissection. Patients suffering from other tumor diseases before or at the same time, having incomplete pathological information, or unavailable standard MRI scans with or without contrast enhancement were excluded.

### **Radiomic feature extraction**

The multiparametric MRI (contrast-enhanced T1-weighted imaging [T1+C], T2-weighted imaging [T2WI], and diffusion-weighted imaging quantitatively measured apparent diffusion coefficient [DWI-ADC]) acquisition protocol across all institutions and MR scanner parameters for patients were described in Supplementary 1 and Supplementary Table 5. All of the MRIs were normalized to obtain a standard normal distribution of image intensities using the N4ITK Bias Correction code. 3D regions of interest (ROIs) of the breast intratumoral area, and peritumoral area (the tumor parenchymal constituting 10-mm extension outward) were semi-automatically segmented by 3D Slicer software method (<https://www.slicer.org/>, version 4.10.2)<sup>23</sup>. The 3D regions of intratumoral and peritumoral (DICOM format) was transferred to the SlicerRadiomics code, the in-house texture extraction platform developed based on the python package “PyRadiomics”. A total of 5,178 quantitative radiomic features, including six groups of radiomic features were extracted separately, including shape, first-order, the gray-level co-occurrence matrix (GLCM), the gray-level size zone matrix (GLSZM), the gray-level dependence matrix (GLDM), and the neighbouring gray tone difference matrix (NGTDM). More details regarding the radiomic feature extraction was described in Supplementary 2.

### **Radiomic signature building and validation**

The Random forest algorithm<sup>24</sup> was applied to select the most predictive candidate radiomic features in the training cohort. The combination of the key features in each sequence were used to constructed T1+C, T2WI, and DWI-ADC single sequence signature of intratumoral or peritumoral. The intratumoral or peritumoral radiomic signature incorporated T1+C, T2WI, and DWI-ADC single sequence signature were calculated by Cox regression. Intratumoral-peritumoral radiomic signature combined intratumoral and

peritumoral radiomic signatures was calculated by Cox regression model with the radiomic scores in the training cohort. The radiomic scores calculated for each patient via a combination of selected features that were weighted by their respective coefficients. The intratumoral-peritumoral radiomic signature predicted RFS was then assessed in the prospective-retrospective validation cohort and the external validation cohort, respectively. The essential radiomic features and formula composition were presented in Supplementary Table 6.

### **Radiomic-clinical signature building and validation**

The univariate analysis was used to assess the association between clinicopathological characteristics and RFS in the training cohort. A multivariate regression analysis was used to test the independent significance of the intratumoral-peritumoral radiomic signature and significant clinical variables in relation to RFS. Finally, to provide the clinician a quantitative tool to predict individual probability of recurrence, a radiomic-clinical signature combined the intratumoral-peritumoral radiomic signature and significant clinical variables was constructed by Cox regression model. The performance of the radiomic-clinical signature was validated in the prospective-retrospective and the external validation cohorts.

### **Radiomic features associated with tumor immune microenvironment and genomics**

To quantify the proportions of tumor immune microenvironment in the 90 breast cancer patients from TCGA and TCIA, the CIBERSORT algorithm<sup>25</sup> and the LM22 gene signature were used for highly sensitive and specific discrimination of 22 human immune cell phenotypes including B cells, T cells, natural killer cells, macrophages, dendritic cells, and myeloid subsets. CIBERSORT is a deconvolution algorithm that uses a set of reference gene expression values (a signature with 547 genes) that is considered a minimal representation for each cell type. Based on those values, CIBERSORT infers cell type proportions in data from bulk tumor samples with mixed cell types using support vector regression. Gene expression profiles were prepared using standard annotation files, the data were uploaded to the CIBERSORT web portal (<http://cibersort.stanford.edu/>), and the algorithm was run using the LM22 signature at 1,000 permutations. The 29 lncRNAs were selected using the univariable Cox proportional hazards regression model and the LASSO algorithm in patients treated with immunotherapy from the IMvigor210 trial<sup>11</sup>.

The unsupervised hierarchical clustering methods (K-means)<sup>26</sup> were used to identify different classes of the radiomic features in 71 patients, who had both T1+C and T2WI sequences from TCGA and TCIA with the ConsensusClusterPlus R package<sup>27</sup>. The *t*-test and R package limma<sup>28</sup> were utilized to identify differentially expressed lncRNAs and genes associated with the key radiomic features, respectively. Differentially expressed radiomic-based lncRNAs and genes were divided into different clusters in 1,082 breast cancer patients with transcriptome RNA sequencing data from TCGA using the ConsensusClusterPlus R package<sup>27</sup>. The GO and KEGG analysis were performed using the clusterProfiler R package<sup>29</sup>. The GO terms and KEGG pathways were considered statistically significant with *P* values and false discovery rates less than .05.

## Statistical analysis

The Fisher's exact test was performed to examine the differences in the occurrence of categorical variables, while the independent *t*-test was used to compare the differences in continuous variables between two groups. Survival was calculated using the Kaplan-Meier method and the log-rank test, and hazard ratios (HRs) and 95% confidence intervals (CIs) were calculated using a Cox regression analysis. Patients were categorized into high and low-risk groups with the optimal cutoff values defined by the R package *ggsurvimier*. The prognostic or predictive accuracy of the signatures was assessed by using operating characteristic curve (ROC) analysis. The area under ROC curve (AUC) was used to assess the sensitivity and specificity, and this was calculated to evaluate the performance of signatures for predicting RFS. For all the analyses, two-sided P-values less than 0.05 were considered statistically significant. Statistical analyses were performed using R software (version 4.0.0). This study is registered with ClinicalTrials.gov, number NCT04003558.

## Declarations

**Ethics:** This multicentre study was conducted in accordance with the Declaration of Helsinki. The study's protocol was approved by the ethics committee of each participating hospital (Sun Yat-sen Memorial Hospital of Sun Yat-sen University, SYSEC-KY-KS-2019-054-001; Sun Yat-sen University Cancer Center, B2020-114-01; Shunde Hospital of Southern Medical University, KYLS-20190579; Tungwah Hospital of Sun Yat-sen University, 2020DHLL018). The requirement for informed consents in retrospective cohorts were waived. Participants from the prospective phase III clinical trials NCT01503905 have signed informed consents, the trial was also approved by the ethics committee with number 2011 EC # (12).

## Article information

**Author Contributions:** All authors had full access to all the data in the study and take responsibility for the integrity of the data and the accuracy of the data analysis. Herui Yao, Yunfang Yu, Wei Ren, Zifan He, Yongjian Chen, Yujie Tan are co-first authors. Erwei Song, Herui Yao, Qiugen Hu, Chuanmiao Xie are co-corresponding authors.

**Concept and design:** All authors.

**Acquisition, analysis, or interpretation of data:** Herui Yao, Yunfang Yu, Wei Ren, Zifan He, Qiugen Hu, Chuanmiao Xie, Erwei Song.

**Drafting of the manuscript:** All authors.

**Critical revision of the manuscript for important intellectual content:** All authors.

**Statistical analysis:** All authors.

**Obtained funding:** Herui Yao.

**Administrative, technical, or material support:** Herui Yao, Yunfang Yu, Wei Ren, Zifan He, Qiugen Hu, Chuanmiao Xie, Erwei Song.

**Supervision:** Herui Yao.

**Conflict of Interest Disclosures:** The authors have no conflicts of interest to declare.

**Funding:** This study was supported by grant 2020ZX09201021 from the National Science and Technology Major Project, grant YXRGZN201902 from the Medical Artificial Intelligence Project of Sun Yat-Sen Memorial Hospital, grants 81572596, 81972471, and U1601223 from the National Natural Science Foundation of China, grant 2017A030313828 from the Natural Science Foundation of Guangdong Province, grant 201704020131 from the Guangzhou Science and Technology Major Program, grant 2017B030314026 from the Guangdong Science and Technology Department, grant 2018007 from the Sun Yat-Sen University Clinical Research 5010 Program, grant SYS-C-201801 from the Sun Yat-Sen Clinical Research Cultivating Program.

## References

1. Bray, F. et al. Global cancer statistics 2018: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries. *CA Cancer J Clin* **68**, 394–424 (2018).
2. Colleoni, M. et al. Annual Hazard Rates of Recurrence for Breast Cancer During 24 Years of Follow-Up: Results From the International Breast Cancer Study Group Trials I to V. *J Clin Oncol* **34**, 927–35 (2016).
3. Goldhirsch, A. et al. Progress and promise: highlights of the international expert consensus on the primary therapy of early breast cancer 2007. *Ann Oncol* **18**, 1133–44 (2007).
4. Cardoso, F. et al. 70-Gene Signature as an Aid to Treatment Decisions in Early-Stage Breast Cancer. *N Engl J Med* **375**, 717–29 (2016).
5. Sparano, J.A. et al. Prospective Validation of a 21-Gene Expression Assay in Breast Cancer. *N Engl J Med* **373**, 2005–14 (2015).
6. Gradishar, W.J. et al. Breast Cancer, Version 3.2020, NCCN Clinical Practice Guidelines in Oncology. *J Natl Compr Canc Netw* **18**, 452–478 (2020).
7. Sun, R. et al. A radiomics approach to assess tumour-infiltrating CD8 cells and response to anti-PD-1 or anti-PD-L1 immunotherapy: an imaging biomarker, retrospective multicohort study. *Lancet Oncol* **19**, 1180–1191 (2018).
8. Park, H. et al. Radiomics Signature on Magnetic Resonance Imaging: Association with Disease-Free Survival in Patients with Invasive Breast Cancer. *Clin Cancer Res* **24**, 4705–4714 (2018).
9. Vaidya, P. et al. CT derived radiomic score for predicting the added benefit of adjuvant chemotherapy following surgery in Stage I, II resectable Non-Small Cell Lung Cancer: a retrospective multi-cohort study for outcome prediction. *Lancet Digit Health* **2**, e116–e128 (2020).

10. Braman, N. et al. Association of Peritumoral Radiomics With Tumor Biology and Pathologic Response to Preoperative Targeted Therapy for HER2 (ERBB2)-Positive Breast Cancer. *JAMA Netw Open* **2**, e192561 (2019).
11. Yu, Y. et al. Association of Long Noncoding RNA Biomarkers With Clinical Immune Subtype and Prediction of Immunotherapy Response in Patients With Cancer. *JAMA Netw Open* **3**, e202149 (2020).
12. Huang, D. et al. NKILA lncRNA promotes tumor immune evasion by sensitizing T cells to activation-induced cell death. *Nat Immunol* **19**, 1112–1125 (2018).
13. Chitalia, R.D. et al. Imaging Phenotypes of Breast Cancer Heterogeneity in Preoperative Breast Dynamic Contrast Enhanced Magnetic Resonance Imaging (DCE-MRI) Scans Predict 10-Year Recurrence. *Clin Cancer Res* **26**, 862–869 (2020).
14. Mazurowski, M.A. et al. Association of distant recurrence-free survival with algorithmically extracted MRI characteristics in breast cancer. *J Magn Reson Imaging* **49**, e231-e240 (2019).
15. Blancas, I. et al. Outcome differences between patients with node-negative breast cancer classified according to the st. Gallen risk categories. *Clin Breast Cancer* **9**, 231–6 (2009).
16. Paik, S. et al. A multigene assay to predict recurrence of tamoxifen-treated, node-negative breast cancer. *N Engl J Med* **351**, 2817–26 (2004).
17. van de Vijver, M.J. et al. A gene-expression signature as a predictor of survival in breast cancer. *N Engl J Med* **347**, 1999–2009 (2002).
18. Sparano, J.A. et al. Clinical and Genomic Risk to Guide the Use of Adjuvant Therapy for Breast Cancer. *N Engl J Med* **380**, 2395–2405 (2019).
19. Rastogi, P. et al. Preoperative chemotherapy: updates of National Surgical Adjuvant Breast and Bowel Project Protocols B-18 and B-27. *J Clin Oncol* **26**, 778–85 (2008).
20. Aerts, H.J. et al. Decoding tumour phenotype by noninvasive imaging using a quantitative radiomics approach. *Nat Commun* **5**, 4006 (2014).
21. von Elm, E. et al. Strengthening the Reporting of Observational Studies in Epidemiology (STROBE) statement: guidelines for reporting observational studies. *Bmj* **335**, 806–8 (2007).
22. Amin, M.B. et al. The Eighth Edition AJCC Cancer Staging Manual: Continuing to build a bridge from a population-based to a more "personalized" approach to cancer staging. *CA Cancer J Clin* **67**, 93–99 (2017).
23. Fedorov, A. et al. 3D Slicer as an image computing platform for the Quantitative Imaging Network. *Magn Reson Imaging* **30**, 1323–41 (2012).
24. Breiman, L. Random forests. *Mach Learn* **45**, 5–32 (2001).
25. Newman, A.M. et al. Robust enumeration of cell subsets from tissue expression profiles. *Nat Methods* **12**, 453–7 (2015).
26. Hartigan, J.A. & Wong, M.A. Algorithm AS 136: A K-Means Clustering Algorithm. *Journal of the Royal Statistical Society. Series C (Applied Statistics)* **28**, 100–108 (1979).

27. Monti, S., Tamayo, P., Mesirov, J. & Golub, T. Consensus Clustering: A Resampling-Based Method for Class Discovery and Visualization of Gene Expression Microarray Data. *Machine Learning* **52**, 91–118 (2003).
28. Ritchie, M.E. et al. limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Res* **43**, e47 (2015).
29. Yu, G., Wang, L.G., Han, Y. & He, Q.Y. clusterProfiler: an R package for comparing biological themes among gene clusters. *Omic*s **16**, 284–7 (2012).

## Tables

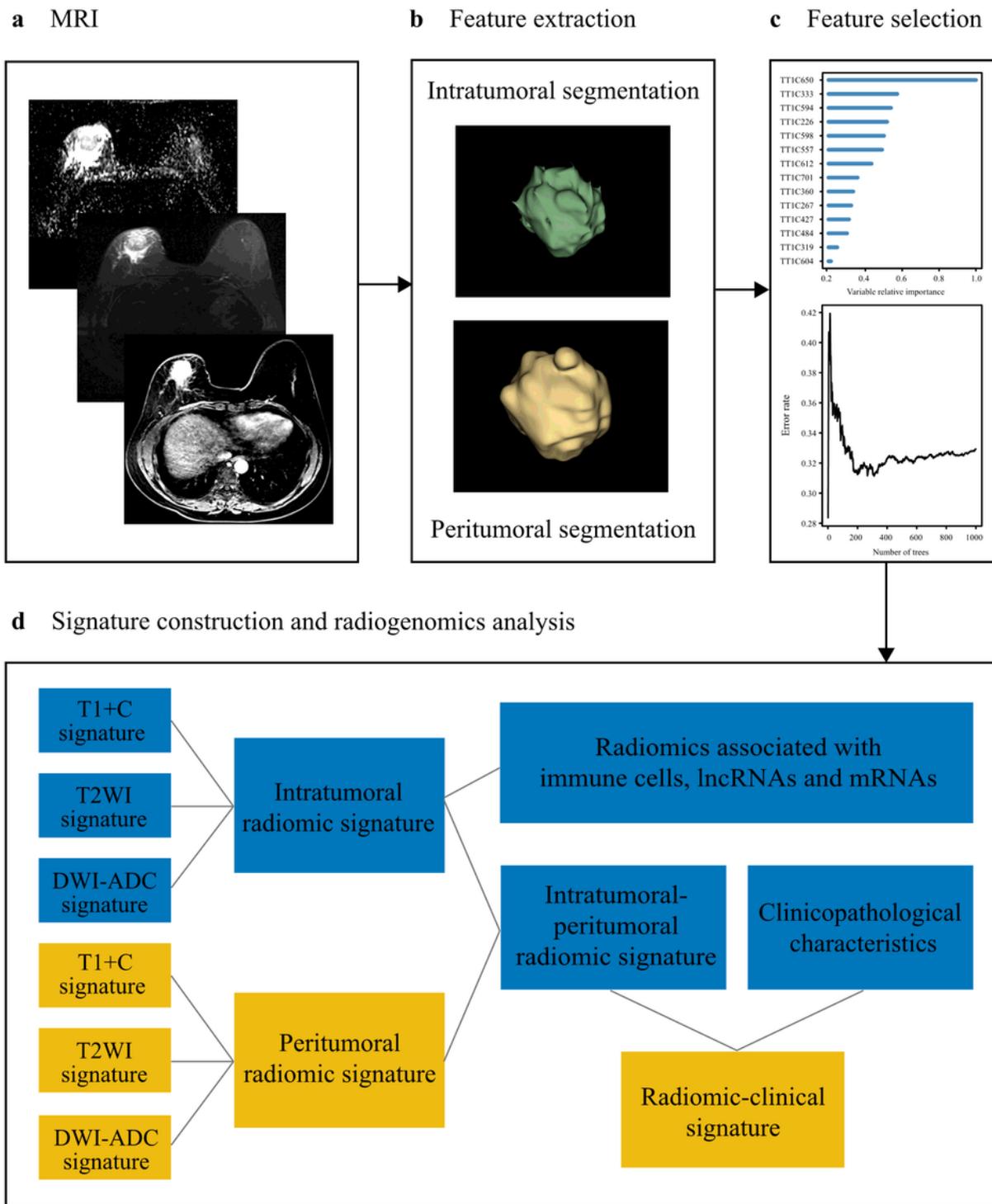
**Table 1.** Clinicopathological Characteristics of Patients in the Training, Prospective-retrospective validation and External validation Cohorts

Characteristics	Training Cohort (No. of patients [n]=799)	Prospective-retrospective validation Cohort (n=105)	External validation Cohort (n=180)
Follow-up time, months (median [IQR])	22.75 [15.49, 35.39]	24.15 [14.30, 33.97]	23.00 [9.70, 48.77]
Age, years (median [IQR])	48 [42, 56]	42 [38, 45]	48 [43, 56]
Number of tumor (%)			
1	701 (87.8)	95 (90.5)	148 (82.2)
>1	97 (12.2)	10 (9.5)	32 (17.8)
Tumor size, cm (median [IQR])	2.2 [1.7, 3.0]	3.5 [2.8, 4.9]	2.4 [1.8, 3.3]
Clinical T stage (%)			
T1	336 (42.1)	10 (9.5)	59 (32.8)
T2	420 (52.6)	66 (62.9)	107 (59.4)
T3	31 (3.9)	21 (20.0)	10 (5.6)
T4	12 (1.5)	8 (7.6)	4 (2.2)
Clinical N stage (%)			
N0	510 (63.8)	27 (26.0)	134 (74.4)
N1	268 (33.5)	70 (67.3)	24 (13.3)
N2	20 (2.5)	6 (5.8)	18 (10.0)
N3	1 (0.1)	1 (1.0)	4 (2.2)
Clinical TNM stage (%)			
I	256 (32.0)	4 (3.8)	55 (30.6)
II	492 (61.6)	77 (74.0)	96 (53.3)
III	51 (6.4)	23 (22.1)	29 (16.1)
Histological grade (%)			
Grade 1 (low)	24 (3.1)	1 (1.2)	9 (5.7)
Grade 2 (intermediate)	382 (50.1)	55 (64.7)	81 (51.3)
Grade 3 (high)	356 (46.8)	29 (34.1)	68 (43.0)
Pathological T stage (%)			
T1	410 (51.3)	53 (51.5)	70 (38.7)
T2	357 (44.7)	36 (35.0)	99 (55.0)
T3	29 (3.6)	5 (4.9)	8 (4.4)
T4	3 (0.4)	9 (8.7)	3 (1.7)
Pathological N stage (%)			
N0	475 (59.4)	28 (26.7)	94 (52.2)
N1	215 (26.9)	35 (33.3)	47 (26.1)
N2	61 (7.6)	32 (30.5)	22 (12.2)
N3	48 (6.0)	10 (9.5)	17 (9.4)
Pathological TNM stage (%)			
I	271 (33.9)	23 (22.3)	45 (25.0)
II	404 (50.6)	34 (33.0)	92 (51.1)
III	124 (15.5)	46 (44.7)	43 (23.9)
ER status (%)			
Negative	115 (14.5)	8 (7.6)	44 (24.6)
Positive	680 (85.5)	97 (92.4)	135 (75.4)
PR status (%)			

Negative	218 (27.4)	32 (30.5)	59 (33.0)
Positive	577 (72.6)	73 (69.5)	120 (67.0)
Her2 status (%)			
Negative	532 (69.3)	57 (66.3)	105 (68.2)
Positive	236 (30.7)	29 (33.7)	49 (31.8)
Ki67 expression (%)			
<30	407 (51.2)	68 (64.8)	101 (58.0)
≥30	388 (48.8)	37 (35.2)	73 (42.0)
Molecular subtypes (%)			
Luminal A	122 (15.6)	21 (22.6)	35 (21.3)
Luminal B	551 (70.6)	65 (69.9)	92 (56.1)
Her2-positive	59 (7.6)	4 (4.3)	17 (10.4)
Triple negative	48 (6.2)	3 (3.2)	20 (12.2)
Type of surgery (%)			
Breast-conserving surgery	373 (46.7)	49 (46.7)	41 (22.8)
Others	425 (53.3)	56 (53.3)	139 (77.2)
Adjuvant chemotherapy (%)			
Yes	709 (88.7)	57 (54.3)	156 (86.7)
No	90 (11.3)	44 (45.7)	24 (13.3)

Abbreviations: IQR, interquartile range; CI, confidence interval; DFS, disease-free survival; TNM, tumor-node-metastasis; Her2, human epidermal growth factor receptors 2; ER, estrogen receptor; PR, progesterone receptors; Ki67, proliferation marker protein Ki-67.

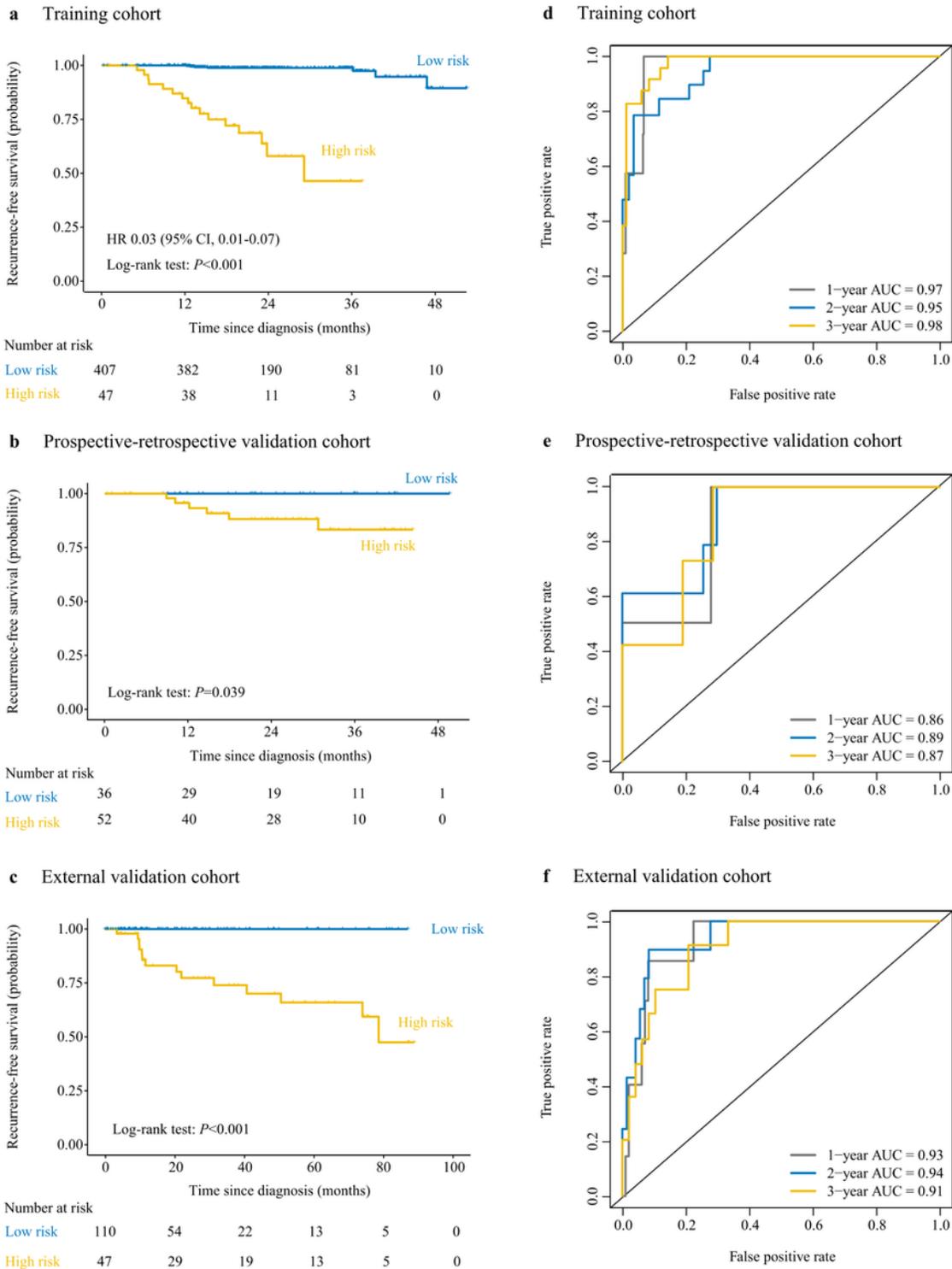
## Figures



**Figure 1**

Study workflow. (a) T1+C, T2WI and DWI-ADC sequences MRI were prepared for (b) radiomic feature extraction from breast intratumoral and peritumoral areas identified and annotated on the MRI scan using 3D Slicer 4.10.2. (c) Key features were selected by the random forest algorithm and then used to (d) construct radiomic signatures and analyse radiogenomics. MRI, Magnetic resonance imaging; T1+C,

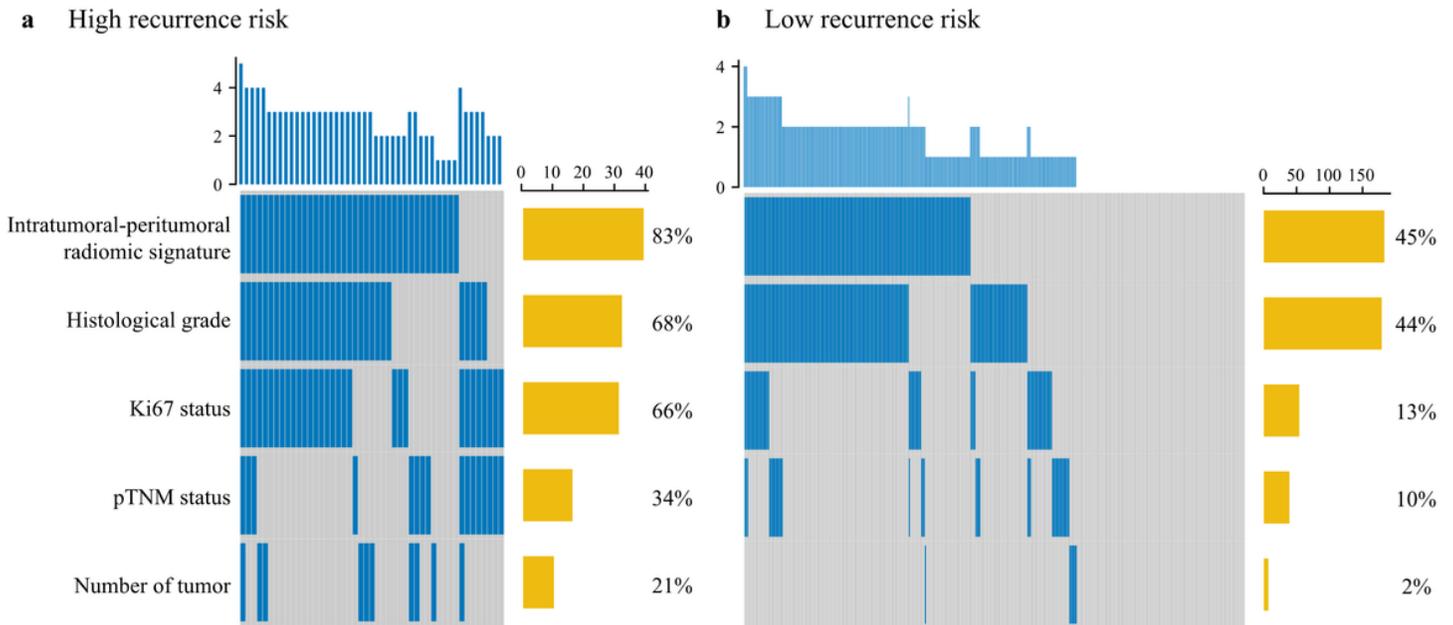
contrast-enhanced T1-weighted imaging; T2WI, T2-weighted imaging; DWI-ADC, diffusion-weighted imaging quantitatively measured apparent diffusion coefficients.



**Figure 2**

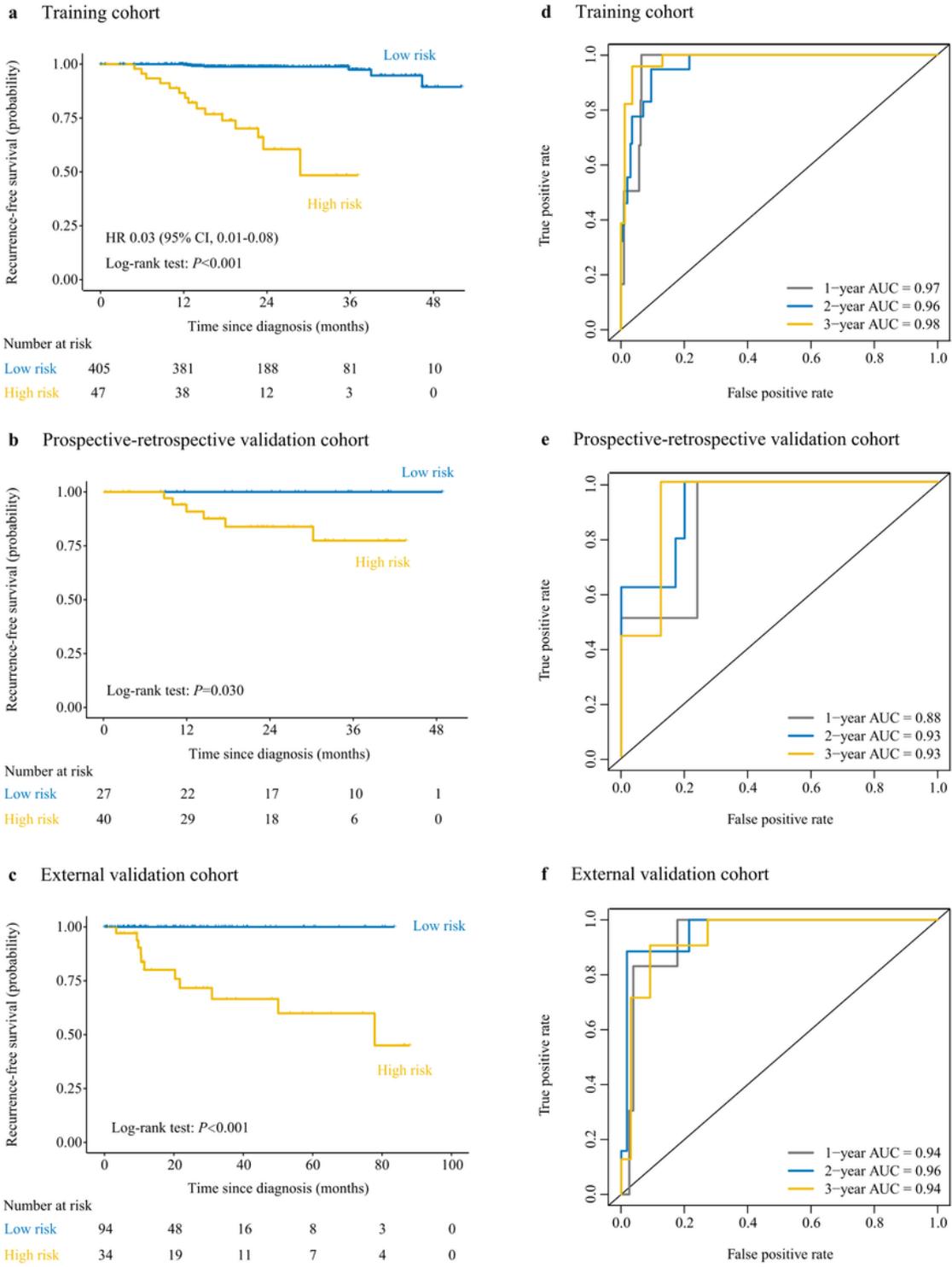
Performance of the intratumoral-peritumoral radiomic signature for predicting the recurrence risk in the training, prospective-retrospective validation and external validation cohorts. Kaplan-Meier curves of RFS according to the intratumoral-peritumoral radiomic signature in the (a) training cohort, (b) prospective-

retrospective validation cohort, and (c) external validation cohort. ROC curves and 1-, 2-, 3-year AUCs were used to assess the prognostic accuracy of the intratumoral-peritumoral radiomic signature in the (d) training cohort, (e) prospective-retrospective validation cohort, and (f) external validation cohort. P values were calculated using the unadjusted log-rank test and hazard ratios were calculated by a univariate Cox regression analysis. RFS, recurrence-free survival; HR, hazard ratio; CI, confidence interval; ROC, receiver operating characteristic; AUC, area under the receiver operating characteristics curve.



**Figure 3**

Intratumoral-peritumoral radiomic signature and clinicopathologic characteristics associated with recurrence risk. pTMN=pathological tumor–node–metastasis stage.



**Figure 4**

Performance of radiomic-clinical signature for predicting the recurrence risk in the training, prospective-retrospective validation and external validation cohorts. Kaplan-Meier curves of RFS according to the radiomic-clinical signature in the (a) training cohort, (b) prospective-retrospective validation cohort, and (c) external validation cohort. ROC curves and 1-, 2-, 3-year AUCs were used to assess the prognostic accuracy of the radiomic-clinical signature in the (d) training cohort, (e) prospective-retrospective

validation cohort, and (f) external validation cohort. P values were calculated using the unadjusted log-rank test and hazard ratios were calculated by a univariate Cox regression analysis. RFS, recurrence-free survival; HR, hazard ratio; CI, confidence interval; ROC, receiver operating characteristic; AUC, area under the receiver operating characteristics curve.

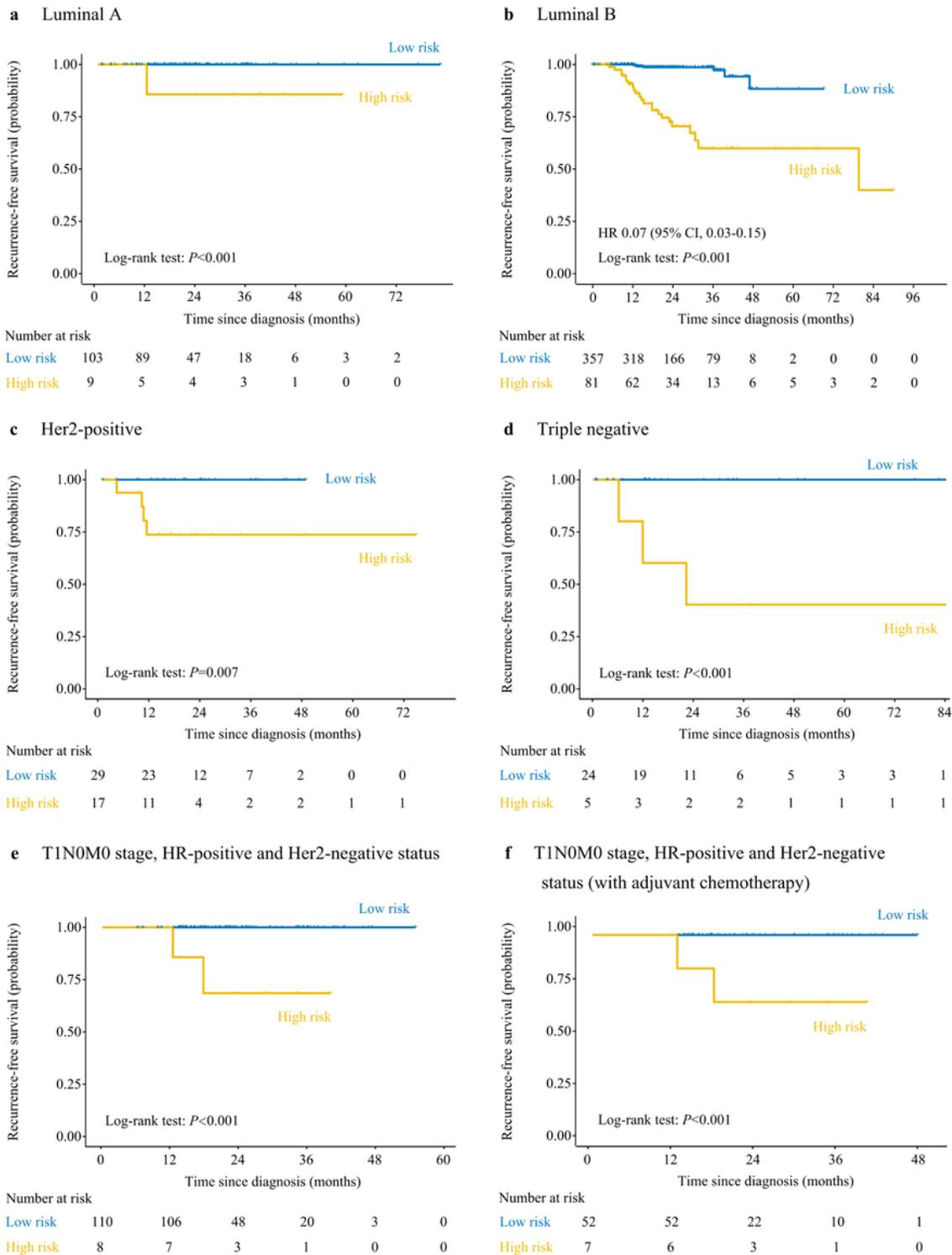


Figure 5

Performance of radiomic-clinical signature for predicting the recurrence risk in different molecular subtype patients. Kaplan-Meier curves of RFS according to the radiomic-clinical signature in the subgroups of (a) Luminal A, (b) Luminal B, (c) Her2-positive, and (d) Triple negative patients. Kaplan-Meier curves of RFS according to radiomic-clinical signature in (e) T1N0M0 stage, HR-positive and Her2-negative status patients and in the subgroups of these (f) patients who received adjuvant chemotherapy. P values were calculated using the unadjusted log-rank test and hazard ratios were calculated by a univariate Cox regression analysis. RFS, recurrence-free survival; HR, hazard ratio; CI, confidence interval; Her2, human epidermal growth factor receptors 2.



eye; GO:0019216, regulation of lipid metabolic process; GO:0071682, endocytic vesicle lumen; GO:0008076, voltage-gated potassium channel complex; GO:0034705, potassium channel complex; GO:0044224, juxtaparanode region of axon; GO:0016614, oxidoreductase activity, acting on CH-OH group of donors; and GO:0016616, oxidoreductase activity, acting on CH-OH group of donors, NAD or NADP as acceptor. P values were calculated using the unadjusted log-rank test and hazard ratios were calculated by a univariate Cox regression analysis. T1+C, contrast-enhanced T1-weighted imaging; T2WI, T2-weighted imaging; LncRNA, long non-coding RNA; HR, hazard ratio; CI, confidence interval; GO, Gene Ontology.

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [Supplementary.docx](#)