

1 (Article)

2 Hormone Receptor-status Prediction in Breast 3 Cancer Using Gene Expression Profiles and Their 4 Macroscopic Landscape

5 Seokhyun Yoon ¹, Hye Sung Won ², Keunsoo Kang ³, Kexin Qiu ¹, Woong June Park ⁴ and
6 Yoon Ho Ko ^{5,*}

7 ¹ Department of Electronics Eng., College of Engineering, Dankook University, Yongin-si, Korea,
8 syoon@dku.edu, qiukexin95@naver.com

9 ² Department of Internal Medicine, Uijeongbu St. Mary's Hospital, College of Medicine, The Catholic
10 University of Korea, Seoul, Korea, woncomet@catholic.ac.kr

11 ³ Department of Microbiology, College of Natural Sciences, Dankook University, Cheonan-si, Korea,
12 kangk1204@dankook.ac.kr

13 ⁴ Department of Molecular Biology, College of Natural Sciences, Dankook University, Cheonan-si, Korea,
14 parkwj@dku.edu

15 ⁵ Department of Internal Medicine, Eunpyeong St. Mary's Hospital, College of Medicine, The Catholic
16 University of Korea, Seoul, Korea, koyoonho@catholic.ac.kr

17 * Correspondence: koyoonho@catholic.ac.kr (Y.H.K.)

18 Received: date; Accepted: date; Published: date

19 **Abstract:** The cost of next-generation sequencing technologies is rapidly declining, making RNA-
20 seq-based gene expression profiling (GEP) an affordable technique for predicting receptor
21 expression status and intrinsic subtypes in breast cancer (BRCA) patients. Based on the expression
22 levels of co-expressed genes, GEP-based receptor-status prediction can classify clinical subtypes
23 more accurately than can immunohistochemistry (IHC). Using data from the cancer genome atlas
24 TCGA BRCA and METABRIC datasets, we identified common predictor genes found in both
25 datasets and performed receptor-status prediction based on these genes. By assessing the survival
26 outcomes of patients classified using GEP- or IHC-based receptor status, we compared the
27 prognostic value of the two methods. We found that GEP-based HR prediction provided higher
28 concordance with the intrinsic subtypes and a stronger association with treatment outcomes than
29 did IHC-based hormone receptor (HR) status. GEP-based prediction improved the identification of
30 patients who could benefit from hormone therapy, even in patients with non-luminal BRCA. We
31 also confirmed that non-matching subgroup classification affected the survival of BRCA patients
32 and that this could be largely overcome by GEP-based receptor-status prediction. In conclusion,
33 GEP-based prediction provides more reliable classification of HR status, improving therapeutic
34 decision making for breast cancer patients.

35 **Keywords:** breast cancer; intrinsic subtype; hormone receptor-status prediction; gene expression
36 profile; LASSO regression

37 1. Introduction

38 Breast cancer (BRCA) is a highly heterogeneous disease that involves several complex molecular
39 networks [1-7]. BRCA can be classified into different subtypes that have distinct clinical behaviors
40 and prognoses and that require different treatment strategies, including targeted therapy and
41 hormone therapy. Therefore, accurate classification of BRCA subtypes is crucial for personalized
42 disease management and for improving patient outcomes [8,9]. Currently, therapeutic decision
43 making in BRCA is based on the expression status of three receptors: estrogen receptor (ER),
44 progesterone receptor (PR), and human epidermal growth factor receptor 2 (HER2) [10, 11]. Although
45 ER, PR, and HER2 status is traditionally determined by immunohistochemistry (IHC), with the
46 advent of high-throughput technologies for gene expression analysis, new molecular subtypes of

47 BRCA have been described. These include luminal A, luminal B, HER2-enriched, basal-like, and
48 normal-like breast tumors [12-14]. The clinical significance of these intrinsic BRCA subtypes has been
49 highlighted by their ability to predict treatment response and prognosis [4-7, 15-21]; hence their use
50 in clinical practice has increased over recent years. Currently, several gene-signature tests based on
51 microarray or quantitative real-time PCR (qRT-PCR) are commercially available [9, 22, 23].

52 The clinicopathological surrogate definitions of the intrinsic BRCA subtypes were endorsed by
53 the 2013 St. Gallen Consensus Recommendations [24]. Luminal A BRCA is hormone receptor (HR)
54 positive, HER2 negative, and expresses low levels of the protein Ki-67. Luminal B BRCA is HR
55 positive and either HER2 positive or HER2 negative, with high levels of Ki-67. The HER2-enriched
56 subtype is HR negative and HER2 positive, and the basal-like subtype is HR negative and HER2
57 negative (triple-negative BRCA) [25 - 27]. Although the expression profiles and clinical features of
58 the four intrinsic BRCA subtypes have been extensively studied in the last few years, discordance has
59 been reported between IHC-based clinical subtypes and intrinsic subtypes in approximately 20–50%
60 of cases [18, 28, 29]. This discordance might be due to intratumoral heterogeneity, the coexistence of
61 cells with different subtypes in the same tumor, as well as measurement inaccuracies in subtype
62 profilers, IHC analysis for ER/PR status, and fluorescence in situ hybridization (FISH) analysis for
63 HER2 status. These inconsistencies could result in administration of the wrong treatment,
64 subsequently leading to poor survival [30]. Therefore, accurate identification of receptor status or the
65 intrinsic BRCA subtype is of high clinical importance.

66 Recently, multi-omics technologies [31], miRNA profiling [32] and principle component
67 analysis-based iterative PAM50 subtyping [33] have helped to improve the accuracy of BRCA
68 subtype classification. However, inconsistencies due to measurement noise remain a challenge in this
69 classification, especially for tumors with receptor expression levels at the boundary between positive
70 and negative [33]. With the development of next-generation sequencing (NGS) technologies, the cost
71 of gene expression profiling (GEP) based on RNA-seq is rapidly decreasing, making it possible to
72 characterize several clinical and molecular features concurrently using RNA-seq-based GEP at a very
73 low cost [34, 35]. Prediction of the intrinsic subtype and receptor status (ER, PR, or HER2) in BRCA
74 using RNA-seq-based GEP would increase the clinical usefulness of RNA-seq technologies in BRCA.
75 In this study, we assessed whether variations in gene expression are reflected in the expression of
76 related genes and whether these changes can be identified by GEP to provide more reliable prediction
77 of the status of the three receptors, thereby improving therapeutic decision making.

78

79 2. Results

80 2.1. Identification of predictor genes

81 In this study, IHC-based characterization of receptor status in BRCA was refined by using co-
82 expressed predictor genes. First, predictor genes were identified; seven genes were selected for ER
83 status prediction, six for PR, and four for HER2 (Table 1). As expected, the *ESR1*, *PGR*, and *ERBB2*
84 genes, which encode the ER, PR, and HER2 proteins, respectively, were among the predictor genes.
85 Model training and receptor-status prediction were then performed using the selected genes. The
86 mismatch rate reported in Table 1 is the percentage of cases in which the IHC-based status differed
87 from the predicted status. Among the predictor genes, *TFF1* and *NAT1* were included in an eighteen-
88 gene set previously reported to predict sensitivity to hormone therapy [36].

Table 1. Summary of mismatch rates and predictor genes for ER, PR, and HER2 status prediction.

Item	Mismatch rate [%]*		Predictor genes
	TCGA	METABRIC	
ER	6.28	6.26	<i>ESR1</i> , <i>AGR3</i> , <i>C1orf64</i> , <i>C4orf7</i> , <i>CLEC3A</i> , <i>SOX11</i> , <i>TFF1</i>
PR	11.43	5.54	<i>PGR</i> , <i>AGR3</i> , <i>ESR1</i> , <i>NAT1</i> , <i>PVALB</i> , <i>S100A7</i>
HER2	11.85	5.17	<i>ERBB2</i> , <i>CPB1</i> , <i>GSTT1</i> , <i>PROM1</i>

* Between the IHC-based and the predicted receptor status

89 2.2. Macroscopic landscape

90 Figure 1 shows uniform manifold approximation and projection (UMAP) plots [38] for receptor
91 status in the TCGA BRCA cohort. Each point represents a sample; the color of the spots corresponds
92 to the (a) subtype (PAM50 class), (b) ER status, (c) PR status, and (d) HER2 status of the sample.
93 Receptor status (ER, PR, or HER2) was provided in the original clinical data based on IHC. The
94 expression of 100 genes selected by LASSO was used to obtain the two-dimensional UMAP projection.
95 The luminal A and B subtypes were mostly HER2- and either ER+ or PR+. However, a small
96 percentage of the luminal A and B subtypes exhibited ER-, PR-, and HER2+. Some patients with
97 HER2-enriched or basal-like subtype BRCA also showed some level of discordance, as some HER2-
98 enriched and basal-like subtype samples were ER+ or PR+. Although most HER2+ and HER2-
99 enriched subtype samples overlapped, some HER2-enriched subtype samples were found to be
100 HER2- BRCA that exhibited basal-like subtype features. As only eight patients exhibited normal-like
101 subtype BRCA in the TCGA dataset, they were not considered in our analyses.

102 On the other hand, the HER2-enriched subtype samples were ER+ and/or PR+, representing a
103 luminal subtype. The UMAP plot of the METABRIC dataset revealed a similar macroscopic
104 landscape (Supplementary Figure 1). Considering that the distance between samples (points) in the
105 UMAP projection is only an approximation of the relative distance in their gene expression profiles
106 and that the receptor status was not clearly defined for all samples, Figure 1 implies that IHC/FISH-
107 based characterization of receptor status might result in inaccuracies in BRCA subtype classification.

108 Figure 2 shows the same UMAP plot based on the predicted values obtained by the linear
109 classifiers. Compared with IHC-based receptor-status characterization, the predicted status was more
110 consistent with the intrinsic BRCA subtype classification, especially for the basal-like and luminal
111 subtypes. Most of the luminal subtypes were ER+ and PR+, and the numbers of ER+ or PR+ samples
112 in the basal-like subtype were much smaller than after IHC-based status characterization. The UMAP
113 plot for the METABRIC dataset based on the predicted receptor status (Supplemental Figure 2) led
114 to the same conclusions, except for PR status, which was not IHC-based in the METABRIC dataset.

115 2.3. GEP-based receptor-status prediction is reliable for the luminal and basal-like subtypes

116 To quantify discordance between the intrinsic subtype and the clinical subtype defined by HR
117 and HER2 status, for each intrinsic subtype, we compared the numbers of positive and negative
118 instances of HR and HER2 status based on IHC with the numbers obtained using GEP-based
119 prediction in the TCGA and METABRIC datasets (Table 2). The rates of discordance for the basal-
120 like, luminal A, and luminal B subtypes were lower using GEP-based prediction than using IHC-
121 based status characterization. Specifically, most samples of the luminal A and B subtypes were
122 characterized as HR+ by GEP-based prediction (except for two samples in the TCGA BRCA cohort),
123 while some luminal A and luminal B BRCA samples were characterized as HR- based on IHC. In
124 BRCA patients with the basal-like subtype, a smaller percentage of tumors was determined to be HR+
125 using GEP-based prediction (10% in TCGA and 13% in METABRIC) than when using IHC-based
126 characterization (17% in TCGA and 20% in METABRIC).

127 On the other hand, considerable discordance was observed in the receptor status of HER2-
128 enriched subtype BRCA patients using both IHC-based characterization and GEP-based prediction.
129 Only 37% and 23% of patients with HER2-enriched subtype BRCA were HR-/HER2+ in the
130 METABRIC and TCGA datasets, respectively. Furthermore, 17% and 18% of tumors were triple
131 negative, and 25% and 9% were luminal-like (HR+ and HER2-) in the METABRIC and TCGA datasets,
132 respectively. Similar findings were obtained for IHC-based characterization of HR and HER2 status.

133 In summary, GEP-based prediction was more concordant with the typical receptor-status pattern
134 of the intrinsic subtypes of patients with the basal-like, luminal A, and luminal B subtypes. However,
135 this does not necessarily mean that receptor-status prediction based on GEP is more accurate than
136 IHC-based characterization. The only way to verify the accuracy of the status predictions is to assess
137 the differences in clinical outcomes among the different clinical subtypes defined by the status of the
138 three receptors.

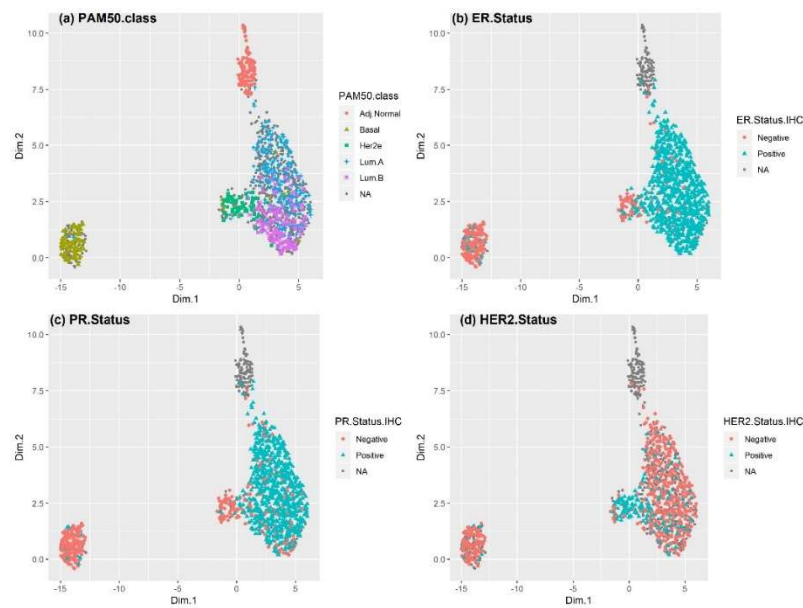


Figure 1. UMAP plot showing the receptor status in the TCGA BRCA cohort. The tumor subtype, as well as the status of ER, PR, and HER2, were based on the available clinical data. Gray points are samples with no available clinical information. A small percentage of the luminal A and B subtypes were ER-/PR- and HER2+. Such discordances were also observed in some BRCA patients with the HER2-enriched and basal-like subtypes. Although most HER2+ and HER2-enriched subtype samples overlapped, some HER2-enriched subtype samples were found to be HER2- BRCA and to exhibit basal-like subtype features. Some samples were ER+ and/or PR+, representing a luminal subtype.

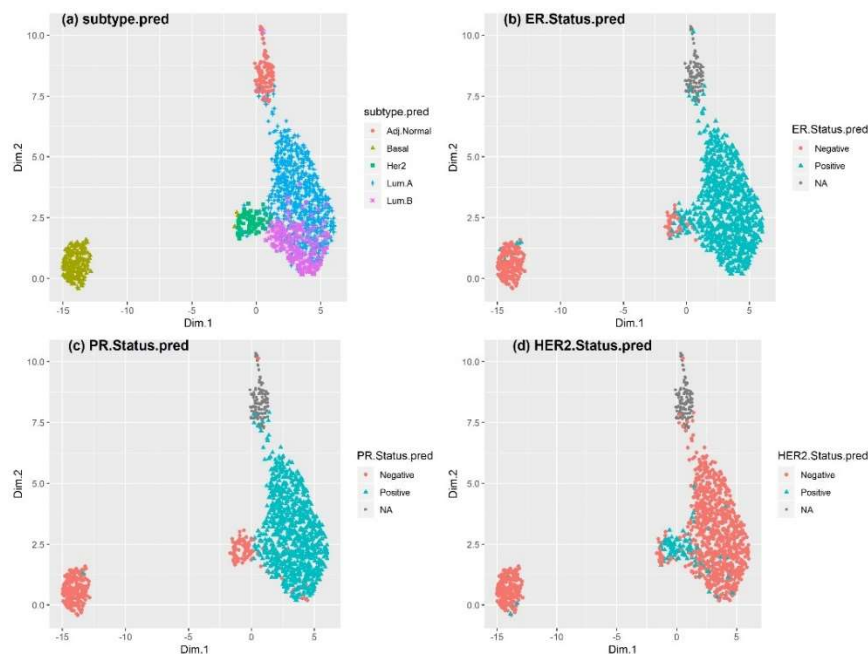


Figure 2. UMAP plot showing GEP-based receptor status in the TCGA BRCA cohort. GEP-based prediction was used to determine the subtype, as well as the status of ER, PR, and HER2. Compared to the case with IHC-based approaches, the predicted status of ER, PR, and HER2 was mostly in accordance with the corresponding pattern of receptor status for basal-like, luminal A, and luminal B. In contrast, HER2-enriched subtype tumors were highly heterogeneous.

Table 2. HR and HER2 status for each intrinsic subtype as determined by (a) IHC- and (b) GEP-based prediction. Patients with no available IHC-based receptor status were excluded.

Dataset	Subtype	(a) IHC-based characterization		(b) GEP-based prediction	
		HR+/-	HER2+/-	HR+/-	HER2+/-
TCGA	Luminal A	222 / 4	24 / 130	229 / 2	4 / 227
	Luminal B	126 / 1	22 / 69	127 / 0	8 / 119
	Basal-like	16 / 78	6 / 59	10 / 87	2 / 95
	HER2-enriched	32 / 24	40 / 10	44 / 14	39 / 19
METABRIC	Luminal A	680 / 6	19 / 283	696 / 0	19 / 677
	Luminal B	465 / 1	23 / 171	474 / 0	29 / 445
	Basal-like	61 / 243	14 / 118	40 / 268	24 / 284
	HER2-enriched	119 / 111	50 / 34	125 / 111	119 / 117
	Normal-like	161 / 21	11 / 51	165 / 19	11 / 173

139 *2.4. GEP-based receptor-status prediction is reliable for the luminal and basal-like subtypes*

140 To verify the accuracy of the receptor-status predictions, survival outcomes for various
 141 combinations of HR and HER2 status were compared. The significance of the prognostic value of
 142 the predicted and IHC-characterized HR and HER2 status was compared. Separate survival
 143 analyses were performed in the following four patient groups:
 144

145 (a) HR+ (either ER+ or PR+) group: This group benefited from hormone therapy. According to
 146 the stage and clinical characteristics, these patients often received a combination of hormone
 147 therapy and chemotherapy. For survival analysis, the patients in this group were stratified based
 148 on administration of hormone therapy.

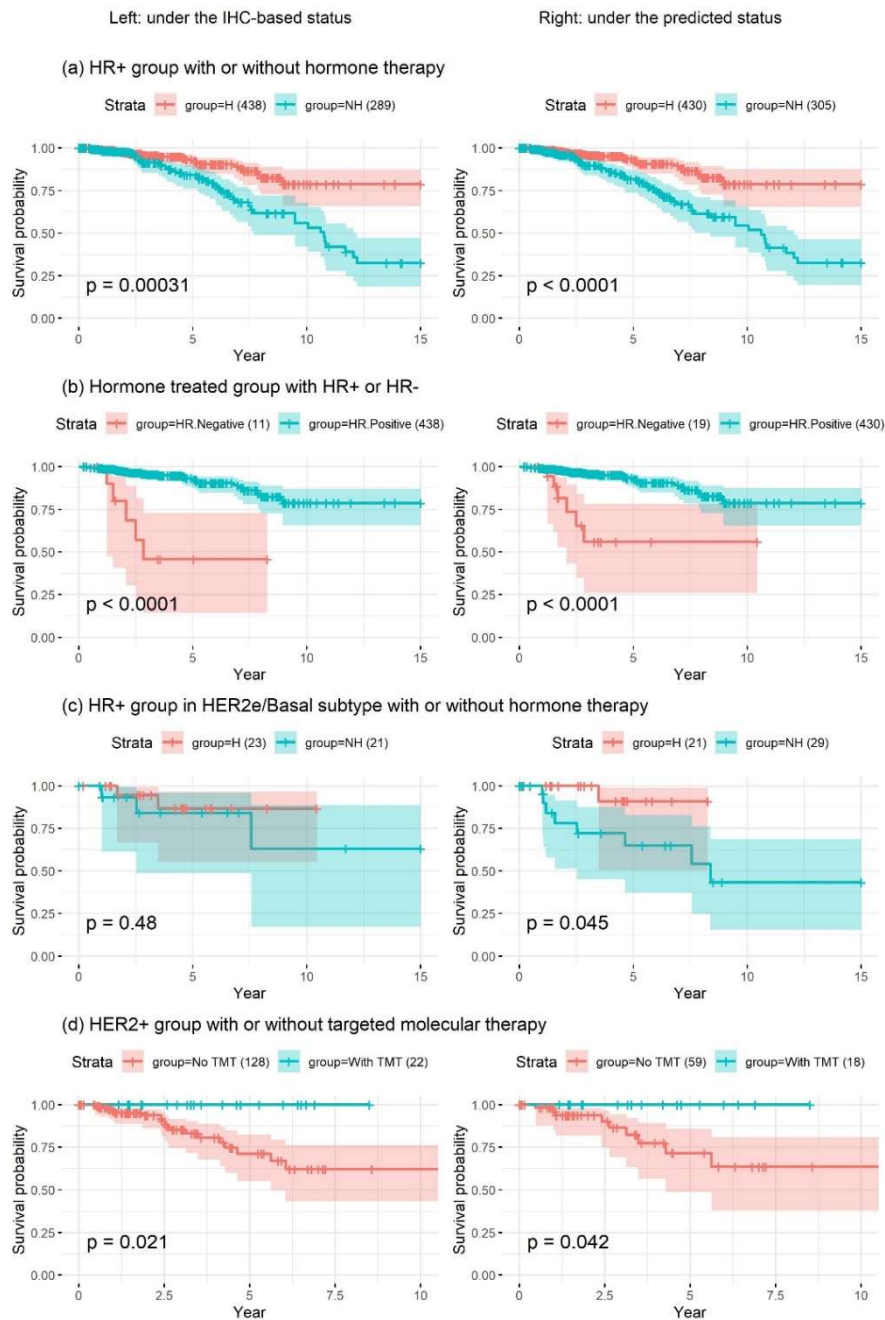
149 (b) Hormone therapy group: To confirm the benefit of hormone therapy for HR+ patients, only
 150 those who received hormone therapy, with or without chemotherapy, were selected, and the
 151 survival of HR+ patients was compared to that of HR- patients.

152 (c) HR+/non-luminal subtype group: As shown in Table 2, there were small percentages of HR+
 153 patients among patients with the HER2-enriched and basal-like subtypes. Hence, we assessed
 154 whether BRCA patients with the HR+ non-luminal subtype benefited from hormone therapy.

155 (d) HER2+ group: BRCA patients with the HER2+ subtype benefited from anti-HER2 targeted
 156 molecular therapy (TMT). We assessed the survival of HER2+ BRCA patients based on TMT. As
 157 no information regarding TMT was available in the METABRIC dataset, this analysis was
 158 performed only for the TCGA BRCA cohort.

159 Among patients in the TCGA BRCA cohort, GEP-based receptor-status prediction provided a
 160 higher hazard ratio with higher significance in HR- patients (a), implying that GEP-based receptor-
 161 status prediction had higher prognostic value than traditional IHC-based HR status characterization.
 162 On the other hand, in the hormone-therapy group (b), IHC-based receptor-status characterization
 163 was found to be more accurate than GEP-based receptor-status prediction. However, the numbers of
 164 samples in the test group (HR- patients) were only 11 and 19 for receptor-status characterization
 165 based on IHC and GEP, respectively. Among patients with HR+ non-luminal subtype BRCA (c), IHC-
 166 based receptor status had no significant prognostic value, in contrast to GEP-based receptor-status
 167 prediction. This finding highlighted that HR+ BRCA patients benefited from hormone therapy, even
 168 if they were diagnosed with non-luminal subtype tumors. Among HER2+ patients (d), IHC-based
 169 receptor-status characterization exhibited higher prognostic value when considering only the p-value.
 170 However, the numbers of patients with IHC-based receptor-status data in the test group (HER2+
 171 patients with TMT) were only 22 and 18 based on IHC and GEP, respectively, and all patients that
 172 received TMT survived; hence, the hazard ratio could not be precisely determined (Figure 3 and Table

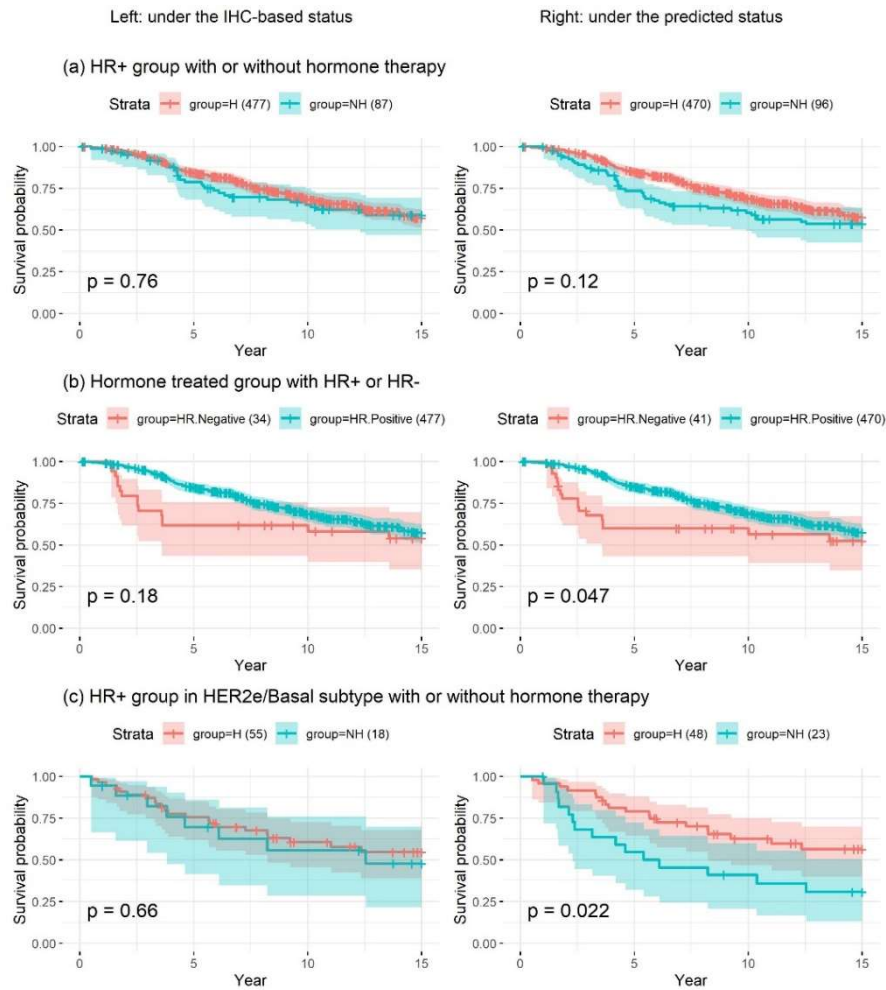
173 3). Survival analyses in the METABRIC cohort (excluding patients with a pathological stage of I)
174 showed similar findings, implying that GEP-based receptor-status prediction had higher prognostic
175 significance in terms of patient survival compared to traditional IHC-based receptor-status
176 characterization (Figure 4 and Table 3).
177



178
179
180
181
182
183
184
185
186
187
188
189

Figure 3. Kaplan–Meier survival analysis of patients from the TCGA dataset using IHC-based (left panel) or GEP-based (right panel) receptor status. Patients were stratified to those who received hormone therapy (H) and those who did not (NH). (a) GEP-based receptor status prediction had higher prognostic significance in terms of patient survival compared to IHC-based HR status. (b) IHC-based receptor-status characterization was found to be more accurate than GEP-based receptor-status prediction. However, the numbers of samples in the test group (HR– patients) for receptor-status characterization based on IHC and GEP were only 11 and 19, respectively. (c) IHC-based receptor status had no significant prognostic value, in contrast to GEP-based receptor-status prediction. (d) The statistical significance of IHC-based receptor-status characterization indicated higher prognostic value. However, the numbers of patients with IHC-based receptor-status data in the test group (HER2+ patients with TMT) were only 22 under IHC and 18 under GEP, and all patients who received TMT survived; hence, the hazard ratio could not be precisely determined.

190



191
192
193
194
195
196

Figure 4. Kaplan–Meier survival analysis in patients of the METABRIC dataset with a pathological stage of II or III (excluding pathological stage I). The analysis was performed using IHC-based receptor status (left panel) or GEP-based receptor status (right panel). GEP-based receptor-status prediction had higher prognostic significance in terms of patient survival compared to traditional IHC-based receptor-status characterization.

Table 3. A summary of the hazard ratios and associated statistical significance obtained from survival analyses using IHC-based receptor status (IHC) or the predicted status (pred.). For the survival analysis, data from the TCGA and METABRIC datasets were used.

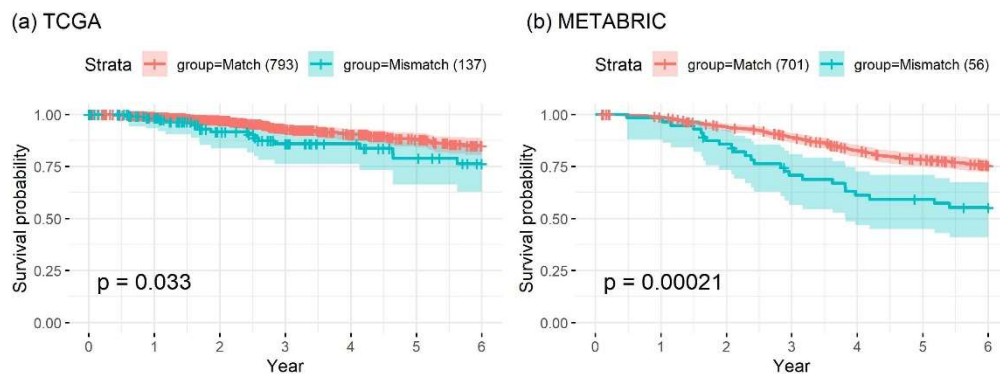
Patient group	Conditions compared	# of samples		p-value		Hazard ratio	
		IHC	Pred.	IHC	Pred.	IHC	Pred.
TCGA							
(a) HR+	H vs. NH	727 (438, 289)	735 (430, 305)	0.00031	2.11·10 ⁻⁰⁵	0.89	1.0
(b) Hormone therapy	HR+ vs. HR-	449 (438, 11)	449 (430, 19)	3.15·10 ⁻⁰⁸	3.38·10 ⁻⁰⁷	2.23	2.0
(c) HR+ in HER2e/Basal	H vs. NH	44 (23, 21)	50 (21, 29)	0.48	0.045	0.65	1.88
(d) HER2+	T vs. NT	150 (22, 128)	77 (18, 59)	0.021	0.042	19.4	19.6
METABRIC							
(e) HR+	H vs. NH	564 (477, 87)	566 (470, 96)	0.76	0.12	0.06	0.28
(f) Hormone therapy	HR+ vs. HR-	511 (477, 34)	511 (470, 41)	0.18	0.047	0.36	0.49
(g) HR+ in HER2e/Basal	H vs. NH	73 (55, 18)	71 (48, 23)	0.66	0.022	0.18	0.77

HR: hormone receptor; H: with hormone therapy regardless of chemotherapy; NH: without hormone therapy; T: with targeted molecular therapy regardless of hormone/chemotherapy; NT: without targeted molecular therapy

197

198 2.5. Patients with non-matching receptor status had significantly worse survival

199 The type of adjuvant therapy is based mainly on the status of the three receptors. Hence, accurate
200 characterization of receptor status is of high clinical importance. As shown in Figure 5, patients with
201 matching receptor status had longer overall survival (OS) compared to those with non-matching
202 status (hazard ratios 0.6 and 0.79 for the TCGA BRCA and METABRIC cohorts, respectively).
203 Assuming higher accuracy for GEP-based receptor-status prediction, these results highlight the
204 impact of inappropriate treatment due to errors in receptor-status characterization. Although it is
205 unlikely that GEP-based receptor-status prediction is 100% accurate, it can identify patients who can
206 benefit from hormone therapy more reliably than the traditional IHC-based method.
207



208 **Figure 5.** Kaplan–Meier survival analysis of patients in the (a) TCGA BRCA cohort and (b) METABRIC dataset
209 with matching and non-matching receptor status. The hazard ratios of patients with non-matching status were
210 0.6 for the TCGA BRCA cohort and 0.79 for the METABRIC dataset.
211

212 **3. Discussion**

213 IHC-based assessment of the expression of a specific protein is undoubtedly an important tool for
214 detecting biomarkers in clinical practice. However, this procedure entails severe limitations,
215 including variations in the IHC procedure that can influence the results and their interpretation. As
216 an alternative, biomarker characterization could be performed at the mRNA level; unfortunately,
217 high mRNA levels do not necessarily translate into high levels of the corresponding protein.
218 Additionally, characterization based solely on the expression levels of a single gene or protein
219 inevitably entails the risk of noise. To overcome these limitations, we considered the potential use of
220 GEP-based receptor-status prediction for molecular characterization of BRCA subtypes. Changes in
221 the expression of a gene should be reflected in those of co-expressed genes; therefore, prediction
222 based on the expression of correlated genes may outperform molecular characterization based on a
223 single gene.

224 In the era of biomarker-assisted targeted therapy, the method used to assess biomarker expression
225 is crucial, as it can improve the prognosis for patients with BRCA and other malignancies. Several
226 challenges remain to be overcome in biomarker-assisted targeted therapies, such as IHC-determined
227 borderline HR-positivity, equivocal HER2 amplification, and discordance between IHC-based
228 subtypes and intrinsic subtypes. Previous studies have shown significant discordance between
229 clinical subtypes and intrinsic subtypes, which affects the prognosis of BRCA patients. Kim et al.
230 reported that discrepancies between the IHC-based subtype and the intrinsic subtype were associated
231 with poor survival, highlighting the limitations of current IHC-based classification methods [30].
232 Consistent with previous results, we confirmed the poor survival of patients with non-matching
233 subgroup classifications in both the TCGA and METABRIC datasets. These results emphasize the
234 clinical importance of establishing more accurate classification methods. Herein, we evaluated the
235 concordance between the intrinsic subtype and the predicted status of ER, PR, and HER2 using GEP.
236 We found a higher concordance rate between the intrinsic subtype and GEP-based receptor-status
237 prediction compared to receptor status as characterized by IHC. This was consistent in all BRCA

238 subtypes except for the HER2-enriched subtype. These findings imply that GEP-based HR status
239 prediction could be a promising alternative approach to IHC.

240 Both IHC-based receptor-status characterization and GEP-based status prediction resulted in
241 considerable discordance between HER2-positivity and the HER2-enriched subtype. Although the
242 HER2-enriched subtype is the predominant type of HER2-positive BRCA, three other subtypes exist.
243 A recent study analyzing data from four prospective neoadjuvant trials reported that the percentages
244 of the luminal A, luminal B, HER2-enriched, and basal-like subtypes among HER2-positive BRCA
245 patients were 24%, 20%, 47%, and 9%, respectively [39]. This finding may be partly explained by high
246 intratumoral heterogeneity. Previous genomic analyses have revealed that HER2-positive BRCA is
247 extremely clinically and biologically heterogeneous [40, 41]. The HER2-enriched subtype is also
248 highly heterogeneous, rendering IHC/FISH- and PAM50-based subtyping challenging.

249 Furthermore, the HER2-enriched subtype can have a distinctive transcriptional landscape
250 independent of HER2 amplification. Analyses in TCGA showed that the HER2-enriched subtype was
251 characterized by the highest number of DNA mutations, including in TP53 and PIK3CA [26].
252 Recently, Daemen A et al. performed genomic and transcriptomic profiling of HER2-enriched tumors;
253 they concluded that HER2 was not a cancer subtype but rather a pan-cancer phenomenon and that
254 HER2-positive tumors are hormonally driven [42]. Even though further stratification of HER2-
255 enriched BRCA might be beneficial, it might be difficult to achieve further characterization based on
256 GEP. To overcome the limitations of macroscopic GEP, different microscopic prediction approaches
257 could be used, including precise reconstruction of transcriptome data and use of single-cell RNA-seq.
258 These approaches might achieve more in-depth characterization of the molecular subtypes.

259 To investigate the clinical relevance of GEP-based prediction of ER, PR, and HER2 receptor status,
260 we performed survival analysis of HR+ patients who did or did not receive hormone therapy, as well
261 as of HR+ and HR- patients treated with hormone therapy. GEP-based receptor-status prediction
262 showed a more significant association between treatment outcomes and HR status compared to IHC-
263 based receptor-status characterization. Of note, some benefit was achieved from hormone therapy by
264 patients who were identified as HR+ non-luminal BRCA using GEP-based prediction, in contrast to
265 when IHC-based HR status characterization was performed. These results imply that GEP-based
266 receptor-status prediction can better identify patients who can benefit from hormone therapy, even
267 in patients with non-luminal subtype BRCA. Some studies have shown that adjuvant or palliative
268 hormone therapy is less effective in patients with HR+ BRCA of the non-luminal subtype [43, 44].
269 However, there is limited evidence regarding which HR+ non-luminal BRCA patients will benefit
270 from hormone therapy. Future studies are needed to determine whether GEP-based receptor-status
271 prediction can address these clinically important questions. In contrast to the HR status, we did not
272 observe improvement in HER2 status prediction; this may be attributed partially to the small number
273 of patients who received targeted molecular therapy for HER2.

274 4. Materials and Methods

275 The workflow of this study is shown in Figure 6. Our analyses were performed in three steps. First,
276 we identified common predictor genes from two different gene-expression datasets. Second, we
277 predicted ER, PR, and HER2 status based on the shared predictor genes. Finally, we compared
278 survival outcomes according to IHC-based and GEP-based predictions of receptor status.

279 4.1. Datasets

280 For this study, we used BRCA patients' gene-expression-profile and clinical data acquired from
281 the cancer genome atlas (TCGA) [<http://firebrowse.org/>] and the Molecular Taxonomy of Breast
282 Cancer International Consortium (METABRIC) databases [<https://www.cbioportal.org/>] [27]. Both
283 datasets include information on the history of adjuvant treatment, which was a critical element in the
284 survival analyses performed in this study. A summary of the data contained in the two datasets is
285 shown in Table 4.

286

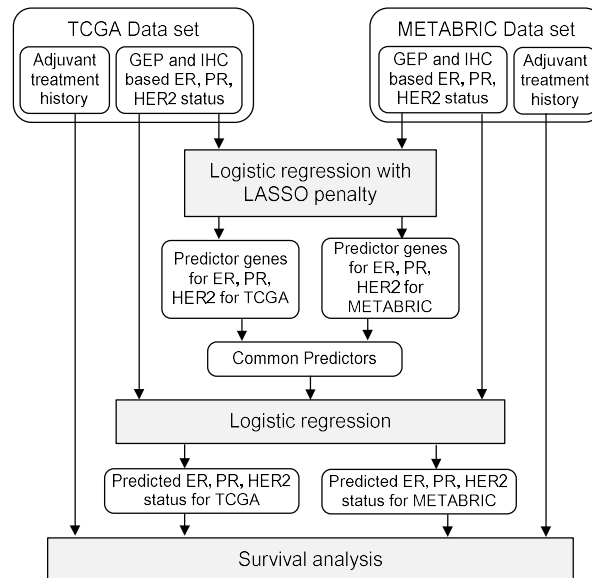


Figure 6. Workflow of gene selection, model training, receptor-status prediction, and survival analysis.

287
288
289
290
291
292
293
294
295
296
297
298

The TCGA BRCA dataset contained data from tumor samples ($n = 1,092$ patients) and adjacent normal tissues ($n = 112$ patients). The METABRIC dataset contained data from 2,506 tumor samples, including GEP data from 1,904 patients. The TCGA and METABRIC datasets also contained clinical data, including ER, PR, and HER2 status, as well as histories of surgery, radiation-therapy, and drug treatments; however, clinical data were not available for all of the patients. Information regarding the tumor subtype was available for some samples in the TCGA BRCA dataset; PAM50 mRNA profile information was available for 523 of 1,092 patients [26]. To ensure consistency between the two datasets, information on ER and HER2 status as determined by IHC was used for patients in the METABRIC dataset. Non-IHC-based PR status was used for the METABRIC cohort because the PR status was not assessed by IHC in these patients.

Table 4. A summary of data availability in the TCGA BRCA cohort and METABRIC dataset.

Item	TCGA BRCA cohort	METABRIC	Comment
Gene expression profile	Yes	Yes	
PAM50-based subtype	Yes (partially)	Yes	
ER status	Yes (IHC)	Yes (IHC, non-IHC)	Used IHC-based status
PR status	Yes (IHC)	Yes (non-IHC)	Used for receptor status
HER2 status	Yes (IHC)	Yes (IHC, non-IHC)	Used IHC-based status
RPPA measurements	Yes	No	
Types of drug treatment	Chemo, hormone and targeted molecular therapy	Chemo and hormone therapy	Used for survival analysis
Age at initial diagnosis	Yes	Yes	Used for sample selection
Pathological stage	Yes	Yes	Used for sample selection

299 4.2. Prediction model and gene selection

300 Based on GEP and the status of the three receptors, logistic regression with LASSO penalty was
301 performed in a supervised mode to identify predictor genes for each of the two datasets. This analysis
302 was performed using the R package glmnet [45-47]. In the TCGA BRCA dataset, the expression levels
303 of 17,202 genes were log₂-transformed and normalized. In the METABRIC dataset, already
304 normalized mRNA expression data were used. To identify the common predictor genes and
305 minimize overfitting-related errors, LASSO penalty weights were selected for a set of predefined
306 genes (e.g., 10, 20, 40, and 60), and for each number, the penalty weight that led to the closest number

307 of selected genes was chosen. This approach was conducted separately for the TCGA and METABRIC
 308 datasets. Common predictor genes between TCGA and METABRIC were then identified to avoid
 309 dataset-related dependencies. After inspecting the overall number of shared genes, 40 genes were
 310 selected; these contained 7, 6, and 4 common predictor genes for ER, PR, and HER2, respectively, as
 311 summarized in Table 1. Subsequently, logistic regression was performed again to train the models
 312 for ER, PR, and HER2 status prediction for both TCGA and METABRIC. The mismatch rate was
 313 obtained by fivefold cross-validation.

314 Pairwise correlations of gene-expression levels between the selected genes are shown in
 315 Supplementary Figures 4, 5, and 6. Of note, PR predictor genes included *ESR1* and *AGR3*, which were
 316 also ER predictor genes. Furthermore, among the four HER2 predictor genes, *CPB1*, *GSTT1*, and
 317 *PROM1* showed only small correlations with *ERBB2*, implying that HER2 status prediction was
 318 determined predominantly by *ERBB2*.

319 4.3. Survival analysis for accuracy evaluation and sample selection

320 The survival analyses were performed for various group/condition pairs; significance (p-value)
 321 was used as an accuracy criterion. Cox's proportional hazard model was used to determine overall
 322 survival [48]; the analysis was repeated using the IHC-based status and the predicted status. For the
 323 survival analysis based on IHC-based receptor status, we used those samples for which IHC-based
 324 receptor status was available. For the survival analysis based on predicted-receptor status, we used
 325 the same set of samples without considering discrepancies between the predicted status and the IHC-
 326 based status. As shown in Table 1, in 5–12% of cases, the predicted status differed from the IHC-
 327 based status.

328 Additionally, for the survival analyses, patients were selected according to the following criteria:
 329 (1) pathological cancer stage I, II, or III and (2) age <80 years at initial diagnosis. Subsequently,
 330 patients were stratified according to adjuvant drug treatments. The characteristics of the patients
 331 included in the survival analyses are summarized in Table 5.
 332

Table 5. A summary of the samples available in the TCGA and METABRIC datasets.

Variable	Conditions	The number of available samples	
		In TCGA	In METABRIC
Age	≤80 years	1,039	1,783
Pathologic stage:	I	170	464
	II	598	736
	III	232	105
Therapy applied:	Chemotherapy	578	393
	Hormone therapy	495	1,084
	Both chemo- and hormone therapy	324	181
	Targeted molecular therapy	30	NA
ER status:	Positive	760	1,339
	Negative	230	418
	NA	2	0
PR status:	Positive	663	946
	Negative	324	837
	NA	4	0
HER2 status:	Positive	159	114
	Negative	524	647
	NA	182	27

* For ER, PR, and HER2 status; 'indeterminate' and 'equivocal' were reported as NA.

333 5. Conclusions

334 Therapeutic decision making in BRCA is heavily based on the clinical subtype defined by HR and
335 HER2 expression status. NGS-based approaches could allow more accurate characterization of the
336 various molecular and clinical features of BRCA. GEP-based receptor-status prediction could provide
337 a better understanding of BRCA pathology and guide physicians in decision making. To improve the
338 performance of GEP-based prediction models, data from larger cohorts are required for
339 standardization of the procedure. In addition, a more comprehensive analysis of receptor status
340 should be performed to identify the characteristics that affect the positivity or negativity of the status
341 of the three receptors, as well as the mechanisms responsible for the discordance between intrinsic
342 subtype and clinical subtype.

343 **Supplementary Materials:** The following materials contain some of TCGA and METABRIC clinical data and the
344 new predictions on the 3-receptor status, which were used for the survival analyses in this work.

- 345 1. TCGA_BRAC_clinical_data_n_pred_status.csv:
- 346 2. METABRIC_clinical_data_n_pred_status.csv

347 **Author Contributions:** Conceptualization, S.Y., H.S.W., K.K., W.J.P. and Y.H.K.; methodology, S.Y., K.K. and
348 W.J.P., software, validation and formal analysis, S.Y., and K.Q.; investigation and data curation, H.S.W. and
349 Y.H.K.; funding acquisition, S.Y.; writing—original draft preparation and visualization, S.Y., H.S.W and K.K.,
350 writing—review and editing, S.Y., H.S.W., K.K., W.J.P. and Y.H.K.; All authors have read and agreed to the
351 published version of the manuscript.

352 **Funding:** This work was supported by Basic Science Research Program through the National Research
353 Foundation of Korea (NRF) grant funded by the Ministry of Education, Science and Technology (NRF-
354 2016R1D1A1B03933651 to S.Y) and by industry-academic cooperation research fund of the Catholic University
355 of Korea (5-2019-D0189-00002 to Y.H.K).

356 **Acknowledgments:** The results here are in part based upon data generated by the TCGA Research Network:
357 <https://www.cancer.gov/tcga>.

358 **Conflicts of Interest:** The authors declare that they have no competing interest.

359 References

- 360 1. Spitale A, Mazzola P, Soldini D, Mazzucchelli L and Bordoni A. Breast cancer classification according to
361 immunohistochemical markers: clinicopathologic features and short-term survival analysis in a
362 population-based study from the South of Switzerland. *Ann Oncol* **2009**, 20: 628-635.
- 363 2. Tang P, Wang J and Bourne P. Molecular classifications of breast carcinoma with similar terminology and
364 different definitions: are they the same? *Hum Pathol* **2008**, 39: 506-513.
- 365 3. Desmedt C, Sotiriou C and Piccart-Gebhart MJ. Development and validation of gene expression profile
366 signatures in early-stage breast cancer. *Cancer Invest* **2009**, 27: 1-10.
- 367 4. Iwamoto T and Pusztai L. Predicting prognosis of breast cancer with gene signatures: are we lost in a sea
368 of data? *Genome Med* **2010**, 2: 81.
- 369 5. Reis-Filho JS, Weigelt B, Fumagalli D and Sotiriou C. Molecular profiling: moving away from tumor
370 philately. *Sci Transl Med* **2010**, 2: 47ps43.
- 371 6. Sotiriou C and Pusztai L. Gene-expression signatures in breast cancer. *N Engl J Med* **2009**, 360: 790-800.
- 372 7. Weigelt B, Baehner FL and Reis-Filho JS. The contribution of gene expression profiling to breast cancer
373 classification, prognostication and prediction: a retrospective of the last decade. *J Pathol* **2010**, 220: 263-280.
- 374 8. Blows FM, Driver KE, Schmidt MK, *et al*. Subtyping of breast cancer by immunohistochemistry to
375 investigate a relationship between subtype and short and long term survival: a collaborative analysis of
376 data for 10,159 cases from 12 studies. *PLoS Med* **2010**, 7: e1000279.
- 377 9. Dai X, Li T, Bai Z, Yang Y, Liu X, Zhan J and Shi B, Breast Cancer Intrinsic Subtype Classification, Clinical
378 Use and Future Trends. *Am J. Cancer Res.* **2015**, 5(10): 2929-2943.
- 379 10. Vallejos CS, Gomez HL, Cruz WR, *et al*. Breast Cancer Classification According to Immunohistochemistry
380 Markers: Subtypes and Association with Clinicopathologic Variables in a Peruvian Hospital Database.
381 *Clinical Breast Cancer* **2010**, 10: 294-300.
- 382 11. Cheang MC, Chia SK, Voduc D, *et al*. Ki67 index, HER2 status, and prognosis of patients with luminal B
383 breast cancer. *J Natl Cancer Inst* **2009**, 101: 736-50.

- 384 12. Sørlie T, Perou CM, Tibshirani R, *et al.* Gene expression patterns of breast carcinomas distinguish tumor
385 subclasses with clinical implications. *Proc Natl Acad Sci USA* **2001**, 98: 10869-10874.
- 386 13. Perou, C., Sørlie, T., Eisen, M. *et al.* Molecular portraits of human breast tumours. *Nature* **2000**, 406, 747-752.
- 387 14. Prat A, Parker JS, Fan C, Perou CM. PAM50 assay and the three-gene model for identifying the major and
388 clinically relevant molecular subtypes of breast cancer. *Breast Cancer Res Treat.* **2012**, 135:301-6.
- 389 15. Puzstai L, Broglio K, Andre F, Symmans WF, Hess KR and Hortobagyi GN. Effect of molecular disease
390 subsets on disease-free survival in randomized adjuvant chemotherapy trials for estrogen receptor-positive
391 breast cancer. *J Clin Oncol* **2008**, 26: 4679-4683.
- 392 16. Paik S, Shak S, Tang G, *et al.* A multigene assay to predict recurrence of tamoxifen-treated, node-negative
393 breast cancer. *N Engl J Med* **2004**, 351: 2817-2826.
- 394 17. Prat A, Perou CM. Deconstructing the molecular portraits of breast cancer. *Mol Oncol.* **2011**, 5:5-23.
- 395 18. Parker JS, Mullins M, Cheang MC, *et al.* Supervised risk predictor of breast cancer based on intrinsic
396 subtypes. *J Clin Oncol* **2009**, 27: 1160- 1167.
- 397 19. Sotiriou C, Neo SY, McShane LM, *et al.* Breast cancer classification and prognosis based on gene expression
398 profiles from a population-based study. *Proc Natl Acad Sci USA* **2003**, 100: 10393-10398.
- 399 20. Fan C, Oh DS, Wessels L, *et al.* Concordance among gene-expression-based predictors for breast cancer. *N*
400 *Engl J Med* **2006**, 355: 560-569.
- 401 21. Lehmann BD, Bauer JA, Chen X, Sanders ME, Chakravarthy AB, Shyr Y and Pietsenpol JA. Identification of
402 human triple-negative breast cancer subtypes and preclinical models for selection of targeted therapies. *J*
403 *Clin Invest* **2011**, 121: 2750-2767.
- 404 22. Reis-Filho JS and Puzstai L. Gene expression profiling in breast cancer: classification, prognostication, and
405 prediction. *Lancet* **2011**, 378: 1812-1823.
- 406 23. Vieira AF, Schmitt F. An Update on Breast Cancer Multigene Prognostic Tests-Emergent Clinical
407 Biomarkers. *Front Med (Lausanne)*. **2018**, 5:248.
- 408 24. Goldhirsch A, Winer EP, Coates AS, Gelber RD, Piccart-Gebhart M, Thürlimann B, Senn HJ, Personalizing
409 the treatment of women with early breast cancer: highlights of the St Gallen International Expert Consensus
410 on the Primary Therapy of Early Breast Cancer 2013, *Ann Oncol.* **2013**, 24(9):2206-23.
- 411 25. Brenton JD, Carey LA, Ahmed AA and Caldas C. Molecular classification and molecular forecasting of
412 breast cancer: ready for clinical application? *J Clin Oncol* **2003**, 23: 7350-7360.
- 413 26. TCGA. Comprehensive molecular portraits of human breast tumours. *Nature* **2012**; 490: 61-70.
- 414 27. Curtis C, Shah SP, Chin SF, *et al.* The genomic and transcriptomic architecture of 2,000 breast tumours
415 reveals novel subgroups. *Nature* **2012** Apr; 486(7403):346-352
- 416 28. Paquet ER and Hallett MT. Absolute assignment of breast cancer intrinsic molecular subtype. *J Natl Cancer*
417 *Inst.* **2015**, 107:357.
- 418 29. Carey LA, Berry DA, Cirincione CT, *et al.* Molecular heterogeneity and response to neoadjuvant human
419 epidermal growth factor receptor 2 targeting in CALGB 40601, a randomized phase III trial of paclitaxel
420 plus trastuzumab with or without lapatinib. *J Clin Oncol.* **2016**, 34:542-9.
- 421 30. Kim HK, Park KH, Kim Y, Absolute assignment of breast cancer intrinsic molecular subtype Discordance
422 of the PAM50 Intrinsic Subtypes Compared with Immunohistochemistry-Based Surrogate in Breast Cancer
423 Patients: Potential Implication of Genomic Alterations of Discordance, *Cancer Res Treat.* **2019**, 51(2):737-747
- 424 31. Tao M, Song T, Du W, *et al.* Classifying Breast Cancer Subtypes Using Multiple Kernel Learning Based on
425 Omics Data. *Genes (Basel)* **2019**, 10(3):200.
- 426 32. Søskilde R, Persson H, Ehinger A, *et al.* Refinement of breast cancer molecular classification by miRNA
427 expression profiles. *BMC Genomics* **2019**, 20, 503.
- 428 33. Raj-Kumar P, Liu J, Hooke JA, *et al.* PCA-PAM50 improves consistency between breast cancer intrinsic and
429 clinical subtyping reclassifying a subset of luminal A tumors as luminal B. *Sci Rep* **2019**, 9, 7956
- 430 34. Park ST and Kim J, Trends in Next-Generation Sequencing and a New Era for Whole Genome Sequencing,
431 *Int Neurolog J.* **2016**, 20(Suppl 2): S76-83.
- 432 35. Li Y, Kang K, Krahn JM, *et al.* A comprehensive genomic pan-cancer classification using The Cancer
433 Genome Atlas gene expression data. *BMC Genomics* **2017**, 18, 508
- 434 36. Sinn, B.V., Fu, C., Lau, R. *et al.* SET_{ER/PR}: a robust 18-gene predictor for sensitivity to endocrine therapy for
435 metastatic breast cancer. *npj Breast Cancer* **2019**, 5:16
- 436 37. Symmans WF, Hatzis C, Sotiriou C, *et al.* Genomic index of sensitivity to endocrine therapy for breast
437 cancer. *J. Clin. Oncol.* **2010**, 28:4111-4119
- 438 38. McInnes L and Healy J. UMAP: Uniform Manifold Approximation and Projection for Dimension
439 Reduction. *ArXiv* **2018**, abs/1802.03426.
- 440 39. Prat A, Pascual T, and Adamo B, Intrinsic molecular subtypes of HER2+ breast cancer, *Oncotarget.* **2017** Sep
441 26; 8(43): 73362-73363.

- 442 40. Ferrari A, Vincent-Salomon A, Pivot X, *et al.* A whole-genome sequence and transcriptome perspective on
443 HER2-positive breast cancers. *Nat Commun* **2016**, 7, 12222.
- 444 41. Montemurro F, Cosimo SD, and Arpino G, Human epidermal growth factor receptor 2 (HER2)-positive
445 and hormone receptor-positive breast cancer: new insights into molecular interactions and clinical
446 implications, *Ann Oncol.* **2013**, 24 (11): 2715-24.
- 447 42. Daemen, A and Manning, G, HER2 is not a cancer subtype but rather a pan-cancer event and is highly
448 enriched in AR-driven breast tumors. *Breast Cancer Res* **2018**, 20, 8.
- 449 43. Prat A, Parker JS, Fan C, *et. al.* Concordance among gene expression-based predictors for ER-positive breast
450 cancer treated with adjuvant tamoxifen. *Ann Oncol* 2012, 23(11):2866-2873
- 451 44. Prat A, Cheang MCU, Galván P, *et al.* Prognostic value of intrinsic subtypes in hormone receptor-positive
452 metastatic breast cancer treated with letrozole with or without lapatinib. *JAMA Oncol.* 2016;2(10):1287–1294.
- 453 45. Friedman J, Hastie T, and Tibshirani R, Regularization Paths for Generalized Linear Models via Coordinate
454 Descent. *Journal of Statistical Software* **2010**, 33 (1): 1–22.
- 455 46. Simon N, Friedman J, Hastie T, and Tibshirani R, Regularization Paths for Cox’s Proportional Hazards
456 Model via Coordinate Descent. *Journal of Statistical Software*, **2011**, 39 (5): 1–13.
- 457 47. Tibshirani R, Bien J, Friedman J, Hastie T, Simon N, Taylor J, and Tibshirani RJ, Strong Rules for Discarding
458 Predictors in Lasso-Type Problems. *Journal of the Royal Statistical Society: Series B* **2012**, 74 (2): 245–66.
- 459 48. Cox DR, Regression models and life-tables. *J. R. Statist. Soc. B* **1972**, 34, 187–220.
- 460



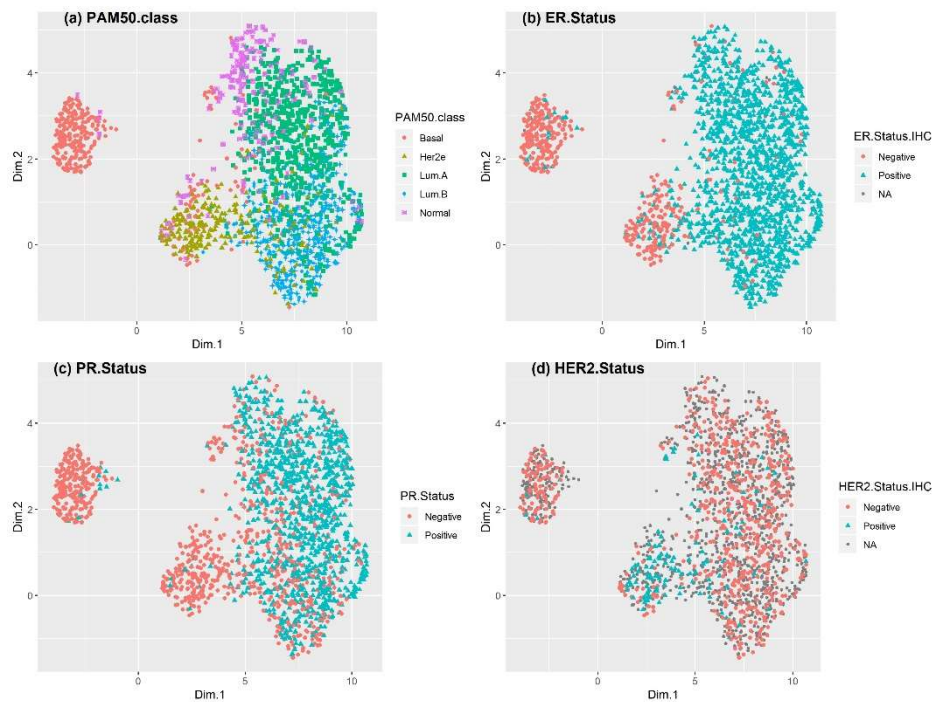
© 2020 by the authors. Submitted for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0/>).

461

462

463 **Supplemental Figures**

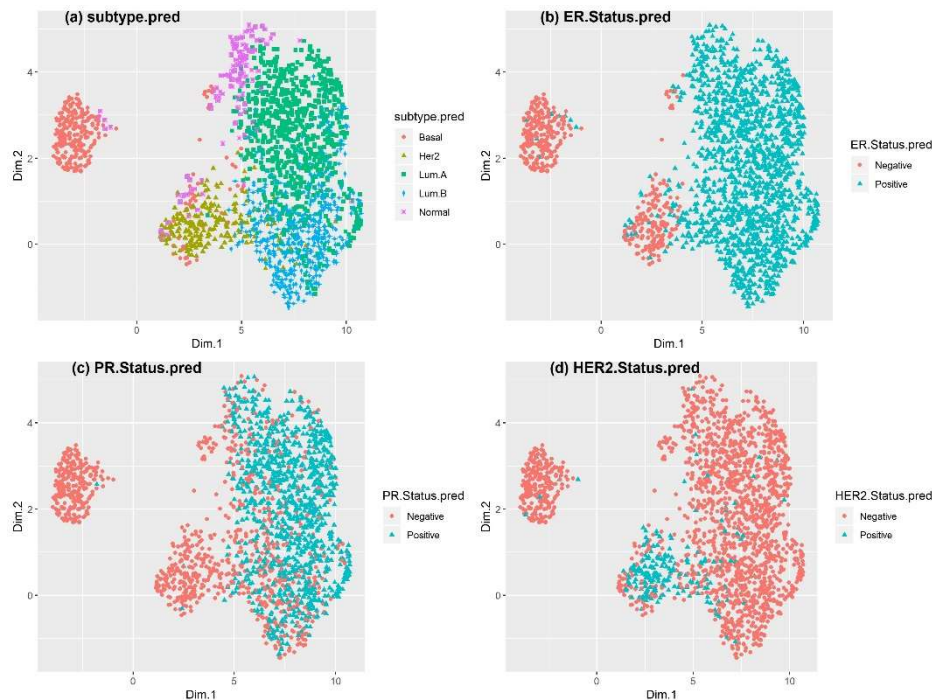
464



465
466
467
468
469

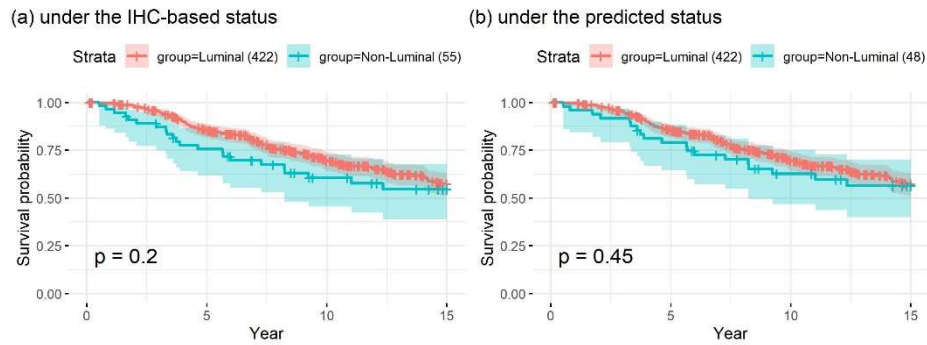
Supplemental Figure 1. UMAP plot showing receptor status of patients in the METABRIC dataset. The tumor subtype and ER, PR, and HER2 status were based on the available clinical data. Gray points are samples with no available clinical information. The UMAP plot of the METABRIC dataset revealed a similar macroscopic landscape to that for TCGA.

470



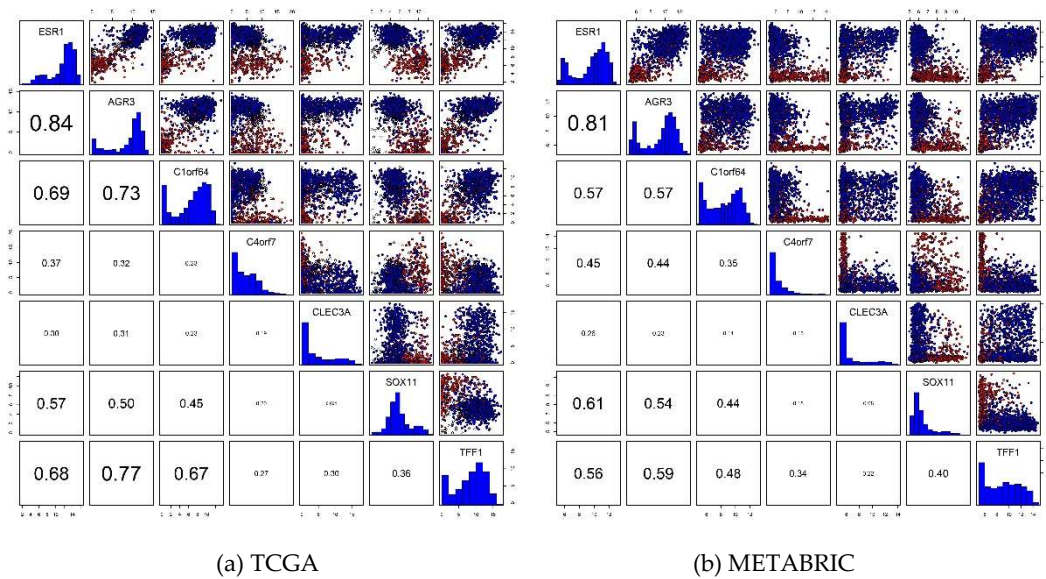
471
472
473
474
475

Supplemental Figure 2. UMAP plot showing receptor status of patients in the METABRIC dataset. GEP-based prediction was used to determine the subtype, as well as the status of ER, PR, and HER2. Similar to TCGA, the predicted ER and HER2 status (but not PR) was mostly in accordance with the corresponding pattern of receptor status for the basal-like, luminal A, and luminal B subtypes.



476

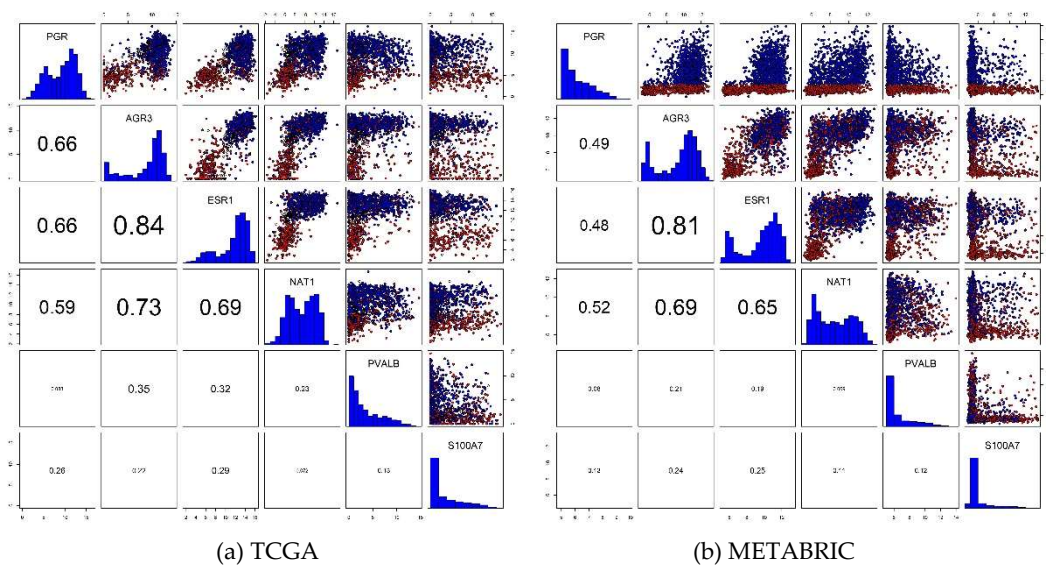
477 **Supplemental Figure 3.** Scatter plots and Pearson's correlation coefficients of seven predictor genes for ER
 478 status prediction. Blue: ER+; red: ER-; empty circle: NA. ER status characterization was based on IHC.



479

480

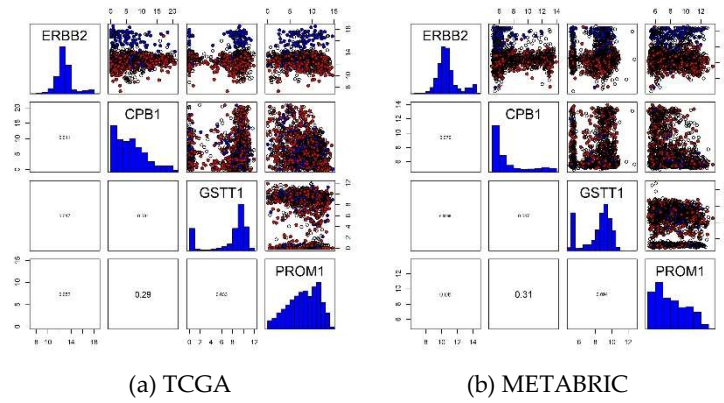
481 **Supplemental Figure 4.** Scatter plots and Pearson's correlation coefficients of seven predictor genes for ER status
 482 prediction. Blue: ER+; red: ER-; empty circle: NA. ER status characterization was based on IHC.



483

484

485 **Supplemental Figure 5.** Scatter plots and Pearson's correlation coefficients of six predictor genes for PR status
 486 prediction. Blue: PR+; red: PR-; empty circle: NA. The PR status of TCGA samples was based on IHC, whereas
 487 that of METABRIC samples was not based on IHC. The PR-status predictor genes included *ESR1* and *AGR3*,
 488 which were also predictor genes for ER status.



489
490

491 **Supplemental Figure 6.** Scatter plots and Pearson's correlation coefficients of four predictor genes for HER2
492 status prediction. Blue: HER2+; red: HER2-; empty circle: NA. *CPB1*, *GSTT1*, and *PROM1* showed weak
493 correlations with *ERBB2*, implying that HER2 status prediction was determined predominantly by *ERBB2*.