

Proteogenomic Metabolism Dependent Signature Identifies Worse Outcome in Breast And Other Cancers

Guisong Wang

Henry M. Jackson Foundation for the Advancement of Military Medicine

Punit Shah

BERG LLC

Richard Searfoss

BERG

Jianfang Liu

Chan Soon-Shiong Institute of Molecular Medicine at Windber

Jamie Leigh Campbell

Henry M. Jackson Foundation for the Advancement of Military Medicine

Brenda Deyarmin

Chan Soon-Shiong Institute of Molecular Medicine at Windber

Rebecca Zingmark

Henry M. Jackson Foundation for the Advancement of Military Medicine

Bradley Mostoller

Chan Soon-Shiong Institute of Molecular Medicine at Windber <https://orcid.org/0000-0002-1242-7162>

Leonid Kvecher

Chan Soon-Shiong Institute of Molecular Medicine at Windber

Stella Somiari

Chan Soon-Shiong Institute of Molecular Medicine at Windber

Praveen-Kumar Raj-Kumar

Chan Soon-Shiong Institute of Molecular Medicine at Windber

Lori Sturtz

Chan Soon-Shiong Institute of Molecular Medicine at Windber

Elder Granger

BERG

Linda Vahdat

Memorial Sloan Kettering Cancer Center

Chas Bountra

University of Oxford

Rangaprasad Sarangarajan

BERG

Mary Cutler

Uniformed Services University of the Health Sciences

Niven Narain

BERG LLC

Jeffrey Hooke

The Henry M Jackson Foundation for the Advancement of Military Medicine

Craig Shriver

Uniformed Services University of the Health Sciences

Hai Hu

Chan Soon-Shiong Institute of Molecular Medicine at Windber <https://orcid.org/0000-0001-5345-8371>

Michael Kiebish

BERG

Albert Kovatich (✉ akovatich@hjfresearch.org)

Uniformed Services University of the Health Sciences

Article

Keywords:

Posted Date: March 17th, 2022

DOI: <https://doi.org/10.21203/rs.3.rs-1443160/v1>

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

[Read Full License](#)

1 **Proteogenomic Metabolism Dependent Signature Identifies Worse**
2 **Outcome in Breast And Other Cancers**

3 Guisong Wang^{1,2}, Punit Shah³, Richard M. Searfoss³, Jianfang Liu⁴, J. Leigh Fan-
4 tacone-Campbell^{1,2}, Brenda Deyarmin⁴, Rebecca N. Zingmark^{1,2}, Bradley Mostoller⁴,
5 Leonid Kvecher⁴, Stella Somiari⁴, Praveen-Kumar Raj-Kumar⁴, Lori A. Sturtz⁴, Elder
6 Granger³, Linda Vahdat⁵, Chas Bountra⁶, Rangaprasad Sarangarajan³, Mary L. Cutler⁷,
7 Niven R. Narain³, Jeffrey A. Hooke^{1,2,7}, Craig D. Shriver^{1,8}, Hai Hu⁴, Michael A. Kiebish³,
8 Albert J. Kovatich^{1,2,*}

9 ¹Murtha Cancer Center / Research Program, Department of Surgery, Uniformed Ser-
10 vices University of the Health Sciences, Bethesda, MD, USA.

11 ²The Henry M Jackson Foundation for the Advancement of Military Medicine Inc., Be-
12 thesda, MD, USA.

13 ³BERG, Framingham, MA, USA.

14 ⁴Chan Soon-Shiong Institute of Molecular Medicine at Windber, Windber, PA, USA.

15 ⁵Memorial Sloan Kettering Cancer Center, New York City, NY, USA.

16 ⁶Department of Clinical Medicine, University of Oxford, Oxford OX3 7LF, UK.

17 ⁷Department of Pathology, Uniformed Services University of the Health Sciences, Be-
18 thesda, MD, USA.

19 ⁸Department of Surgery, Uniformed Services University of the Health Sciences, Bethes-
20 da, MD, USA.

21 **Corresponding Author:**

22 Albert J. Kovatich

23 Murtha Cancer Center Research Program-USU Walter Reed Surgery

24 6720A Rockledge Dr.

25 Bethesda, MD 20817

26 240-694-2557

27 akovatich@hifresearch.org

28

29 **Abstract**

30 Breast cancer (BrCA) therapeutic selection routinely incorporates clinicopathologic in-
31 formation along with immunohistochemistry (IHC) for ER/PR/HER2/Ki-67. However, this
32 is incomplete and has shortcomings that are seen in clinical outcome differences even
33 within the same subtype. Herein, we analyzed the proteome of 116 HER2-negative pri-
34 mary BrCA samples and subsequently validated a 34-proteogenomic signature in 5,963
35 BrCA tumor samples from TCGA, METABRIC, and GSE96058 that demonstrated a
36 metabolic enrichment signature impacting overall survival, progression free survival,
37 and response to therapy. The 34-proteogenomic signature selected ER+ BrCA tumors
38 for upstaging to a more triple negative pathophysiological phenotype, herein referred to
39 as Luminal/TN-like (L/T), impacting likelihood for chemotherapy consideration and other
40 therapeutic modalities rather than hormonal therapy alone. Further, analysis of 9,530
41 tumors across 33 types of cancers in TCGA demonstrated the 34 proteogenomic signa-
42 ture utility in the reclassification of other cancer types into different risk groups.

43 **Introduction**

44 Precision classifying of BrCA subgroups is critical to support appropriate therapeutic in-
45 terventions. In practice, IHC markers are incorporated with clinicopathologic information
46 for this clinical decision. Estrogen receptor (ER), progesterone receptor (PR), human
47 epidermal growth factor receptor 2 (HER2), and Ki-67 are important IHC markers well
48 established in clinical practice. As previously reported, the existence of ER-negative
49 PR-positive (ER-PR+) BrCA is highly unlikely and possibly inaccurate¹. Therefore, BrCA
50 tumors for this study were classified into five IHC subtypes: Luminal A (LA, ER+/HER2-
51 /Ki-67-), Luminal B1 (LB1, ER+/HER2-/Ki-67+), Luminal B2 (LB2, ER+/HER2+), HER2+

52 (ER-/PR-/HER2+) and Triple Negative Breast Cancer or TNBC subtype (ER-/PR-
53 /HER2- or ER-/HER2-)². Among them, Luminal (LA, LB1) BrCA is considered the least
54 aggressive and has the highest relative survival rate. TNBC BrCA is the most aggres-
55 sive and has the lowest relative survival rate.

56 Although IHC-based subtyping is the gold standard in clinical diagnosis and
57 treatment selection, emerging high throughput technologies provide a more compre-
58 hensive molecular characterization of breast tumors. Numerous studies have demon-
59 strated that at the genomic and transcriptomic levels, patients with the same IHC sub-
60 type have different molecular profiles and are associated with different clinical outcomes
61 and treatment responses³⁻⁶. A 50-gene signature (PAM50), developed from gene ex-
62 pression microarray platform, has been applied to classify BrCA into five intrinsic mo-
63 lecular (gene-based) subtypes: Luminal A, Luminal B, HER2-enriched, Basal-like, and
64 Normal-like. It has been reported that the PAM50 subtypes and IHC-based subtypes
65 are discrepant^{3,4}. In addition, five heterocellular subtypes were identified from Luminal
66 breast cancer samples which were associated with varying clinical outcomes and treat-
67 ment responses⁵. Also, a Luminal immune-positive subgroup in TNBC samples was as-
68 sociated with an improved prognosis⁶. This provides an opportunity for additional granu-
69 larity with survival implications.

70 Currently available molecular signatures developed for BrCA subtyping are most-
71 ly based on mRNA expression or genomic mutation profiles. Multigene expression as-
72 says such as PAM50³, BluePrint⁷ and MammaPrint⁸, Oncotype Dx⁹ have been estab-
73 lished to stratify patients into different risk groups. It has been reported that the correla-
74 tion between a gene and its corresponding protein expression was relatively low and

75 that the mRNA expression level of a given gene may not be suitable to represent the
76 corresponding protein expression level. For example, studies showed that the median
77 Pearson correlation value of mRNA-to-protein was 0.39¹⁰, the mean Spearman's corre-
78 lation coefficient between mRNA and protein abundance was 0.31 in tumors and 0.19 in
79 adjacent non-cancerous tissues¹¹. As proteins, not genes, convey the actual functional
80 properties of cells, the classification based on protein expression may be more accurate
81 to reflect the cell functional heterogeneity, thus appropriate to investigate a biomarker
82 signature at the protein expression level. A few researchers have already investigated
83 BrCA subtyping using quantitative proteomics data. The Super-SILAC mix technique for
84 quantitative proteomics of BrCA tumors was used to identify unique features not ob-
85 servable by genomic approaches¹². LC-MS/MS-based protein quantification was utilized
86 to investigate the PAM50 subtype and showed Luminal B and HER2 subtypes were in-
87 termixed at the protein level and basal-like subtype was separated into two groups¹³.

88 Moreover, most of the molecular subtyping was generally achieved through un-
89 supervised clustering methods on the molecular profiles of tumors. BluePrint, an 80-
90 gene signature developed from microarray gene expression data, is the only signature
91 using the IHC-based clinical subtype as a guide⁷. Considering that IHC-based BrCA
92 subtypes have been widely used and accepted for stratifying patients into different
93 treatment groups in clinical practice, there is clear benefit to using available IHC sub-
94 type information to identify differentially expressed biomarkers between BrCA subtypes
95 as the basis of developing biomarker panels. In previous research, the use of a label-
96 free protein quantification technique on formalin-fixed, paraffin-embedded HER2- breast
97 tumors to generate global proteomics data and perform differential analysis (Signifi-

98 cance Analysis of Microarrays) identified 224 differentially expressed proteins between
99 ER+ BrCAs and TNBCs and demonstrated that a subset of ER+/PR+ BrCA tumors had
100 a protein expression profile and clinical outcomes similar to those of TNBCs¹⁴.

101 In this study, available clinical IHC results from core biopsies were used to select
102 the cohort of patients' tumors from flash-frozen surgical samples. These 116 unifocal
103 primary tumors characterized by IHC as HER2-, including Luminal (LA, LB1) and TNBC
104 tumors, were selected for proteomic analysis. Liquid chromatography-tandem mass
105 spectrometry (LC-MS/MS) based quantitative proteomic analysis with the tandem mass
106 tags (TMT) labeling was conducted on these tumors to generate the global proteomics
107 data. Consensus altered proteins between TNBC and Luminal subtypes were first iden-
108 tified, followed by analyses of mRNA expression, gene-protein correlation, pathway en-
109 richment, and proteogenomic characteristics using TCGA-BRCA and CPTAC proteoge-
110 nomics data to identify a proteogenomic metabolism dependent predictive and prognos-
111 tic signature. This signature demonstrates the potential to enhance the clinical IHC sub-
112 typing-based diagnosis for HER2- patients, enabling stratification of patients into differ-
113 ent risk groups and providing potential targetable interventions. Moreover, this study
114 demonstrates the potential to stratify the clinical subgroups of all breast cancers into dif-
115 ferent risk groups, as well as separate the patients associated with worse survival from
116 the patients associated with good survival. In addition, the signature could apply to mul-
117 tiple cancer types and stratify the cancer patients into groups with significantly different
118 outcomes demonstrated from The Cancer Genome Atlas (TCGA) data sets.

119 **Results**

120 **Demographic and clinicopathological characteristics of the study cohort**

121 The demographic and clinicopathological characteristics of patients are shown in Table
122 1. There were 32 (27.6%) LA, 69 (59.5%) LB1 and 15 (12.9%) TN in the cohort. We
123 designated ER+/HER2- including LA and LB1 as Luminal and defining the cohort with
124 two subtypes (Luminal and TN) as Luminal-TN cohort in this paper. Of the 116 cases,
125 94 cases were ER+ (ER>10%), 15 cases were ER- (ER<1%) and 7 cases were low
126 ER+ (ER between 1% and 10%). To consolidate our findings, we split our cohort into a
127 training cohort and a testing cohort using the stratified random sampling method based
128 on IHC-based subtype. To avoid any analysis biases caused by low ER+ cases, these
129 low ER+ cases were removed from our training and testing cohorts (see Results sec-
130 tion: Low ER+ BC tumors are closer to ER- BC tumors than ER+ BC tumors). We ana-
131 lyzed 70 cases with 61 ER+ and 9 ER- in the training cohort and 39 cases with 33 ER+
132 and 6 ER- in the testing cohort. There were no statistically significant differences be-
133 tween training and testing cohort among each characteristic shown in Table 1.

134 **MS-based proteomics quantification of the study cohort**

135 TMT-labeled LC-MS/MS-based proteomics quantification was performed as described
136 in the section of Methods. A total of 7990 proteins were detected at a 1% false-
137 discovery rate (FDR) in our sample cohort, with 4422 proteins expressed across all
138 samples. To avoid any data analysis bias caused by estimating the abundance of the
139 undetected proteins, we focused our further analyses on these 4422 proteins.

140 **Low ER+ BrCA tumors are closer to ER- BrCA tumors than ER+ BrCA tumors**

141 The 1,521 most variably expressed proteins were used in CPTAC-BRCA subtyping
142 analysis¹⁰. 901 proteins were within our common detected 4,422 proteins from our co-
143 hort. To investigate the association between the IHC-based Luminal-TN subtypes and
144 the proteomic clusters, these 901 proteins were utilized for unsupervised clustering
145 analysis of our 116 cases. Unsupervised agglomerative hierarchical clustering analysis
146 demonstrated that most of the low ER+ BrCAs clustered with ER- BrCAs instead of ER+
147 BrCAs (Supplementary Figure S1), which is very consistent with the previous reports
148 from gene expression experiments¹⁵⁻¹⁷.

149 **Integrated and consensus data analyses identify metabolism dependent 34-gene**
150 **panel**

151 To avoid subtype bias and identify the significantly differentially expressed proteins be-
152 tween TN and ER+/HER2- cases, comparative analyses were performed for TN versus
153 LA and TN versus LB1 separately. To obtain robust and stable significantly altered pro-
154 teins from the training cohort, the consensus differential analysis (see Methods) were
155 performed. For each comparison, the significance was reported at Benjamini-Hochberg
156 (BH) adjusted p-value < 0.05 and (fold change (FC) > 1.5 or FC <0.67). The consensus
157 differential analyses identified 512 significantly differentially enriched proteins for the
158 comparison of TN versus LA BrCAs from the training dataset (Supplementary Table S1
159 and Figure 1A), where 242 proteins were up-regulated and 270 proteins were down-
160 regulated in TNBCs. Similarly, the comparison between TN and LB1 BrCAs generated
161 226 significantly differentially expressed proteins with 108 proteins up-regulated and
162 118 proteins down-regulated in TNBCs (Supplementary Table S1 and Figure 1B). The

163 Venn diagram shows that 164 significantly differentially expressed proteins were de-
164 tected from both comparisons, where 75 proteins are up-regulated and 89 proteins
165 down-regulated in TNBCs compared to luminal BrCAs (LA/LB) (Supplementary Table
166 S1 and Figure 1C). Among these 164 proteins, 153 proteins or their coding genes were
167 significantly detected in all four independent public datasets: Clinical Proteomic Tumor
168 Analysis Consortium (CPTAC), TCGA, Molecular Taxonomy of Breast Cancer Interna-
169 tional Consortium (METABRIC) and GSE96058.

170 Further filtering was performed (Figure 1D). A total of 811 HER2- primary female
171 BrCA samples with at least 30 days follow-up and measured by RNA sequencing (642
172 luminal and 169 TNs) were identified from the TCGA cohort for overall survival (OS) and
173 progression-free interval (PFI) analyses. Survival analyses unveiled 22 genes whose
174 higher expression is significantly associated with favorable outcomes (OS and PFI with
175 log-rank p -value <0.05). Compared with the results from the differential analyses, 20 of
176 the 22 genes have relative expression levels that are aligned with the relative outcomes
177 of the two subtypes (Luminal and TN). However, 2 of the 22 genes were actually ex-
178 pressed higher in Luminal but their higher expression was associated with worse out-
179 comes, contradicting the outcome results of the two subtypes. Therefore, these 2 genes
180 were removed from subsequent study. Similarly, of the 18 genes whose higher expres-
181 sion is significantly associated with unfavorable outcomes, 3 were removed from the
182 subsequent study because of the contradicting results from differential and survival
183 analyses. We further investigated the correlation of the selected proteins from our train-
184 ing dataset since highly correlated proteins are generally functionally related and the
185 linear model could benefit from reducing the level of correlation between the predictors.

186 The high correlation was determined by the Pearson correlation $> 0.7^{18}$. Correlation
187 analysis showed two genes (PLOD1, COLGALT1) were highly correlated (Pearson cor-
188 relation=0.78) in the training dataset; COLGALT1, which had a higher p-value from the
189 differential analysis, was removed from the list while PLOD1 was retained. Consequent-
190 ly, a set of 34 genes constituted the genes and proteins of interest for this study. A total
191 of 20 genes were down-regulated in TN breast tumors and associated with good OS
192 and PFI, whereas 14 genes were up-regulated in TN breast tumors and associated with
193 poor OS and PFI (Supplementary Table S1 and Figure 1E). Further investigation of the
194 gene-protein expression correlation showed 31 of the 34 proteins had moderate or high
195 gene-protein expression correlation (Pearson correlation > 0.39). The 3 biomarkers
196 (SLC2A1, NCBP1 and PHPT1) with low gene-protein expression correlation were inves-
197 tigated in the literature. These 3 biomarkers were potential biomarkers for cancer thera-
198 py or related to cancer development and were retained in our protein panel list¹⁹⁻²². The
199 moderate or high gene-protein expression correlation in our identified biomarker set in-
200 dicates the 34 biomarkers may not only be a protein signature but also a gene signature
201 for subtype prediction.

202 KEGG pathways involved in any 34 genes were extracted. Among 24 of 34
203 genes involved in KEGG pathways, 14 genes were involved in metabolic pathways
204 (Supplementary Table S2). The KEGG pathway over-representation analysis also
205 demonstrated that some metabolic pathways are significant at $p < 0.05$. These significant
206 metabolic pathways are involved in amino acid metabolism (alanine, aspartate, gluta-
207 mate, valine, leucine, isoleucine), one carbon metabolism, purine metabolism, pyrimi-

208 dine metabolism, nicotinate and nicotinamide metabolism, biosynthesis of cofactors and
209 fatty acid metabolism (Supplementary Table S2).

210 **Proteogenomic characterization of 34 genes reveals high fractions of positive cis** 211 **effects of CNVs on mRNA and protein**

212 The proteogenomic characteristics of 34 genes were investigated utilizing CPTAC pro-
213 teogenomic data analysis results from Mertins et al.¹⁰. While 24 single amino acid vari-
214 ants (SAAVs) and 6 novel splice isoforms involved in 18 of the 34 proteins were detect-
215 ed (Supplementary Table S3, corresponding to supplementary Table 5 and 10 by Mer-
216 tins et al.¹⁰), the number of variants at peptides were low. The consequence analyses of
217 copy number alterations (CNVs) on RNA and protein showed 27 of the 31 genes (87%)
218 have significant positive cis effects on their mRNA expression and 17 of the 30 genes
219 (57%) have significant positive cis effects on their protein abundance. The fractions of
220 the significant positive cis effects on mRNA and protein in 34 genes are both significant
221 compared with the fractions of all significant positive cis effects on mRNA (64%) and
222 protein (31%) (one-sided fisher test p-value=0.0037 for mRNA and 0.0035 for protein).
223 The observations are consistent with the BrCA and colon cancer analysis that metabolic
224 functions were enriched in genes with the positive cis effect of CNVs on mRNA^{10, 23}.

225 **34-biomarker signature distinguishes two distinct tumor subtypes**

226 Unsupervised hierarchical clustering heatmaps of the training cohort with 34 proteins
227 (Figure 2A) demonstrated that there were two distinct clusters, one cluster consisted
228 mostly of IHC-based luminal tumors, while the other was mostly TN tumors. The unsu-
229 pervised clustering heatmaps of testing cohort and CPTAC HER2- cohort (53 tumors,
230 after removing the known low ER+ tumors) with 34 proteins also demonstrated the

231 same pattern from independent proteomic datasets: there were two distinct clusters,
232 one cluster mapped well with IHC-based luminal subtype and PAM50 Luminal subtype
233 (Luminal A and Luminal B) and another mapped well with IHC-based TN subtype and
234 PAM50 Basal-like subtype. Unsupervised clustering heatmaps of TCGA HER2- cohort
235 (799 tumors), METABRIC HER2- cohort (1645 tumors) and GSE96058 HER2- cohort
236 (2535 tumors) after removing known low ER+ tumors using 34 protein-coding genes as
237 features also demonstrated similar patterns from independent datasets at the gene ex-
238 pression level (Figure 2B). These findings demonstrated that the 34 proteins or protein-
239 coding genes are a strong predictive protein or gene signature to stratify HER2- patients
240 into Luminal-like patients and TN-like patients from both protein abundance and gene
241 expression profiles.

242 To define the solid novel proteomic subtypes, consensus clustering analysis of
243 the training cohort using 34 proteins was performed to investigate the optimal number of
244 clusters and the corresponding clusters from the training cohort. The cluster results
245 demonstrated that two clearly distinct groups were identified (Figure 3). One was de-
246 fined as a Luminal-like subtype and another one as a TN-like subtype by the Fisher's
247 exact test for significance. We designated Luminal-like and TN-like subtypes as LT34
248 subtypes.

249 **Centroid model was used to predict LT34 subtype**

250 The centroid of each LT34 subtype was determined by calculating the median of the
251 normalized protein abundance values of the samples within the subtype for each of 34
252 proteins from the training dataset and were defined as LT34 centroids (Supplimentary
253 Table S4). The LT34 subtype for each sample in all cohorts was determined by the

254 nearest centroid method through comparing the Spearman's rank correlation between
255 the sample's 34-protein profile and the centroid profile of LT34 subtypes. Four en-
256 hanced subtypes (referred to as IHC-LT34 here) based on both IHC and LT34 subtype
257 were further defined for each sample: L/L (Luminal determined by IHC and Luminal-like
258 determined by LT34), L/T (Luminal determined by IHC and TN-like determined by
259 LT34), T/L (TN determined by IHC and Luminal-like determined by LT34), and T/T (TN
260 determined by IHC and TN-like determined by LT34). The Spearman correlation with
261 LT34 centroids, LT34 subtypes, four enhanced subtypes, the IHC-based subtype,
262 PAM50/CLAUDIN subtype, grade, stage, survival and treatment information, as well as
263 34 proteins/genes normalized expression values for all samples across datasets in
264 breast cancer are shown in Supplementary Table S5. To be consistent with the availa-
265 ble public clinical data (see Methods), we used available PAM50 subtypes in TCGA and
266 GSE96058, and kept the PAM50 + Claudin-low subtypes in METABRIC cohort down-
267 loaded from cBioPortal.

268 **L/T subtype was associated with worse prognosis compared with L/L subtype**

269 To have a robust survival estimate, generate appropriate survival results and reduce the
270 uncertainty of a survival estimate caused by a small number of patients at risk at the
271 censoring timepoint and incomplete follow-up data, data maturity analysis was per-
272 formed for each survival curve to investigate if censoring at 5 years was appropriate for
273 each survival analysis (see Methods). The data maturity results of OS analyses by IHC-
274 LT34 subtypes in each cohort and the merged cohort (Supplementary Table S6)
275 demonstrated that all of the OS analyses had robust survival estimates under our crite-
276 ria. Therefore, the OS analyses among IHC-LT34 subtypes are robust and are shown in

277 Figure 4. A higher percentage of L/T subtype patients were deceased compared with
278 the percentage of L/L subtype in each of TCGA, METABRIC, GSE96058, and the
279 merged cohort and were shown in the contingency tables in Figure 4A (Fisher test p-
280 value=0.03 in TCGA, p-value=8.08E07 in METABRIC, p-value=8.55E-07 in GSE96058,
281 p-value=1.24E-12 in merged cohort). There is an equal distribution of survival status in
282 L/T compared with T/T subtype patients in each independent cohort (Fisher test p-
283 value=0.58 in TCGA, p-value=0.28 in METABRIC, p-value=0.09 in GSE96058). The OS
284 Kaplan-Meier (K-M) plots among L/L, L/T and T/T subtypes for Luminal-TN cohort are
285 shown in Figure 4B for each of TCGA, METABRIC, GSE96058 and the merged cohort.
286 The hazard ratios for paired comparison between IHC-LT34 subtypes are shown in Fig-
287 ure 4C and Supplementary Table S7. The survival curves and hazard ratios demon-
288 strate that T/T tumors have the worst outcome whereas L/L have the most favorable
289 outcome. L/T tumors have a statistically significant worse outcome than L/L tumors (p-
290 value <0.05), however, the survival difference between T/T and L/T tumors is not statis-
291 tically significant except in the merged cohort. Significantly, these findings demonstrate
292 that the IHC-based Luminal subtypes contain two distinct subtypes associated with dif-
293 ferent survival and that the signature from our study can distinguish them. One subtype
294 is aggressive, similar to the survival of the T/T subtype. PFI and progression-free sur-
295 vival (PFS) in TCGA cohort and relapse-free survival (RFS) difference in METABRIC
296 cohort among the IHC-LT34 subtypes are shown in Supplementary Figure S2 and Sup-
297plementary Table S7 demonstrated L/T subtype patients were associated with worse
298 PFI/PFS compared with L/L subtype patients but there was no significant difference be-
299tween them, a significant RFS difference between L/T and L/L subtypes was found in

300 the METABRIC cohort but there was no significant PFI/PFS/RFS difference between
301 T/T and L/T subtypes.

302 **IHC-based ER+/HER2- subtype contains at least 3 distinct subtypes**

303 Of 116 cases in our cohort, all 7 low ER+ cases (2 LA and 5 LB1) and all 15 TN cases
304 were identified as TN-like. Among the remaining 94 Luminal cases, 25 of 30 LA cases
305 (83.3%) were predicted as Luminal-like (L/L) while 5 of them (16.7%) were identified as
306 TN-like (L/T), 49 of 64 LB1 cases (76.6%) were identified as Luminal-like (L/L) and 15 of
307 them (23.4%) were identified as TN-like (L/T) (Supplementary Table S8). L/L (or L/T)
308 subtype patients were equally enriched in LA and LB1 subtype patients (Fisher exact p-
309 value=0.59). This finding indicates that cell proliferation as measured by Ki-67 percent-
310 ages and growth may not distinguish the L/T from L/L subtype patients. Significantly, the
311 consensus clustering analysis of L/L subtype patients demonstrates that there are two
312 distinct groups identified in L/L subtype patients (Supplementary Figure S3). The clus-
313 ters were also consistent with the LA/LB1 subtype distribution (Fisher exact p-
314 value=0.0003). Therefore, there are at least three subtypes in ER+/HER2- cases, two
315 subtypes in L/L Luminal-like, and one in the L/T TN-like.

316 **Survival outcome of L/T subtype is similar to T/T rather than L/L subtype with or** 317 **without treatment**

318 The choice of therapy is influenced by numerous factors. In the GSE96058 cohort,
319 available treatments for the patients are endocrine or hormone therapy (ET or HT, we
320 called it HormT or HT here), chemotherapy (ChemoT or CT), or combined treatments of
321 CT+HT. In the METABRIC cohort, available treatments for the patients are HT, radio-
322 therapy treatment (RT), CT, combined treatments of CT+HT, HT+RT, CT+RT and

323 CT+HT+RT. Our results and conclusions about the response to treatment are based on
324 available treatment information from these two datasets.

325 The data maturity results of OS analyses for IHC-LT34 subtypes under each
326 treatment per cohort are shown in Supplementary Table S6. The survival difference of
327 5-year OS among the IHC-LT34 subtypes within each treatment in METABRIC and
328 GSE96058 cohorts are shown Figure 5 and SupplementaryTable S9, where each sur-
329 vival curve has sufficient numbers of samples and follow-ups satisfying our data maturi-
330 ty criteria (see Methods). The results demonstrated a statistically significant OS differ-
331 ence between L/T and L/L subtype patients under each of the following treatments:
332 HT($p=1.1E-8$ in GSE96058 and 0.04 in METABRIC), RT ($p=0.039$ in METABRIC), the
333 combined treatments of CT+HT ($p=0.00014$ in GSE96058), HT+RT ($p=0.0066$ in
334 METABRIC) and CT+HT+RT ($p=0.0022$ in METABRIC). These demonstrated that the
335 L/T subtype patients were still associated with poor survival compared with L/L subtype
336 patients for each treatment. This suggests that L/T subtype patients were resistant to
337 the provided treatments compared to L/L subtype patients.

338 There is no statistically significant difference in OS between L/T subtype and T/T
339 subtype patients under each comparable treatment: CT ($p=0.72$ in GSE96058), HT+RT
340 ($p=0.71$ in METABRIC) and CT+HT+RT ($p=0.28$). This finding suggests that the L/T
341 subtype is closer to T/T when compared with the L/L subtype with or without treatments.
342 The molecular profile and survival outcome of L/T subtype patients are more similar to
343 those of T/T subtype patients.

344 **Survival differs significantly between two LT34 subtypes within each clinical**
345 **group**

346 After removing low ER+ cases from the 5780 samples in the merged cohort (TCGA +
347 METABRIC + GSE96058), there are 5716 samples including 4370 IHC-based Luminal,
348 609 TN, 508 ER+/HER2+ and 229 ER-/HER2+. Of the 4370 IHC-based Luminal tumors,
349 83.6% (3653) tumors were L/L, and 16.4% (717) were L/T. Of the 609 IHC-based TN
350 tumors, 96.1% (585) were T/T, and 3.9% (24) were T/L. Of the 508 IHC-based
351 ER+/HER2+ tumors, 51.2% (260) were predicted as Luminal-like, and 48.8% (248)
352 were predicted as TN-like. Of the 229 IHC-based ER-/HER2+ tumors, 94.3% (216) were
353 TN-like and 5.7% (13) were Luminal-like (Supplementary Table S8). In summary, 16.4%
354 of IHC-based Luminal tumors and 48.8% of ER+/HER2+ tumors were predicted as ag-
355 gressive tumors more similar to IHC-based TN than Luminal tumors. 3.9% of IHC-based
356 TN tumors and 5.7% of ER-/HER2+ tumors were predicted to be the favorable tumors
357 more similar to IHC-based Luminal than TN tumors.

358 Taking into consideration tumor grade, there were 640 G1, 2175 G2 and 1847
359 G3 in the combined cohort. In G1 tumors 93.1% (596) and 6.9% (44) were predicted as
360 Luminal-like and TN-like, respectively. In G2 tumors 85.5% (1859) and 14.5% (316)
361 were predicted as Luminal-like and TN-like. In G3 tumors 43.7% (807) and 56.3%
362 (1040) were predicted as Luminal-like and TN-like (Supplementary Table S8)

363 Next, considering tumor stage, there were 2148 stage I, 2072 stage II, 346 stage
364 III and 52 stage IV (or above) in the combined cohort. In stage 1 tumors, 75.6% (1624)
365 and 24.4% (524) were predicted as Luminal-like and TN-like, respectively. In stage 2

366 tumors, 64.8% (1343) and 35.2% (729) were predicted as Luminal-like and TN-like. In
367 stage 3 tumors, 60.7% (210) and 39.3% (136) were predicted as Luminal-like and TN-
368 like. In stage 4 tumors, 63.5% (33) and 36.5% (19) were predicted as Luminal-like and
369 TN-like (Supplementary Table S8).

370 Next, comparing PAM50 and CLAUDIN gene expression based subtypes, there
371 were 2769 Luminal A, 1346 Luminal B, 366 Normal-like, 522 HER2-enriched, 522 Ba-
372 sal-like and 191 Claudin-low in the combined cohort. 93.1% (2579) and 6.9% (190) of
373 Luminal A tumors were predicted as Luminal-like and TN-like respectively; 71.8% (966)
374 and 28.2% (380) of Luminal B tumors were predicted as Luminal-like and TN-like;
375 17.8% (93) and 82.2% (429) of HER2-enriched tumors were predicted as Luminal and
376 TN-like; 1.3% (7) and 98.7% (515) of Basal-like tumors were predicted as Luminal-like
377 and TN-like; 26.7% (51) and 73.3% (140) of Claudin-low tumors were predicted as Lu-
378 minal-like and TN-like (Supplementary Table S8). In summary, 6.9% of Luminal A tu-
379 mors and 28.2% of Luminal B tumors were predicted as TN-like, whereas only 1.3% of
380 Basal-like tumors and 17.8% HER2-enriched tumors were predicted as Luminal-like.

381 The data maturity results of OS analyses by LT34 subtypes under each clinical
382 characteristic (IHC-based subtype, grade, stage, and PAM50 [CLAUDIN-LOW] per co-
383 hort are shown in Supplementary Table S6. The OS differences by Luminal-like and
384 TN-like within each clinical group are shown in Figure 6 and Supplementary Table S10,
385 where each survival curve satisfies our data maturity criteria. Both the K-M plots and
386 hazard ratio table demonstrated that there was a significant survival difference between
387 Luminal-like and TN-like subtypes within each clinical group, not just from an IHC-based
388 perspective.

389 **Survival differs significantly between two LT34 subtypes in other cancers**

390 9530 primary tumors with follow-up of at least 30 days across 33 different TCGA can-
391 cers were used for pan-cancer survival analyses. This clinical data, the normalized
392 RNA-Seq V2 expressed data of 34 genes at z-score level, as well as the predicted sub-
393 type by LT34 subtype, are shown in Supplementary Table S11. The data maturity re-
394 sults in OS analyses by LT34 subtypes for each cancer showed there were 15 cancers
395 passing our criteria that could be used to perform OS analysis (Supplementary Table
396 S6). There was a significant survival difference between the Luminal-like subtype and
397 TN-like subtype in 9 of these 15 cancers. The OS K-M plots for these 9 cancers are
398 shown in Figure 7 and hazard ratios are shown in Supplementary Table S12.

399 **Discussion**

400 Utilizing LC-MS/MS proteomics data analysis for discovery followed by analyses
401 from mRNA expression, gene-protein correlation, collinearity, pathways, proteogenomic
402 characteristics, we identified a 34 metabolism pathway enriched protein/gene novel bi-
403 omarker panel. This resulted in an easily applied classifier to separate HER2- BrCA pa-
404 tients into Luminal-like and TN-like BrCA patients. We successfully validated our bi-
405 omarker panel and classifier with the use of large external cohorts across different plat-
406 forms, patient treatment responses and survival outcomes. This approach suggests that
407 the 34-biomarker panel and its centroid profile are not technology-dependent and could
408 be adapted to multiple molecular platforms to serve as a solid predictive and prognostic
409 signature. The signature provides additional robust risk information, enhances the accu-
410 racy of patient survival stratification by incorporating available IHC-based biomarker sta-

411 tus and clinical characteristics. By validating treatment responses for the different en-
412 hanced subtypes, this signature provides the potential for personalized treatment strat-
413 egies.

414 It is clinically significant that we identified two subtypes (L/L and L/T) within the
415 IHC-based Luminal subtype. We went on to demonstrate that the L/T subtype patients
416 were significantly associated with worse overall survival and greater resistance to
417 treatments when compared with the L/L subtype patients. We also observed that L/L (or
418 L/T) subtype patients were equally enriched in LA and LB1 subtype patients which sug-
419 gest that the Ki-67 biomarker alone does not capture the L/T subtype from the L/L sub-
420 type. These observations are in agreement with previous reports that a stem-like sub-
421 type was detected from Luminal BrCA samples regardless of Luminal A and Luminal B
422 subtype status^{5,24}. This suggests L/T subtype patients might be undertreated and more
423 aggressive treatment might possibly be considered for them. Further, we identified two
424 distinct clusters in the L/L subtype. These clusters are consistent with the distribution of
425 Ki-67 biomarker-based LA and LB1 subtypes. These findings indicate there are at least
426 three subtypes in ER+/HER2- (Luminal) cases: two L/L Luminal-like and one L/T (TN-
427 like). Further investigation towards developing a biomarker signature of the two sub-
428 types of L/L subtype breast cancers will also be important for patient stratification in the
429 future.

430 Moreover, we identified two subtypes in TNBCs: T/L and T/T subtypes. We could
431 not investigate the T/L subtype currently because of the small number of cases. Re-
432 searchers have reported that a luminal immune-positive subtype with favorable progno-

433 ses was detected in the TNBC subtype⁶. We suspect that the T/L subtype may be an
434 independent subtype associated with better outcomes compared to the T/T subtype,
435 possibly preventing overtreatment for this subtype in the future or consider more target-
436 ed approaches.

437 The significant survival difference between the Luminal-like subtype and TN-like
438 subtype across 9 of 15 suitably analyzed TCGA pan-cancers indicates our LT34 signa-
439 ture can be applied to several other TCGA cancers. Two distinct IHC-based subtypes
440 (ER+/HER2- and TN subtypes) selected for the study are involved in two distinct tumor
441 cell types: Luminal cells and Basal cells. The previous findings demonstrate the tumor
442 cell-of-origin impacts the potential development of the tumor, plays a dominant role for
443 cancers and determines the distinct cancer subtypes within an organ²⁵⁻²⁸. The cell-of-
444 origin mechanism of the LT34 subtypes will be investigated in the future to understand
445 its applications to pan-cancers. Moreover, the identified 34 genes are enriched in bi-
446 omarkers of metabolism. Cancer metabolism is being widely investigated and previous
447 findings show that the activities of oncogenes and tumor suppressor genes are associ-
448 ated with metabolic reprogramming²⁹⁻³². As example, previous research strongly sup-
449 ports our findings regarding ABAT and alanine metabolism in ER positive BrCA³³. The
450 future integrated analyses of the abundances of the metabolites and proteins involved in
451 the metabolic pathways may provide us a deep understanding of the metabolic path-
452 ways employed by BrCA as well as other cancers.

453 **Methods**

454 **BrCA sample selection**

455 Fresh frozen HER2- BrCA tissue samples were obtained from The Clinical Breast Care
456 Project (CBCP). Clinical immunohistochemistry (IHC) subtyping of formalin-fixed paraf-
457 fin-embedded (FFPE) core biopsies was used to select a flash-frozen surgical sample
458 cohort of 116 HER2- BrCA patients. The tumors were all primary single-focal BrCA tu-
459 mors and all surgical samples were collected after immediate surgery. The posi-
460 tive/negative status of ER/PR/HER2 was defined using updated ASCO 2020 guide-
461 lines³⁴. ER status of a sample was determined by the percentage of tumor cell nuclei
462 staining positive by ER immunohistochemistry. The sample is considered ER- if less
463 than 1% of cells stain ER-positive, whereas ER+ samples have $\geq 10\%$ of cells staining
464 positive, the samples with tumor cells having ER positive staining between 1% and 10%
465 are regarded as low ER+. The sample is considered as HER2+ if the HER2- values are
466 3+ and HER2- if the HER2 values are 0, 1+. When HER2 value is 2+, FISH was used to
467 further determine its status following ASCO 2020 guidelines. Ki-67 status was deter-
468 mined using the 2011 St. Gallen's International Expert Consensus recommendations³⁵,
469 where a cut-off of 14% was used to denote Ki67+ or Ki67-.

470 Samples and data for this research were collected from study participants who
471 consented to the protocol, 'Tissue and Blood Library Establishment for the Molecular,
472 Biochemical and Histologic Study of Breast Disease,' at Walter Reed National Military
473 Medical Center (WRNMMC) or at Anne Arundel Medical Center (AAMC), and who
474 agreed to the use of their samples and data in future cancer research.

475 **LC-MS/MS proteomic analysis**

476 **Tissue lysis**

477 Tissue samples were lysed with a 200 uL lysis buffer containing 7M Urea, 2M thiourea,
478 0.1% SDS, 1% Protease and Phosphatase Inhibitor Cocktail, and Optima LC/MS Water.
479 Samples were Homogenized using Omni bead rupter. Homogenized samples were cen-
480 trifuged at 17,000 x g for 10 minutes, and the supernatant was then used in the Brad-
481 ford Assay to determine the protein concentration.

482 **Trypsin digestion**

483 Samples were prepared following the method previously described in Sturtz et al.³⁶. In
484 brief, proteins were reduced and alkylated with 10 mM Tris (2-carboxyethyl) phosphine
485 (TCEP) and 18.75 mM iodoacetamide, respectively. Proteins were then precipitated at -
486 20 °C overnight using cold acetone and pellets were reconstituted in 200 mM tri-
487 ethylammonium bicarbonate (TEAB) and trypsin digested overnight at 37 °C.

488 **TMT labeling of peptides**

489 20 µg of peptides from samples were aliquoted and labeled with 10-plex TMT (Tandem
490 Mass Tag) reagents (Thermo Fisher Scientific) using the manufacture's protocol. Fur-
491 ther, an equal amount of peptides from all samples was pooled to create a reference
492 sample in TMT Channel 126. Samples were incubated at room temperature for 1 hour
493 before being quenched with 5% hydroxylamine for 30 mins. TMT-labeled peptides were
494 mixed and dried in a speedvac (Thermo Fisher Scientific). Dried samples were desalted
495 on C18 spin columns and dried again to be stored at -20°C until LC-MS/MS analysis.

496 **Mass spectrometry**

497 TMT-labeled peptides were analyzed via LC-MS/MS using a Waters nanoAcquity online
498 2-dimensional reversed-phase LC system and a Thermo Q Exactive Plus mass spec-
499 trometer. Nine fractions were created from a single injection of 5 µg of TMT-labeled
500 peptides in the first dimension using 20 mM ammonium formate as Buffer A and 100%
501 acetonitrile as Buffer B, and sequential elutions with 16, 20, 24, 26, 28, 30, 32, 36, and
502 50% of Buffer B. Fractions were further separated in the second dimension over a 170
503 minute gradient using 0.1% formic acid in water as Buffer A and 0.1% formic acid in ac-
504 etonitrile as Buffer B, and a gradual change of 20-23% from Buffer A to Buffer B. MS
505 survey scans were performed at a resolution of 70,000 with a scan range of 400–1800
506 Thomsons (Th; $Th = Da/z$) to select peptides for fragmentation. MS/MS fragment scans
507 were performed at 35,000 resolution consisting of an isolation window of 1.2 Th. Only
508 ions of +2 to +4 charge were ultimately selected for fragmentation.

509 **Protein quantification**

510 Data generated was processed using Proteome Discoverer v1.4 (Thermo Scientific).
511 The database search algorithm SEQUEST was used to search spectra against the Ref-
512 Seq protein database and the reporter ion node to provide relative quantitation for all
513 matching spectra. Protein quantification was reported using only unique peptides, with a
514 minimum of two unique peptides required to identify a protein. Specific search param-
515 eters included peptides of ≥ 6 amino acids and no more than two missed cleavages per
516 peptide. The search utilized a 10 ppm precursor mass tolerance, a 0.02 Da fragment
517 mass tolerance, static N-terminal TMT-10Plex and cysteine carbamidomethylation

518 modifications, and dynamic lysine TMT-10Plex, asparagine/glutamine deamidation, and
519 methionine oxidation modifications.

520 **Statistics Methods**

521 **Sample quality control investigation and normalization**

522 Log₂-transformed raw TMT ratios at the protein level were employed for data analysis.
523 Density plot and dip statistics demonstrated that the protein expression profile of each
524 sample followed an expected unimodal Gaussian distribution. A 2-component Gaussian
525 mixture model-based normalization algorithm used in CPTAC-BRCA was applied to our
526 data for normalization^{10,37,38}. Briefly, the z-scores were calculated for each sample
527 where the center was the median of protein expression values and the standard devia-
528 tion was calculated from the expression abundance of non-changed proteins in the
529 sample compared to the reference pool sample. The non-changed proteins were deter-
530 mined by a 2-component Gaussian mixture model-based method. The z-score method
531 centered on the distribution of the log₂-transformed TMT ratio to zero and utilized the
532 standard deviation of non-regulated proteins compared with the reference pool sample
533 to nullify the effect of different protein loading and systematic MS variation.

534 **Low ER+ BrCA cluster investigation with ER+ and ER- BrCAs**

535 901 proteins common to 1500+ protein-coding genes used in CPTAC-BRCA subtyping
536 analysis¹⁰ were utilized for unsupervised clustering analysis of 116 cases and *Com-*
537 *plexHeatmap* bioconductor package (version 2.8.0) was used for heatmap visualiza-
538 tion³⁹. The Spearman rank correlation distance was used as a distance matrix and

539 Ward's criterion was used as a linkage criterion in the unsupervised hierarchical cluster-
540 ing algorithm.

541 **Independent public BrCA datasets**

542 Independent public BrCA cohorts were extracted as evaluation datasets to evaluate our
543 identified protein signatures. The TCGA normalized RNA-Seq expression data and
544 CPTAC normalized protein abundance data at the z-score level relative to all of the
545 samples were extracted from Bio Cancer Genomics Portal through *cgdsr* Bioconductor
546 packages (version 1.3.0). The status of ER/HER2 for the cases in the CPTAC cohort
547 and TCGA cohort was obtained using the same method as reported in TCGA-BRCA
548 Nature 2012 paper and Huo et al.^{40,41}, whereas the OS, PFI and PFS survival infor-
549 mation were extracted from the Pan-Cancer clinical data resource⁴². TCGA treatment
550 information was processed internally. Primary tumors with at least 30 days survival fol-
551 low-up^{43,44} in the METABRIC study were used as one independent evaluation dataset.
552 The clinical data and normalized expression data at the z-score level of METABRIC co-
553 hort were extracted from *cgdsr* Bioconductor package. Another independent large RNA-
554 Seq validation cohort is Sweden Cancerome Analysis Network – Breast Initiative study
555 (SCAN-B): GSE96058⁴⁵. The primary tumor samples with at least 30 days survival fol-
556 low-up and their normalized expression data were extracted from the GEO data reposi-
557 tory (Reference link: <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE96058>).

558 **Biomarker panel selection**

559 **Consensus differential analyses between TN and Luminal subtypes**

560 Comparative analyses were first performed for TN versus LA and TN versus LB1 in
561 training dataset using Linear Models for Microarray Data (LIMMA)^{46,47} Bioconductor

562 package (version 3.38.3) separately. For each comparison, the significance was report-
563 ed at Benjamini-Hochberg (BH) adjusted p-value < 0.05 and (fold change (FC) > 1.5 or
564 FC < 0.67). The proportional stratified randomly subsampling technique was further em-
565 ployed 100 times, where each subsample cohort was 80% of the training cohort and
566 was stratified by IHC-based tumor subtype. Differential analysis was performed on each
567 sub-sample cohort to compare TN vs. LA and TN vs. LB1 cases separately. The con-
568 sensus significance of one protein is reported if the protein is significant in the training
569 dataset and all of its sub-sample cohort per comparison. The final initial biomarker can-
570 didate pools for TN vs. LA+LB1 consist of the common significantly differentially ex-
571 pressed proteins consensus in both comparisons.

572 **Biomarker panel reduction**

573 TNBCs are more aggressive and associated with poor OS and PFI. To investigate if
574 each coding gene of the identified significantly altered proteins was significantly associ-
575 ated with survival outcome, TCGA cohort and their expressed RNA-Seq data were se-
576 lected for survival analysis. The mapping system in DAVID 6.8⁴⁸ was used to map gene-
577 protein names to avoid any mismatched biomarker names. For each coding gene, its
578 expression values across the cohort were first categorized into low and high expression
579 groups, where the optimal cutoff was determined using the method implemented in the
580 *survMisc* R package^{49,50}. Next, the univariate PFI analysis and OS analysis with the cor-
581 responding optimal cutoff were performed on the cohort and the significance of the as-
582 sociation of each gene with survival outcome was reported by log-rank p-value < 0.05 .
583 The K-M plots were generated to visualize the survival association using the *survival*
584 and *survminer* R package^{51,52}. The biomarkers significantly associated with survival

585 analysis were further selected based on the concordant altered direction so that the se-
586 lected biomarkers up-regulated in TN were associated with poor survival or selected bi-
587 omarkers up-regulated in Luminal were associated with good survival.

588 We further investigated the correlation of the selected proteins from our training dataset
589 since highly correlated proteins are generally functionally related and the linear model
590 could benefit from reducing the level of correlation between the predictors. The high cor-
591 relation was determined by the Pearson correlation > 0.7 . In the identified highly corre-
592 lated proteins, we selected the one with the most significance from the comparative
593 analysis in the training dataset as the representative protein.

594 The selected biomarkers were further investigated with gene-protein correlation, KEGG
595 pathway analysis (see KEGG pathway enrichment analysis) and proteogenomic charac-
596 teristics utilizing CPTAC-BRCA data analyses generated by Mertin et al..

597 **Cluster investigation of 34-biomarker signature**

598 Consensus hierarchical clustering analysis implemented in *ConsensusClusterPlus* R
599 package^{53,54} with the identified 34 proteins was employed to investigate the optimal
600 number of clusters and the corresponding clusters from the training cohort, where
601 spearman correlation was used to generate distance matrix and ward.D was used as
602 the linkage method.

603 To evaluate if the 34-protein/coding-gene panel is a reliable multigene classifier
604 to separate the cohorts into two distinguished clusters and if they could discriminate TN
605 breast tumors from luminal breast tumors, unsupervised hierarchical clustering analysis
606 was applied to the expression data of all of our cohorts separately with these bi-

607 biomarkers and *ComplexHeatmap* bioconductor package was used for heatmap visuali-
608 zation.

609 **LT34 subtype prediction**

610 Each sample's LT34 subtype was defined by the nearest centroid through comparing
611 the Spearman's rank correlation between the sample's 34-protein profile and the cen-
612 troid profile of LT34 subtypes. In short, we calculated the Spearman's rank correlation
613 between one sample's 34 proteins/coding genes profile and the centroids profile of two
614 LT34 subtypes, then assigned the subtype with the higher correlation to the sample.

615 **Data maturity analysis and survival analysis**

616 Data maturity analysis was performed for each survival curve using criterion 1 and crite-
617 rion 2 proposed by GebSKI et al.⁵⁵ to investigate if censoring at 5 years was appropriate
618 for each survival analysis. In short, we used the threshold of the acceptable decrease in
619 the estimated percentage of survival as 5% for individual cohort (or 2.5% for merged
620 cohort) (Criterion 1) and within one-sided 95% CI (Criterion 2) if one extra event oc-
621 curred at the interested time point.

622 The K-M plot was generated using *survminer* R package for each univariate sur-
623 vival analysis. Only the survival curve satisfying all three data mature criteria was shown
624 in the K-M plot. The Cox proportional hazard model implemented in *survival* R package
625 was used to calculate the hazard ratios. The follow-up time was censored at 5 years to
626 investigate the early breast cancer survival outcomes. The significance of the survival
627 difference was reported at log-rank P-value <0.05.

628 **TCGA pan-cancer data**

629 TCGA pan-cancer clinical data across 33 types of cancer were retrieved from TCGA
630 Pan-Cancer Clinical Data Resource (TCGA-CDR, table S1) generated by Liu et al.⁴².
631 Normalized RNA-seq V2 gene expression data at z-score level, median-centered and
632 relative to all samples, were extracted from cBioPortal through *cgdsr* Bioconductor
633 package (version 1.3.0) for each of 33 types of cancers. Three cancers (COAD, READ
634 and UCEC) have few samples comparing the RNA-seq V2 data stored in Broad GDAC
635 firehose. Therefore, the normalized RNA-seq V2 RESM data were downloaded from
636 Broad GDAC firehose for these three cancers and processed by z-score method with
637 median-center and relative to all of the samples. The expressed primary samples with at
638 least 30 days' follow-up were filtered for further data analysis. The nearest centroid
639 method using 34-genes (Spearman's rank correlation with our simple-centroids as dis-
640 tance) was applied to each sample and used to predict the sample's LT34 subtype.

641 **KEGG pathway enrichment analysis**

642 KEGG pathways with genes were downloaded using *ClusterProfiler* Bioconductor pack-
643 age on Feb. 1, 2022 (version 3.18.1)⁵⁶. Pathways involved in any gene of the 34 genes
644 were extracted. Gene set over-representation enriched analysis was performed using
645 the method implemented in the same package to identify significant gene sets.

646 **Data availability**

647 TMT proteomics data for our 116 cases were generated in our internal lab. The TCGA
648 clinical data were downloaded from the supplemental information Table S1 of the Cell
649 paper published by Liu J. et al.. CPTAC normalized protein expression data at z-score

650 level relative to all samples, TCGA normalized RNA-seq V2 gene expression data at z-
651 score level relative to all samples, and METABRIC clinical and normalized microarray
652 data at z-score level relative to all samples were downloaded from Bio Cancer Ge-
653 nomics Portal through *cgdsr* Bioconductor package. GSE96058 clinical and normalized
654 RNA-Seq gene expression data were downloaded through Gene Expression Omnibus
655 (GEO). A total of 5,963 samples across our internal cohorts, TCGA, METABRIC and
656 GSE96058 were processed and the IHC/PAM50/Claudin/LT34/IHC-LT34 subtypes, the
657 OS/PFI/PFS survival information with the normalized 34 biomarker expression values
658 for each sample were generated and available in Supplementary Table S5. TCGA nor-
659 malized RNA-seq V2 gene expression data at z-score level relative to all samples
660 across 33 cancers were extracted from Bio Cancer Genomics Portal through *cgdsr* Bio-
661 conductor package. A total of 9,530 samples across 33 TCGA cancers were processed
662 and LT34 subtype, clinical information as well as normalized gene expression values of
663 34 genes for each sample were available in Supplementary Table S11. The data pro-
664 cessing, analysis and visualization were performed with R/Bioconductor packages.

665 **References**

- 666 1. Foley, N.M. et al. Re-Appraisal of Estrogen Receptor Negative/Progesterone Re-
667 ceptor Positive (ER-/PR+) Breast Cancer Phenotype: True Subtype or Technical
668 Artefact?. *Pathol Oncol Res.* 24, 881–884 (2018).
- 669 2. Raj-Kumar, P.K. et al. PCA-PAM50 improves consistency between breast cancer
670 intrinsic and clinical subtyping reclassifying a subset of luminal A tumors as lu-
671 minal B. *Sci Rep* 9, 7956 (2019).

- 672 3. Parker, J.S. et al. Supervised risk predictor of breast cancer based on intrinsic
673 subtypes. *J Clin Oncol.* 27, 1160–1167 (2009).
- 674 4. Kim, H.K. et al. Discordance of the PAM50 Intrinsic Subtypes Compared with
675 Immunohistochemistry-Based Surrogate in Breast Cancer Patients: Potential Im-
676 plication of Genomic Alterations of Discordance. *Cancer Res Treat.* 51, 737–747
677 (2019).
- 678 5. Poudel, P. et al. Heterocellular gene signatures reveal luminal-A breast cancer
679 heterogeneity and differential therapeutic responses. *NPJ breast cancer* 5, 21
680 (2019).
- 681 6. Prado-Vázquez, G. et al. A novel approach to triple-negative breast cancer mo-
682 lecular classification reveals a luminal immune-positive subgroup with good
683 prognoses. *Sci Rep.* 9, 1538 (2019).
- 684 7. Krijgsman, O. et al. A diagnostic gene profile for molecular subtyping of breast
685 cancer associated with treatment response. *Breast Cancer Res. Treat.* 133, 37–
686 47 (2012).
- 687 8. Van't Veer, L. et al. Gene expression profiling predicts clinical outcome of breast
688 cancer. *Nature* 415, 530-536 (2002).
- 689 9. Paik, S. et al. A multigene assay to predict recurrence of tamoxifen-treated,
690 node-negative breast cancer. *N Engl J Med.* 351, 2817-1826 (2004).
- 691 10. Mertins, P. et al. Proteogenomics connects somatic mutations to signalling in
692 breast cancer. *Nature* 534, 55–62 (2016).

- 693 11. Tang, W. et al. Integrated proteotranscriptomics of breast cancer reveals globally
694 increased protein-mRNA concordance associated with subtypes and survival.
695 *Genome Med.* 10, 94 (2018).
- 696 12. Yanovich, G. et al. Clinical Proteomics of Breast Cancer Reveals a Novel Layer
697 of Breast Cancer Classification. *Cancer Res.* 78, 6001–6010 (2018).
- 698 13. Johansson, H.J. et al. Breast cancer quantitative proteome and proteogenomic
699 landscape. *Nat Commun* 10, 1600 (2019).
- 700 14. Gámez-Pozo, A. et al. Functional proteomics outlines the complexity of breast
701 cancer molecular subtypes. *Sci* 7, 10100 (2017).
- 702 15. Iwamoto, T. et al. Estrogen receptor (ER) mRNA and ER-related gene expres-
703 sion in breast cancers that are 1% to 10% ER-positive by immunohistochemis-
704 try. *J Clin Oncol.* 30, 729–734 (2012).
- 705 16. Deyarmin, B. et al. Effect of ASCO/CAP guidelines for determining ER status on
706 molecular subtype. *Ann Surg Oncol.* 20, 87–93 (2013).
- 707 17. Prabhu, J.S. et al. A Majority of Low (1-10%) ER Positive Breast Cancers Be-
708 have Like Hormone Receptor Negative Tumors. *J Cancer* 5, 156–165 (2014).
- 709 18. Dormann, C.F. et al. Collinearity: a review of methods to deal with it and a simu-
710 lation study evaluating their performance. *Ecography* 35, 1– 20 (2012).
- 711 19. Wu, Q. et al. GLUT1 inhibition blocks growth of RB1-positive triple negative
712 breast cancer. *Nat Commun.* 11, 4205 (2020).
- 713 20. Wang, L. et al. Novel RNA-Affinity Proteogenomics Dissects Tumor Heterogenei-
714 ty for Revealing Personalized Markers in Precision Prognosis of Cancer. *Cell*
715 *Chem Biol.* 25, 619–633 (2018).

- 716 21. Zhang, H. et al. NCBP1 promotes the development of lung adenocarcinoma
717 through up-regulation of CUL4B. *J Cell Mol Med.* 23, 6965–6977 (2019).
- 718 22. Shen, H. et al. Nuclear expression and clinical significance of phosphohistidine
719 phosphatase 1 in clear-cell renal cell carcinoma. *J Int Med Res.* 43, 747–757
720 (2015).
- 721 23. Zhang, B. et al. Proteogenomic characterization of human colon and rectal can-
722 cer. *Nature* 513, 382–387 (2014).
- 723 24. Sørli, T. et al. Gene expression patterns of breast carcinomas distinguish tumor
724 subclasses with clinical implications. *Proc Natl Acad Sci USA.* 98, 10869–10874
725 (2001).
- 726 25. Rycaj, K. & Tang, D. G. Cell-of-Origin of Cancer versus Cancer Stem Cells: As-
727 says and Interpretations. *Cancer Res.* 75, 4003–4011 (2015).
- 728 26. Hoadley, K. A. et al. Cell-of-Origin Patterns Dominate the Molecular Classifica-
729 tion of 10,000 Tumors from 33 Types of Cancer. *Cell* 173, 291–304 (2018).
- 730 27. Visvader J. E. Cells of origin in cancer. *Nature* 469, 314–322 (2011).
- 731 28. Bhat-Nakshatri, P. et al. A single-cell atlas of the healthy breast tissues reveals
732 clinically relevant clusters of breast epithelial cells. *Cell Rep Med.* 2, 100219
733 (2021).
- 734 29. Frezza, C. Metabolism and cancer: the future is now. *Br J Cancer* 122, 133–135
735 (2020).
- 736 30. Ghaffari, P., Mardinoglu, A. and Nielsen, J. Cancer Metabolism: A Modeling Per-
737 spective. *Front Physiol.* 6, 382 (2015).

- 738 31. Martínez-Reyes, I., Chandel, N.S. Cancer metabolism: looking forward. *Nat Rev*
739 *Cancer* 21, 669–680 (2021).
- 740 32. Levine, A.J., and Puzio-Kuter A.M. The control of the metabolic switch in cancers
741 by oncogenes and tumor suppressor genes. *Science* 330, 6009 (2010).
- 742 33. Budczies, J. et al. Comparative metabolomics of estrogen receptor positive and
743 estrogen receptor negative breast cancer: alterations in glutamine and beta-
744 alanine metabolism. *J. Proteomics* 94, 279-288 (2013).
- 745 34. Allison, K.H. et al. Estrogen and Progesterone Receptor Testing in Breast Can-
746 cer: ASCO/CAP Guideline Update. *J Clin Oncol.* 38, 1346–1366 (2020).
- 747 35. Goldhirsch, A. et al. Strategies for subtypes--dealing with the diversity of breast
748 cancer: highlights of the St. Gallen International Expert Consensus on the Prima-
749 ry Therapy of Early Breast Cancer 2011. *Ann Oncol.* 22, 1736-1747 (2011).
- 750 36. Sturtz, L.A. et al. Comparative analysis of differentially abundant proteins quanti-
751 fied by LC-MS/MS between flash frozen and laser microdissected OCT-
752 embedded breast tumor samples. *Clin Proteomics* 17, 40 (2020).
- 753 37. Scrucca, L. et al. mclust 5: Clustering, Classification and Density Estimation Us-
754 ing Gaussian Finite Mixture Models. *R J.* 8, 289–317 (2016).
- 755 38. Benaglia, T. et al. mixtools: An R package for analyzing finite mixture models. *J*
756 *Stat Softw.* 32, 1-29 (2009).
- 757 39. Gu, Z. et al. Complex heatmaps reveal patterns and correlations in multidimen-
758 sional genomic data. *Bioinformatics* 32, 2847–2849 (2016).
- 759 40. Cancer Genome Atlas Network. Comprehensive molecular portraits of human
760 breast tumours. *Nature* 490, 61-70 (2012).

- 761 41. Huo, D. et al. Comparison of Breast Cancer Molecular Features and Survival by
762 African and European Ancestry in The Cancer Genome Atlas. *JAMA Oncol.* 3,
763 1654-1662 (2017).
- 764 42. Liu, J. et al. An Integrated TCGA Pan-Cancer Clinical Data Resource to Drive
765 High-Quality Survival Outcome Analytics. *Cell* 173, 400–416 (2018).
- 766 43. Curtis, C. et al. The genomic and transcriptomic architecture of 2,000 breast tu-
767 mours reveals novel subgroups. *Nature* 486, 346–352 (2012).
- 768 44. Pereira, B. et al. The somatic mutation profiles of 2,433 breast cancers refines
769 their genomic and transcriptomic landscapes. *Nat Commun.* 7, 11479 (2016).
- 770 45. Brueffer, C. et al. Clinical Value of RNA Sequencing-Based Classifiers for Predic-
771 tion of the Five Conventional Breast Cancer Biomarkers: A Report From the
772 Population-Based Multicenter Sweden Cancerome Analysis Network-Breast Initi-
773 ative. *JCO Precis Oncol* 2, PO.17.00135 (2018).
- 774 46. Ritchie, M.E. et al. limma powers differential expression analyses for RNA-
775 sequencing and microarray studies. *Nucleic Acids Res.* 43, e47 (2015).
- 776 47. Smyth, G.K. et al. limma: Linear Models for Microarray and RNA-Seq Data User's
777 Guide. *R package version 3.38.3* (2019).
- 778 48. Huang, D., Sherman, B. T., and Lempicki, R. A. Systematic and integrative anal-
779 ysis of large gene lists using DAVID bioinformatics resources. *Nat Protoc* 4, 44–
780 57 (2009).
- 781 49. Clark, T. G. et al. Survival analysis part IV: further concepts and methods in sur-
782 vival analysis. *British journal of cancer* 89, 781–786 (2003).

- 783 50. Mandrekar, J. N. et al. Cutpoint Determination Methods in Survival Analysis us-
784 ing SAS. *Proceedings of the 28th SAS Users Group International Conference*
785 *(SUGI)* 261-28 (2003).
- 786 51. Therneau, T. *A Package for Survival Analysis in R*. R package version 3.2-11
787 (2021).
- 788 52. Kassambara, A. et al. survminer: Survival Analysis and Visualization. R package
789 version 0.4.9 (2021).
- 790 53. Monti, S. et al. Consensus Clustering: A Resampling-Based Method for Class
791 Discovery and Visualization of GeneExpression Microarray Data. *Machine Learn-*
792 *ing*, 52, 91–118 (2003).
- 793 54. Wilkerson, D. M. and Hayes, N.D. ConsensusClusterPlus: a class discovery tool
794 with confidence assessments and item tracking. *Bioinformatics* 26, 1572-1573
795 (2010).
- 796 55. GebSKI, V. et al. Data maturity and follow-up in time-to-event analyses. *Int. J. Ep-*
797 *idemiol* 47, 850–859 (2018).
- 798 56. Yu, G. et al. clusterProfiler: an R package for comparing biological themes
799 among gene clusters. *OMICS* 16, 284-287 (2012).

800 **Acknowledgements**

801 We thank the patients who participated in this study. We thank Carl Berg for research
802 support, and Joseph A. Hooke from Wake Forest University for developing a reference
803 library for our work.

804 **Author contributions**

805 A.J.K., M.A.K., H.H., J.A.H., M.L.C., N.R.N., J.L.F.-C. and GW designed research study.
806 P.S. and R.M.S. conducted experiments. A.J.K., B.D., R.N.Z., B.M., L.K., S.S. and
807 L.A.S. acquired data. A.J.K., G.W., M.A.K., P.S., J.L., P.K.-R.K. performed data anal-
808 yses. P.S., R.M.S. and L.A.S. provided reagents. G.W. and P.S. wrote the draft manu-
809 script. A.J.K., M.A.K., H.H., N.R.N., R.S., J.A.H., M.L.C., R.N.Z., J.L.F.-C., C.D.S. and
810 G.W. reviewed and edited the manuscript. A.J.K. and H.H. administrated project. C.D.S.
811 and N.R.N. supported funding acquisition. All authors reviewed and approved the final
812 version of the manuscript.

813 **Competing interests**

814 Michael A. Kiebish, Rangaprasad Sarangarajan, Punit Shah, Rick M. Searfoss, Elder
815 Granger, Niven R. Narain are employees of BERG, LLC. Parts of this work are the sub-
816 ject of a patent application: 119992_22501.

817 **Materials & Correspondence**

818 Further information and requests for resources will be fulfilled by the corresponding au-
819 thor, Albert J. Kovatich (akovatich@hifresearch.org).

820

821 **Funding and support**

822 The study was supported by DOD-BERG Breast Cancer CRADA from the
823 U.S.Department of Defense.

824 **Disclaimer**

825 The contents of this publication are the sole responsibility of the author(s) and do not
826 necessarily reflect the views, opinions or policies of Uniformed Services University of
827 the Health Sciences (USUHS), The Henry M. Jackson Foundation for the Advancement
828 of Military Medicine, Inc., the Department of Defense (DoD), the Departments of the
829 Army, Navy, or Air Force. Mention of trade names, commercial products, or organiza-
830 tions does not imply endorsement by the U.S. Government.

831

832

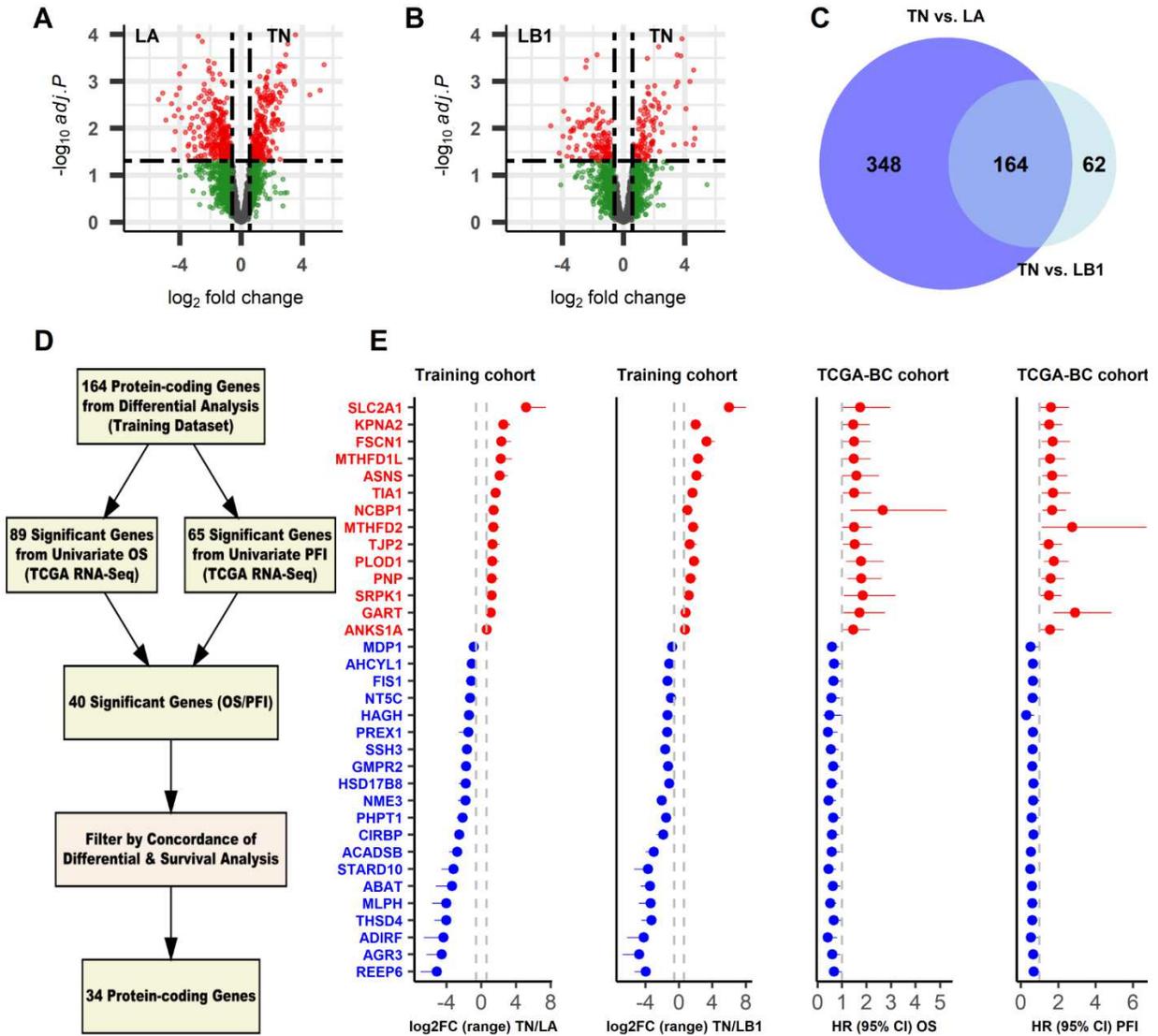
833 **Main Tables**834 **Table 1: Demographic and clinicopathological characteristics of our study cohort**

	Training (N=70)	Testing (N=39)	Low ER+ (N=7)	Total (N=116)
Race				
African American	11 (15.7%)	2 (5.1%)	1 (14.3%)	14 (12.1%)
White	56 (80.0%)	32 (82.1%)	5 (71.4%)	93 (80.2%)
Asian	2 (2.9%)	1 (2.6%)	0 (0%)	3 (2.6%)
Other	1 (1.4%)	4 (10.3%)	1 (14.3%)	6 (5.2%)
Age at Diagnosis (years)				
Mean (SD)	58.6 (11.3)	57.7 (14.1)	49.9 (9.89)	57.8 (12.3)
Median [Min, Max]	59.5 [34, 86]	54 [30, 85]	53 [35, 66]	57 [30, 86]
Grade				
G1	10 (14.3%)	9 (23.1%)	0 (0%)	19 (16.4%)
G2	31 (44.3%)	13 (33.3%)	1 (14.3%)	45 (38.8%)
G3	27 (38.6%)	17 (43.6%)	6 (85.7%)	50 (43.1%)
Missing	2 (2.9%)	0 (0%)	0 (0%)	2 (1.7%)
Tumor Size (mm)				
Mean (SD)	25.4 (10.5)	23.5 (10.1)	21.6 (6.24)	24.5 (10.1)
Median [Min, Max]	24 [7, 55]	20 [11, 51]	23 [12, 28]	22 [7, 66]
Lymph Node Status				
Negative	42 (60%)	20 (51.3%)	5 (71.4%)	67 (57.8%)
Positive	26 (37.1%)	19 (48.7%)	2 (28.6%)	47 (40.5%)
Missing	2 (2.9%)	0 (0%)	0 (0%)	2 (1.7%)
ER Status				
Negative	9 (12.9%)	6 (15.4%)	0 (0%)	15 (12.9%)
Positive	61 (87.1%)	33 (84.6%)	0 (0%)	94 (81%)
Low ER+	0 (0%)	0 (0%)	7 (100%)	7 (6%)
PR Status				
Negative	19 (27.1%)	7 (17.9%)	6 (85.7%)	32 (27.6%)
Positive	51 (72.9%)	32 (82.1%)	1 (14.3%)	84 (72.4%)
HER2 Status (Negative)				
0	11 (15.7%)	8 (20.5%)	1 (14.3%)	20 (17.2%)
1+	47 (67.1%)	23 (59%)	5 (71.4%)	75 (64.7%)
2+	12 (17.1%)	8 (20.5%)	1 (14.3%)	21 (18.1%)
Ki-67 Status				
Negative	18 (25.7%)	10 (25.6%)	2 (28.6%)	30 (25.9%)
Positive	46 (65.7%)	26 (66.7%)	5 (71.4%)	77 (66.4%)
unknown	6 (8.6%)	3 (7.7%)	0 (0%)	9 (7.8%)
IHC Subtype				
LA	20 (28.6%)	10 (25.6%)	2 (28.6%)	32 (27.6%)
LB1	41 (58.6%)	23 (59%)	5 (71.4%)	69 (59.5%)
TN	9 (12.9%)	6 (15.4%)	0 (0%)	15 (12.9%)

835 ER: Estrogen receptor; PR: progesterone receptor; HER2: human epidermal growth factor receptor 2;

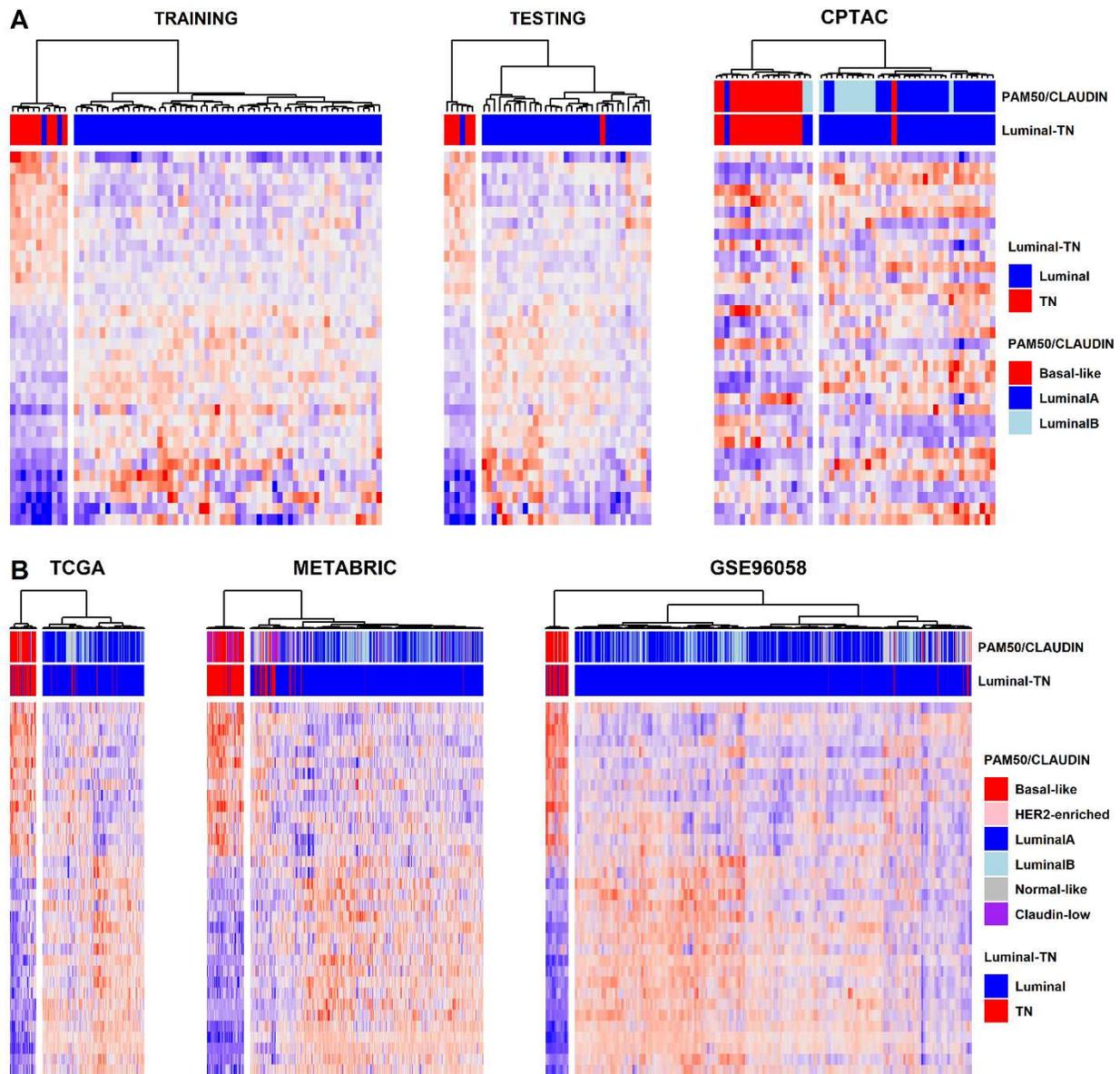
836 LA: ER+/HER2-/Ki-67-; LB1: ER+/HER2-/Ki-67+; TN: Triple-negative

837 **Main figures, titles and legends**



838

839 **Figure 1. LT34 proteomic biomarker panel identification.** 840 **A-B** Volcano plots showing the consistently 841 differential analysis results of the comparison between IHC-based TN and LA subtypes (**A**) and the com- 842 parison between TN and LB1 (**B**) subtypes separately from our training dataset. The significantly altered 843 proteins shown in red were reported at $FDR < 0.05$ and ($FC > 1.5$ or $FC < 0.667$) consistently across 101 dif- 844 ferential analyses. **C** Venn Diagram showing 164 consistently significantly altered proteins detected from 845 TN versus Luminal (LA, LB1). **D** Workflow showing the steps filtering 164 protein-coding genes corre- 846 sponding to 164 significantly altered proteins to 34 protein-coding genes from TCGA transcriptomic data. 847 **E** Forest plots showing $\log_2(\text{fold change})$ from the differential analyses from our training dataset, as well 848 as hazard ratio of 34 protein-coding genes from Cox proportional hazard model using TCGA HER2-co-

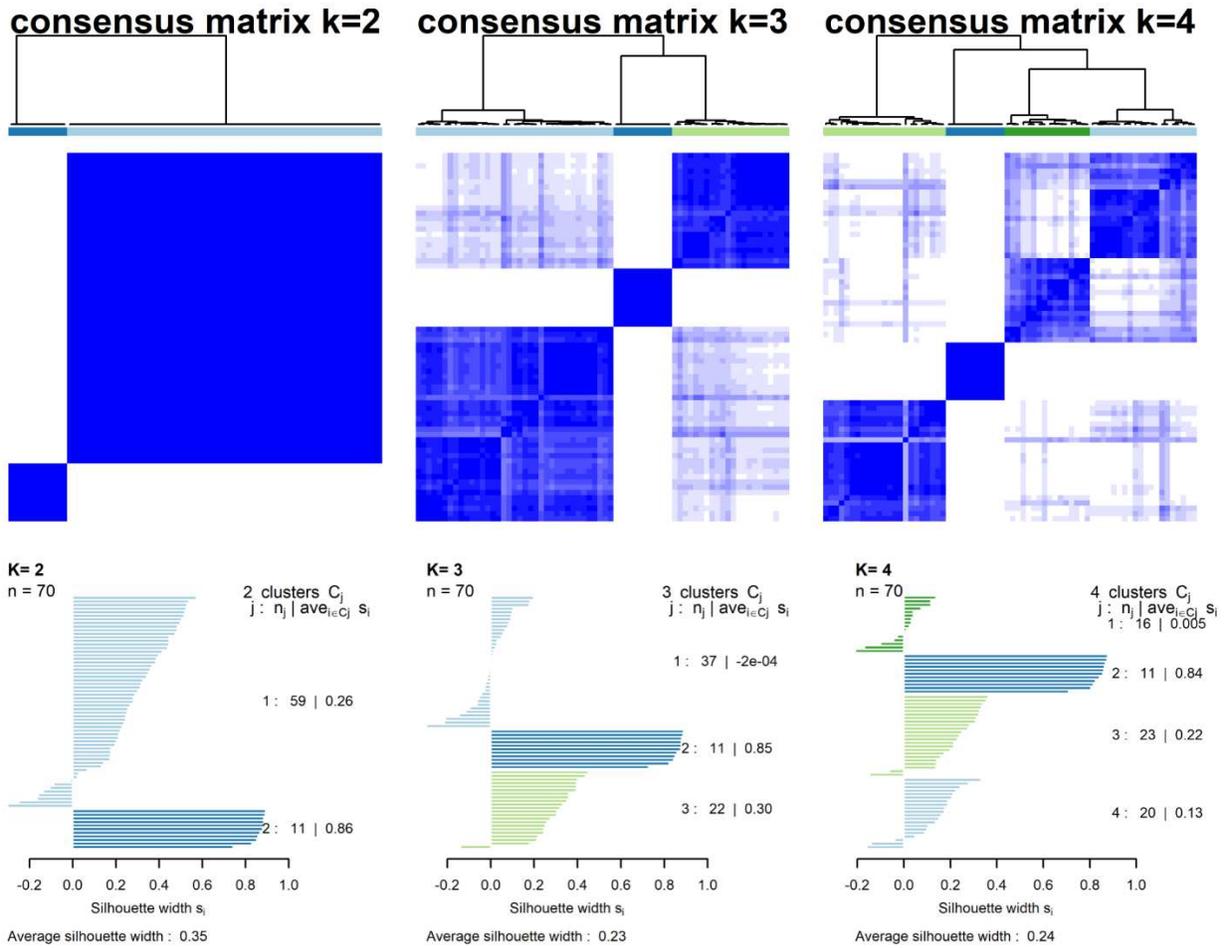


849

850 **Figure 2. Hierarchical clustering heatmaps across cohorts using 34 proteins/genes.** **A** Hierarchical
 851 clustering heatmaps for our internal training cohort (70 cases), our internal testing cohort (39 cases) and
 852 CPTAC HER2- cases (53 cases) using 34 proteins. **B** Hierarchical clustering heatmaps of TCGA HER2-
 853 cohort (799 cases in RNA-seq data), METABRIC HER2- cohort (1645 cases in Microarray data) and
 854 GSE96058 HER2- cohort (2435 cases in RNA-seq data) using 34 coding-genes. The heatmaps demon-
 855 strating that two distinct clusters were derived from both proteomic and transcriptomic platforms using 34
 856 proteins/genes.

857

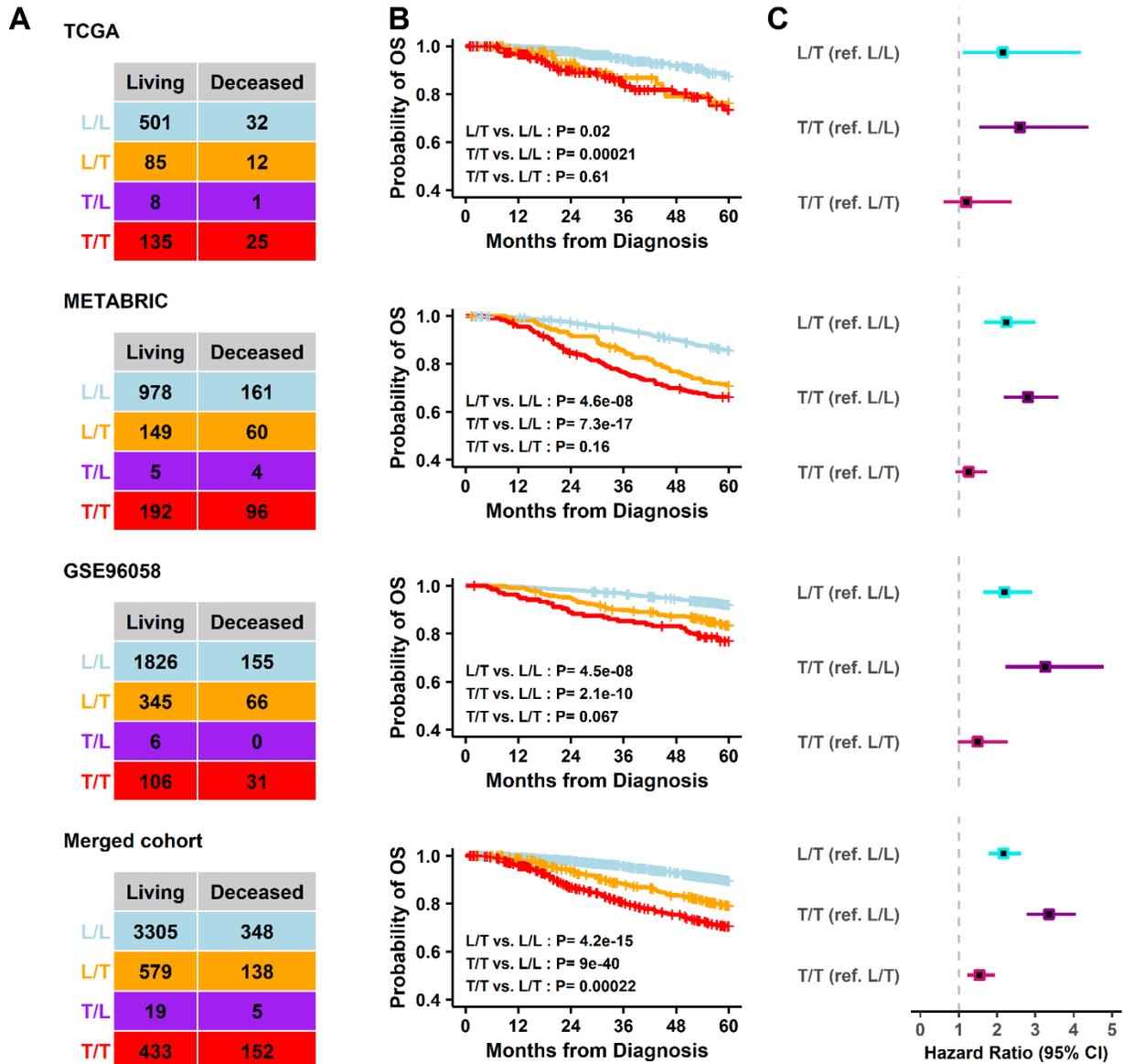
858



859

860 **Figure 3. Consensus clustering analysis for training cohort using 34 proteins.** Two novel proteomic
861 subtypes (LT34) were clearly identified using consensus clustering analysis with 34 proteins from training
862 cohort. One cluster was defined as a TN-like subtype, another one was defined as a Luminal-like subtype
863 based on Fisher's exact test.

864



866

867

868

869

870

871

872

873

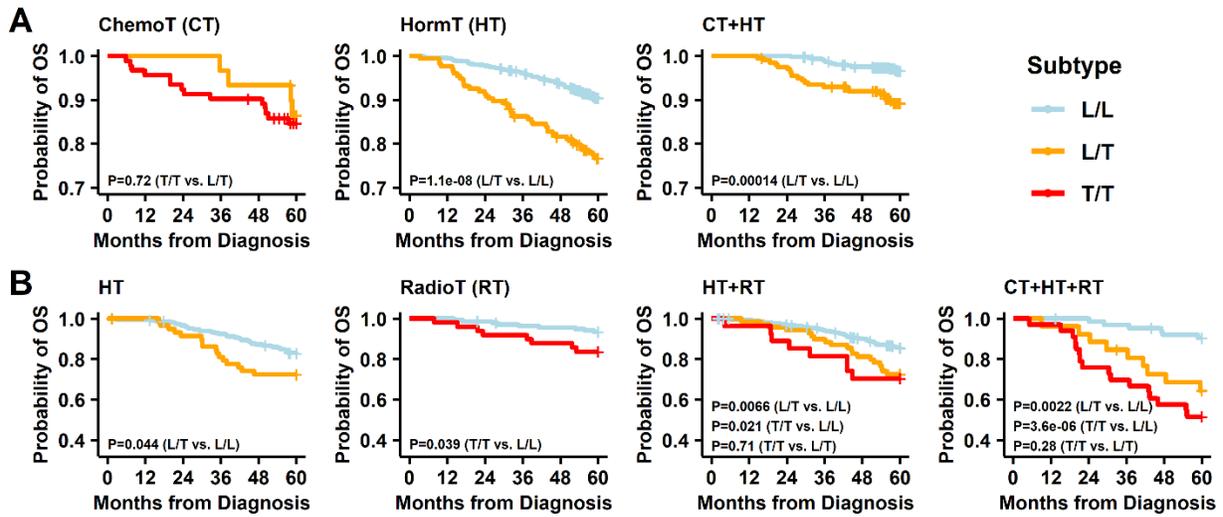
874

875

876

Figure 4. OS differences by IHC-LT34 subtypes. **A** The contingency tables between IHC-LT34 subtypes and living status for TCGA, METABRIC, GSE96058 and the merged cohort showing more percentages of L/T subtype patients were deceased compared with the percentage of L/L subtype in each cohort respectively. **B** OS K-M plots among L/L, L/T and T/T subtypes in Luminal-TN cohort without low ER+ cases demonstrating that T/T tumors had the worst outcome whereas L/L had the most favorable outcome, and L/T tumors had a statistically significant worse outcome comparing with L/L tumors (p -value < 0.05), however, the survival difference between T/T and L/T tumors is not statistically significant except in the merged cohort. **C** The hazard ratio forest plots corresponding to each K-M plot and the hazard ratios were calculated using Cox Proportional Regression Model and also shown in Supplementary Table S7.

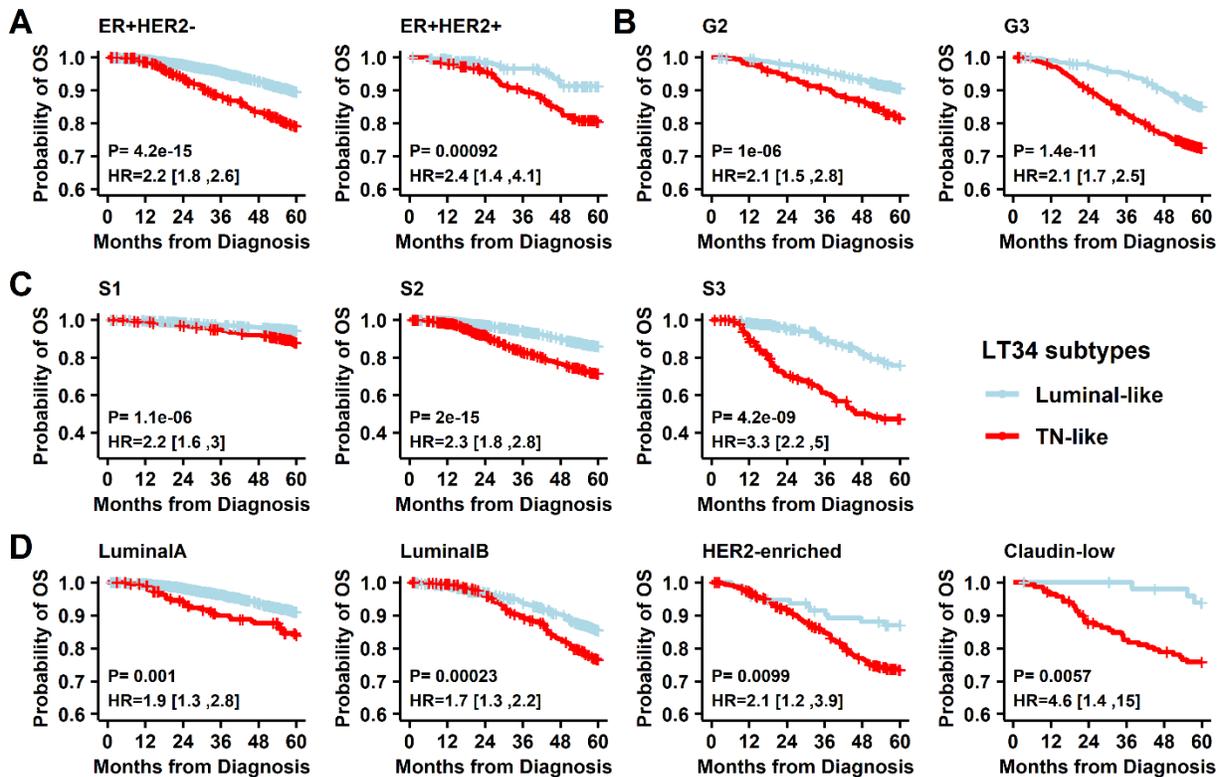
877



878

879 **Figure 5. K-M plots by IHC-LT34 subtypes within each treatment.** K-M plots of IHC-LT34 subtype un-
880 der each treatment in GSE96058 cohort (**A**) and METABRIC cohort (**B**). Only survival curve that passed
881 our data maturity criteria (see Supplementary Table S6) were shown. They all demonstrate that the L/T
882 subtype patients were still associated with poor survival compared with L/L subtype patients under each
883 treatment and imply that L/T subtype patients were resistant to the provided treatments compared to L/L
884 subtype patients. The L/T subtype has a similar OS as the T/T subtype compared to the L/L subtype.

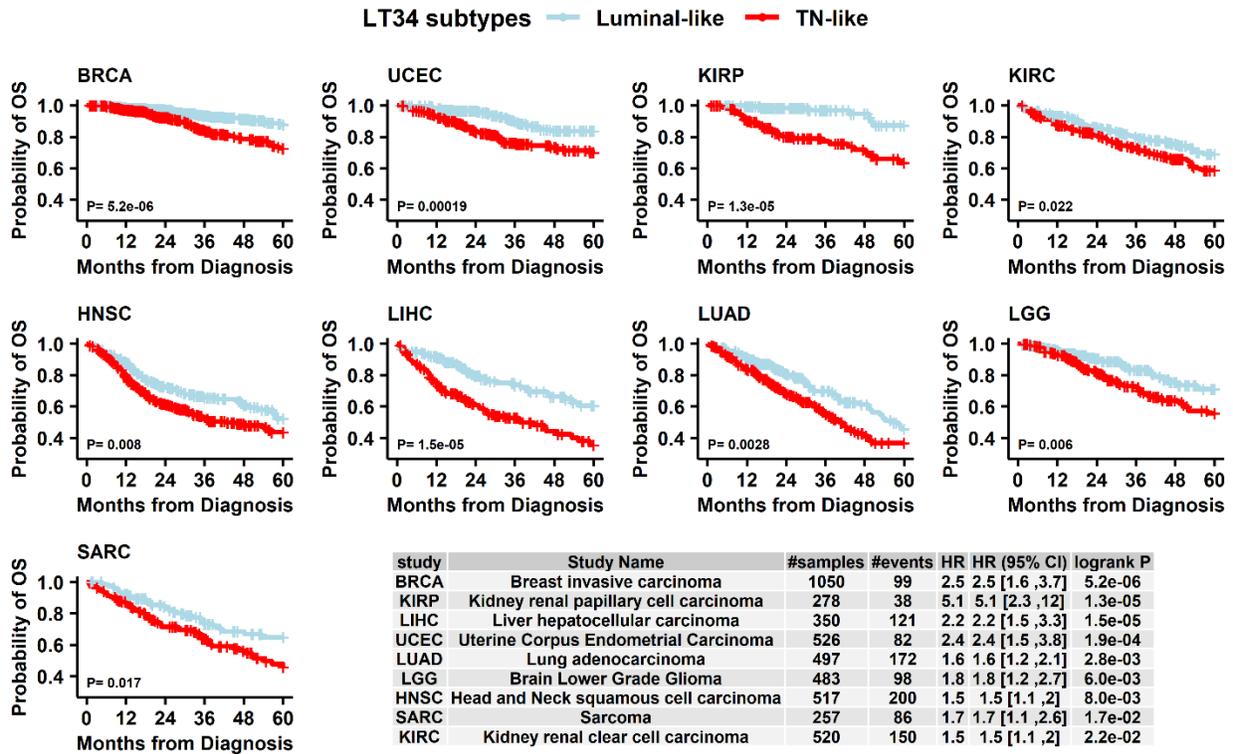
885



886

887 **Figure 6. K-M plots by LT34 subtypes within each clinical group.** K-M plots of LT34 subtype within
 888 each clinical group: IHC-based subtype (A), grade (B), stage (C) and PAM50 or Claudin-low subtype (D)
 889 respectively in the merged cohort (TCGA + METABRIC + GSE96058). Only survival curves that passed
 890 our data maturity criteria were shown (see Supplementary Table S6). They all demonstrated that there
 891 was a significant OS difference between TN-like subtype patients and Luminal-like subtype patients and
 892 TN-like subtype patients were associated with poor OS compared with Luminal-like subtype patients.

893



895

896 **Figure 7. K-M plots by LT34 subtypes within each TCGA cancer significantly associated with sur-**
 897 **vival.** K-M plots of LT34 subtype within 9 TCGA cancers. Only K-M plots with log-rank p-value <0.05 and
 898 survival curves that passed our data maturity criteria are shown (see Supplementary Table S6). They
 899 demonstrate that there is a significant OS difference between TN-like subtype patients and Luminal-like
 900 subtype patients in each of 9 cancers, and TN-like subtype patients are associated with poorer OS com-
 901 pared to Luminal-like subtype patients.

Supplementary Files

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

- [SupplementaryTableS4.xlsx](#)
- [SupplementaryTableS9.xlsx](#)
- [SupplementaryTableS2.xlsx](#)
- [SupplementaryTableS8.xlsx](#)
- [SupplementaryTableS7.xlsx](#)
- [SupplementaryTableS1.xlsx](#)
- [SupplementaryTableS3.xlsx](#)
- [ManuscriptSupplementaryinformation.pdf](#)
- [SupplementaryTableS12.xlsx](#)
- [SupplementaryTableS6.xlsx](#)
- [SupplementaryTableS10.xlsx](#)
- [SupplementaryTableS5.xlsx](#)
- [SupplementaryTableS11.xlsx](#)