

1 **MIRS: an AI scoring system for predicting the prognosis**
2 **and therapy of breast cancer**

3

4

5 **Chen Huang^{1,3†}, Min Deng^{2†}, Dongliang Leng^{2†}, Elaine Lai-Han Leung^{1,3},**

6 **Baoqing Sun⁴, Peiyan Zheng⁴, Xiaohua Douglas Zhang^{2*}**

7

8 ¹ Dr. Neher's Biophysics Laboratory for Innovative Drug Discovery, Macau

9 University of Science and Technology, Macau, SAR, China

10 ² CRDA, Faculty of Health Sciences, University of Macau, Taipa, Macau

11 ³ Stat Key laboratory of Quality Research in Chinese Medicine, Macau Institute For

12 Applied Research in Medicine and Health, Macau University of Science and

13 Technology, Macau, SAR, China

14 ⁴ Department of Allergy and Clinical Immunology, State Key Laboratory of

15 Respiratory Disease, National Clinical Research Center of Respiratory Disease,

16 Guangzhou Institute of Respiratory Health, First Affiliated Hospital of Guangzhou

17 Medical University, Guangzhou, Guangdong, China

18 [†]These authors contributed equally to this work

19 *Correspondence: Xiaohua Douglas Zhang douglaszhang@um.edu.mo

20

21

22 **Abstract**

23 Current scoring systems for prognosis of breast cancer are available but usually
24 consider only one prognostic feature. We aim to develop a novel prognostic scoring
25 system based on both immune-infiltration and metastatic features to not only assess the
26 patient prognoses more accurately but also guide therapy for patients with breast cancer.
27 Computational immune-infiltration and gene profiling analysis identified a 12-gene
28 panel firstly characterizing immune-infiltrating and metastatic features. Neural network
29 model yielded a precise prognostic scoring system called metastatic and
30 immunogenomic risk score (MIRS). The influence of MIRS on the prognosis and
31 therapy of breast cancer was then comprehensively investigated. MIRS significantly
32 stratifies patients into high risk-group (MIRS^{high}) and low risk-group (MIRS^{low}) in both
33 training and test cohorts. The MIRS^{low} patients exhibit significantly improved survival
34 rate compared with MIRS^{high} patients. A series of analyses demonstrates that MIRS can
35 well characterize the metastatic and immune landscape of breast cancer. Further
36 analysis on the usage of MIRS in chemotherapy suggests that MIRS^{high} patients may
37 benefit from three chemotherapeutic drugs (Cisplatin, Tamoxifen and Vincristine).
38 Higher immune infiltration and significantly prolonged survival are observed in
39 MIRS^{low} patients, indicating a better response in immune checkpoint inhibitor therapy.
40 Our analysis demonstrates that MIRS could effectively improve the accuracy of
41 prognosis for patients with breast cancer. Also, MIRS is a useful webtool, which is
42 deposited at <https://lva85.github.io/MIRS/>, to help clinicians in designing personalized
43 therapies for patients with breast cancer.

44 **Keywords:** Breast cancer, Metastasis, Immune infiltration, Prognosis, Personalized
45 treatment

46 **Introduction**

47 Cancer has long history in mankind and remains the leading cause of death, with
48 breast cancer being one of the most common malignancies in women worldwide (1, 2).
49 Breast cancer is also the second most common cause of death in cancer-related deaths
50 among women. (3, 4). Despite tremendous advancement of medicine over the years has
51 lowered the mortality rate, the high level of heterogeneity in breast cancer still makes
52 the prognosis and treatment challenging.

53 Over the decade, a considerable amount of work has been done to develop
54 prognostic measures on the progression of breast cancer (5). The majority (~80%) of
55 breast cancer becomes invasive (6) and approximately 20~30% of them results in
56 distant metastasis after treatment (7). Metastasis is thereby the most fatal development
57 of breast cancer, which greatly reduces the rate of long-term survival from 90% to 5%
58 (8). However, most metastasis-based signatures were developed based on organ-
59 specific metastatic events, yet breast cancer consists of tumors with extremely
60 heterogeneous cell types, resulting in the discrepancy between prognosis and survival
61 (9, 10). Hence currently available metastasis-based prognostic measures have poor
62 performance (11). On the other side, tumor-infiltrating lymphocytes have already been
63 reported to be inextricably linked to therapeutic efficacy and patient survival in various
64 cancers (12, 13). Many prognostic predictors were developed by assessing the level of

65 the infiltration of immune cells into tumor and were preferably adopted for prognosis
66 in cancers (14, 15). These histological strategies based on the analysis of a small
67 proportions of immune cell marker genes support the prognostic significance of
68 immune infiltration but still have limitations. Firstly, strategies for describing the level
69 of immune infiltrate are the first limitation of the current studies (16). Specifically, each
70 immune cell subset is computationally estimated by reference profiles based on bulk
71 analysis of tissue samples. This is the main drawback because the transcriptional
72 program of immunocytes exhibits high plasticity under tumor microenvironments (17).
73 Secondly, while most studies were used the immune-related characteristics to improve
74 cancer prognosis, only one or two subsets of immunocyte are included and these subsets
75 lack functional variation, thus the treatments based on these indicators fail to achieve
76 satisfactory immune response effects (18). Therefore, prognostic indicator based on
77 only one characteristic without considering other crucial features is insufficient to
78 accurately assess risk stratification and direct treatment strategies.

79 Given the limitations of the aforementioned work a more comprehensive approach
80 should be developed to assess prognostic value and translate it into clinical practice.
81 For the first time, we develop a prognostic signature for breast cancer patients,
82 integrating immune-related gene signatures involved in metastasis, to classify patients
83 with breast cancer into groups of high and low risk for potential therapeutic strategies.
84 We construct a Neuron network to estimate gene weights, which exhibit outstanding
85 performance in binary classification. A metastatic and immunogenomic risk score
86 (MIRS) is then established, which has conspicuous power to predict survival status

87 compared with previously published indicators based on single feature. Ultimately, the
88 ability of MIRS to predict the treatment is identified, suggesting its potential to guide
89 therapeutic tactics in breast cancer.

90 **Materials and Methods**

91 **Collection and pre-processing of breast cancer data**

92 All analyzed expression profiles and the corresponding clinical datasets were
93 collected from Gene Expression Omnibus (GEO), The Cancer Genome Atlas (TCGA,
94 <https://www.cbioportal.org/datasets>), and Molecular Taxonomy of Breast Cancer
95 International Consortium (METABRIC, <https://www.cbioportal.org/datasets>). Only the
96 datasets available with sufficient overall survival information were included, consisting
97 of 8,424 patients from 14 cohorts. The detailed information of each cohort is presented
98 in Table S1 and S2.

99 Raw series matrix files generated by Affymetrix were downloaded from GEO
100 database. The R package GEOquery (19) was used to process raw matrix data.
101 Duplicated genes detected by multiple probes were retained by taking the maximum
102 expression value of the probe sets. Gene expression value was normalized by log₂
103 transformation. Each GEO and RNA-seq dataset were processed independently.

104 **Construction of immune cell infiltration groups**

105 A set of biomarkers is derived from Charoentong et al (20), comprising 45 immune
106 signatures related to immune cell types, immunogenomic pathways and functions. The

107 concrete gene signatures for each immune cell type were obtained from (21), and the
108 immune-related pathways and functions were downloaded from database ‘ImmPort’
109 (22). Single sample gene set enrichment analysis (ssGSEA) implemented in R package
110 GSVA was used to quantify the infiltration level of different immune cells,
111 immunogenomic pathways and the activity of immune-related functions via expression
112 data of breast cancer (23). Based on the results of ssGSEA, patients in TCGA breast
113 cancer cohort (TCGA-BRCA) were divided into high and low immune cell infiltration
114 groups using hierarchical clustering analysis (Figure S1) (24).

115 **Identification of immune and metastatic candidate** 116 **genes**

117 Using Wilcoxon rank-sum (Wilcoxon) test, the differentially expressed genes
118 (DEGs) related to tumor immune infiltration were detected from high and low immune
119 infiltration conditions according to the filtering criteria $|\log_2FC| > 0.5$ and adjusted $p <$
120 0.05 using Benjamini and Hochberg (BH) method (25). Meanwhile, utilizing the
121 Wilcoxon test with the same criteria in the comparison between metastasis and primary
122 breast cancer groups from the union of GSE10893 and GSE3521, the DEGs involved
123 in metastatic mechanism were then identified. For these two DE analyses, Venn
124 analysis found 52 metastatic and immunogenomic candidate genes. The heatmap of
125 these DEGs are visualized in Supplementary Figures S3-S4.

126 **Establishment of prognostic risk score**

127 Univariate Cox proportional hazard regression analysis was designed to screen
128 features related to overall survival (OS) from 52 candidate genes in TCGA BRCA
129 cohort. The filtered gene list is provided in Table S2. Subsequently, only the genes with
130 absolute Hazard ratio (HR) larger than 1 and p-value less than 0.05 were retained. To
131 eliminate collinearity, the eligible candidate genes were further filtered depending on
132 the criteria that the square root of Variance Inflation Factor (VIF) was less than 2 and
133 the Pearson Correlation Coefficient was smaller than 0.5. Ultimately, 12 prognostic
134 genes that were significantly correlated with patients' OS were identified.

135 These 12 prognostic signatures were classified into binary status. One was defined
136 as the protective status in which HR was less than 1 whereas another was the dangerous
137 status in which the corresponding HR was greater than 1. The expression status of each
138 protective mRNA was assigned as 1 if the expression level of this mRNA was above
139 the median of the expression values of all samples, otherwise it would be assigned as
140 0. In contrast, the expression of dangerous mRNA was assigned as 1 if it had expression
141 value below median, otherwise assigned as 0. This approach not only allows the risk
142 score, which is based on protective and dangerous genes, to simultaneously contribute
143 to consistent survival outcome, but also avoids the influence of inconsistent sequencing
144 platforms. To date, several machine learning methods were found to be successful in
145 various data mining problems, including those with transcriptomic data (26, 27).
146 Therefore, a multilayer perceptron neuron network was built to estimate the weights of
147 the 12 prognostic genes. In the Figure S4, the $net_{n_1} = W_{1,1}i_1 + W_{2,1}i_2 + \dots +$

148 $W_{12,1}i_{12} + b_1$ was defined, where W is the weight of each input node and i_j ($j =$
149 $1, 2, \dots, 12$) is the '0-1' status of gene. Then we exploited rectified linear unit (ReLU):

150

$$151 \quad \text{ReLU}(net_{n_1}) = Out_{n_1} = \begin{cases} net_{n_1}, & net_{n_1} > 0 \\ 0, & net_{n_1} \leq 0 \end{cases},$$

152

153 and $net_{01} = W'_{1,1} * Out_{n_1} + W'_{2,1} * Out_{n_2} + W'_{3,1} * Out_{n_3} + W'_{4,1} * Out_{n_4} + b_2$
154 as an activation function in the hidden layer. In the output layer, we applied the Softmax
155 function to each node and designated probability of death:

$$156 \quad \text{softmax}(net_{01}) = Out_{01} = \frac{e^{net_{01}}}{e^{net_{01}} + e^{net_{02}}} \in (0,1).$$

157 We then created two nodes $a_0 = 0$ and $a_1 = 1$ for alive and dead, respectively.
158 Cross entropy error is computed as:

$$159 \quad E = \sum_{i=1}^N E_i = -a_0^i * \log(out_{01}^i) - a_1^i * \log(out_{02}^i), \quad \text{where } i \text{ is } i\text{th sample.}$$

160 Finally, the value of each weight was optimized by minimizing E using gradient
161 descent. The R packages Tensorflow and Keras were employed to construct neuron
162 network. After training, the coefficient of each prognostic gene was then determined as
163 the maximum weight in the hidden layer (26).

164 Lastly, the risk score that consists of 12 metastatic and immunogenomic prognostic
165 genes (MIRS) for each patient is defined as the following:

$$166 \quad MIRS_i = \sum_{j=1}^m \text{weight}_j \times I_{\{\text{protective gene } j\}} + \sum_{k=1}^n \text{weight}_k \times I_{\{\text{dangerous gene } k\}}$$

167 where m and n denote the number of protective and dangerous genes, respectively,
168 weight is the maximum weight from the hidden layer. Additionally, $I_{\{\text{protective gene } j\}}$
169 and $I_{\{\text{dangerous gene } k\}}$ denote the following indicator functions:

170

171

$$I_{\{\text{protective gene } j\}} = \begin{cases} 1, & \text{Protective gene } j < \text{Median expression value across all samples,} \\ 0, & \text{otherwise.} \end{cases}$$

$$I_{\{\text{dangerous gene } j\}} = \begin{cases} 1, & \text{Dangerous gene } k > \text{Median expression value across all samples,} \\ 0, & \text{otherwise.} \end{cases}$$

174 **Statistical analysis**

175 All statistical analyses were performed using R software. The R packages pheatmap
176 and ggplot2 were used to plot heatmap and other graphs. The R package forestplot was
177 used to draw forest plot. The pROC package was employed to generate the Receiver
178 Operating Characteristic (ROC) curve and calculate the Area Under Curve (AUC),
179 which was an indicator to evaluate the predictive performance of risk score.

180 Based on the risk score, breast cancer patients in the investigated cohort were
181 stratified into subtypes of high risk or low risk depending on whether the value of
182 $(MIRS \text{ in each patient}) / (\text{median of } MIRS \text{ in all patients})$ was greater or
183 less than 1. This stratification method allows reasonable comparisons between different
184 data platforms. OS curves were established by Kaplan–Meier survival (KM) curve
185 function ggsurvplot, as implemented in R package survminer, and the difference in
186 survival distributions between risk subgroups was estimated by two-side log-rank test.
187 Based on univariate Cox proportional hazard regression analysis, the targeted
188 prognostic genes which were significantly correlated with OS were disclosed and the
189 Hazard ratio (HR), 95% confident interval of HR and p-value were also evaluated.
190 Multivariate Cox proportional hazard regression model was implemented to assess
191 whether the risk score is an independent prognosis factor when compared with other
192 important clinical features. All statistical tests were considered significant with p-value

193 < 0.05.

194 Full details about data and methods descriptions, including data information, gene
195 set enrichment analysis and mutation landscape analysis.

196 **Results**

197 **Screening of candidate genes from three public datasets**

198 To obtain significant prognostic biomarkers in breast cancer, we proposed a
199 systematic scheme of bioinformatic analysis (Figure 1). Given that the processes of
200 metastasis and immune infiltration in tumor play various important roles in cancer
201 development, we hypothesize that the expression of genes which were associated with
202 metastasis and immune infiltration in tumor should be correlated to the OS of cancer
203 patients. We thereby identified prognostic signatures based on these two characteristics.
204 Concretely, using ssGSEA method, the expression profile of 1,100 patients from TCGA
205 cohort were used to construct groups of high and low immune cell infiltration. Then the
206 patients were classified into the high immune infiltration group and low immune
207 infiltration group (Figure 2A and Figure S1). Furthermore, to validate the reliability of
208 the above grouping tactic, we investigated the expression level of two immune-related
209 gene families between these two groups: *CDI* and *ILI*. As expected, the expression of
210 both immune-related gene families in the high immune infiltration group is
211 significantly higher than that in the low immune infiltration group (Figure 1B and
212 Figure S5). Additionally, compared with low immune cell infiltration group, high
213 immune cell infiltration group exhibits a higher fraction of immune cell, stromal cell

214 but lower tumor purity using ESTIMATE (28) algorithm (Figure 2C). Furthermore, we
215 found that high immune cell infiltration group had significantly higher proportions in
216 most immune cell types than low immune infiltration group (Figure 2D) using
217 CIBERSORT algorithm under the permutation test with 1000 times. These findings
218 support that our immune cell infiltration grouping is highly confident to be used in
219 downstream analyses. Next, 1,222 differentially expressed genes were identified via
220 differential expression (DE) analysis between these two groups, which represents a
221 high-confidence dataset of genes related to immune infiltration (Table S4).

222 On the other part, aimed at identification of metastasis-related candidates, DE
223 analysis between metastasis and primary patients with breast cancer were performed
224 using two GEO cohorts (GSE10893 and GSE3521). The reason why we only chose
225 these two GEO datasets is that they have relatively balanced sample sizes between the
226 metastasis and primary groups when compared with other cohorts (Table S1). For
227 instance, TCGA breast cancer cohort contains 1,165 primary individuals but only 23
228 metastatic individuals. There is no doubt that such an extremely imbalanced data would
229 lead to biased result in DE analysis. This step yielded a union of 2,159 DE genes from
230 the results of these two GEO datasets (Table S4). Finally, a total of 52 genes was
231 obtained by intersecting 1,222 immune-infiltration-related genes and 2,159 metastasis-
232 related genes (Figure 2E), which represents prognostic candidates associated with both
233 tumor-immune infiltration (Figure 2F) and metastasis (Figure 2G).

234 **Construction and validation of MIRS in breast cancer** 235 **cohorts**

236 Univariate Cox regression analyses were performed to estimate the prognostic
237 relationship between candidate genes and overall survival in TCGA cohort. Among
238 these 52 candidate genes, 15 genes with p-value less than 0.05 were selected for follow-
239 up study (Table S2). Given that too many redundant variables would result in
240 overfitting in the linear model, we employed the analyses of Variance Inflation Factor
241 and Pearson Correlation Coefficient to eliminate the redundant genes (Figure 3A and
242 3B). As a result, a panel of 12 genes is reserved to establish the predictive model.

243 The TCGA-BRCA data (N = 1100 patients) were randomly classified into training
244 data (N = 770 patients) and testing data (N = 330 patients) at a ratio of 7:3. We then
245 optimized the weights for each gene with Neuron network in the training TCGA data.
246 The MIRS for each patient was built by summation of $Weight \times$
247 $I_{\{protective\ or\ dangerous\ gene\}}$ of all 12 genes (Table 1). MIRS was initially used to
248 predict patient's survival status, which yielded great predictive performance with AUC
249 accuracy of 0.875 in the training TCGA cohort (Figure 3C). In addition, all the patients
250 were classified into MIRS^{high} group and MIRS^{low} group using the median value of
251 MIRS as risk cut-off. As shown in Figure 3D, patients in MIRS^{low} group had
252 significantly longer OS or disease-free survival (DFS) time than those in MIRS^{high}
253 group (log-rank p<0.001) (Figure 3D and Figure S6-A).

254 To further examine the robustness and feasibility of this MIRS model, a
255 comprehensive survival analysis with KM method was performed in three independent

256 testing cohorts. Notably, MIRS exhibited robust predictive capacity with AUC of 0.934,
257 0.901, and 0.904 in GSE96058, GSE86166 and GSE20685, respectively (Figure 3E and
258 3G, Supplementary Figure S6-C). Regarding the survival analyses, consistent with the
259 result of the training data, the patients who are divided into MIRS^{high} group have
260 significantly worse OS than those in MIRS^{low} group (Figure 3F, Figure 3H and
261 Supplementary Figure S6-B). These analyses indicated that MIRS had precisely
262 prognostic ability in breast cancer. The higher score of MIRS corresponds to poor
263 outcome, and the lower score of MIRS refers to favorable outcome.

264 **Correlation of MIRS with the metastatic and** 265 **immunogenomic landscape between the high and low** 266 **subtypes**

267 We want to further scrutinize the correlation of metastatic and immunogenomic
268 landscape with MIRS in breast cancer patients. Initially, we investigated the correlation
269 between MIRS and the fraction of immune cell, stromal cell, as well as tumor purity
270 via ESTIMATE in the GSE86166 cohort. The results showed that MIRS^{low} group had a
271 higher fraction of immune cell and stromal cell cell but a lower tumor purity (Figure
272 4A). Similar situations were observed in GSE96058 (Figure S7). Reasonably, a higher
273 fraction of immune cell and lower tumor purity reflects a high level of infiltrating T-
274 lymphocytes in the patients of MIRS^{low} group, which is consistent with previous
275 survival analysis.

276 Moreover, 730 genes were identified to be correlated to the 12 genes of MIRS

277 (Spearman Correlation Coefficient ≥ 0.4) using GSE86166, subsequently, functional
278 enrichment analysis achieved via METASCAPE, indicating various immune-related
279 processes and pathways were significantly enriched, including T cell activation,
280 Cytokine-cytokine receptor interaction and B cell activation (Figure 4B). This
281 observation discloses a strong correlation of MIRS with immune activity. Alternatively,
282 we applied ssGSEA analysis to evaluate the immune infiltration level in GSE86166
283 using 17 immune-related biological functions and pathways derived from the immune-
284 related database 'ImmPort' (22). The result illustrates that most of the 17 items show
285 significant difference between MIRS^{high} and MIRS^{low} group (Figure 4C). Notably, all
286 immune-related biological processes and pathways exhibit significantly higher level of
287 immune infiltration in MIRS^{low} group (Figure 4C), which is consistent with our
288 previous analysis. Moreover, we estimated the correlation of MIRS with three
289 important immune checkpoint molecules: PD-1, PD-L1 and CTLA4. As illustrated in
290 Figure 4E, compared with MIRS^{high} group, MIRS^{low} group shows significantly higher
291 expression (Wilcoxon test $P < 0.0001$). MIRS scores are moderately correlated to the
292 expression levels of PD-1, PD-L1 and CTLA4 (Figure 4D). Overall, the differences in
293 tumor immunogenicity between the MIRS groups are significant, MIRS^{high} group has
294 relatively low immune infiltration level while MIRS^{low} group has relatively high
295 immune infiltration level. Similar results were also observed in TCGA and GSE96058
296 cohort (Supplementary Figure S8). This finding further suggested MIRS^{low} group
297 might have better response in therapy of immune checkpoint blockade.

298 To investigate the correlation between MIRS score and metastatic mechanism, we

299 firstly downloaded the metastasis breast cancer (METABRIC) cohort from human
300 cancer metastasis database <https://hcmdb.i-sanger.com/>, which contains primary tumor
301 and metastatic tumor. Then the functional analysis achieved by GSEA detects 23
302 qualified metastasis-related gene sets ($NES| > 1$, NOM p-value < 0.05 and FDR q-value
303 < 0.25). After that, ssGSEA analysis was used to evaluate the above significant
304 metastatic pathways. We observe that the metastatic pathways exhibit significant
305 difference between two MIRS groups, and the majority of MIRS^{high} group had higher
306 ssGSEA score (Figure 4G). A higher ssGSEA score suggests high activity of metastatic
307 processes. Similar results are found in TCGA and GSE96058 cohort (Figure S9-S10).
308 Furthermore, the expression discrepancy of three well-known genes (DCC, MMP9 and
309 ETS) were found to be correlated to the invasion and metastasis in breast cancer (29)
310 (Figure 4F), and MIRS exhibits moderately negative correlation with the expression of
311 these genes (Figure 4H).

312 We also examined the relationship between intrinsic molecular subtypes and MIRS.
313 In breast cancer, major subtypes based on the ER, PR and HER2 exist on tumor cells.
314 As shown in Figure S15, although the expression levels of ER, PR and HER2 were
315 moderately correlated with MIRS, the differences in the expression levels of ER, PR
316 and HER2 between high and low MIRS subtypes were statistically significant in TCGA
317 and GSE86166. Additionally, for TCGA cohort, we noticed the imbalanced proportions
318 of intrinsic molecular subtypes between MIRS^{high} and MIRS^{low} groups (Figure 4I).
319 48.09% of LumA tumor and 22.5% of Normal-like tumor are present in MIRS^{high} group
320 whereas 32.33% of LumB tumor in MIRS^{low} subtype. However, we found that higher

321 proportion of Basal-like tumor was present in MIRS^{low} subtype. In Muenst et al.' study
322 (30), they pointed out that the number of tumor-infiltrating lymphocytes was the highest
323 in the basal-like subtype, which may support a high enrichment of basal-like tumor in
324 MIRS^{low} group. We also found that the normal-like had significantly lower MIRS than
325 other molecular subtypes, in contrast to the LumB subtype had the highest MIRS
326 (Figure 4J). In addition, a statistically significant difference was detected among these
327 five intrinsic molecular subtypes by using Kruskal-Wallis method (Figure 4J).
328 Similar results were found in METABRIC cohort (Supplementary Figure S11). These
329 analyses indicate that MIRS group exhibited chaotic correlation with classic molecular
330 subtypes, which could be attributed to the high tumorous heterogeneity in breast cancer.

331 **Identification of MIRS related biological characteristics in** 332 **prognosis of breast cancer**

333 The above analyses implied high correlations between MIRS and tumor-infiltration
334 microenvironment as well as tumor metastasis. We further explore the molecular
335 mechanism of 12-gene panel underlying the prognosis of breast cancer. Initially,
336 through literature, we found that the majority of those prognosis-related genes, except
337 for *APOA5*, has been reported to be involved in the processes of tumorigenesis (Table
338 S6). It is worth mentioning that *APOA5*, encoding an apolipoprotein, is associated with
339 cardiovascular diseases (31, 32), but little work studies its roles in tumorigenesis and
340 prognosis. To delineate its potential prognostic role in breast cancer, we divided
341 *APOA5* expression into four quartiles, then GSEA analysis between the highest and

342 lowers quartiles in TCGA-BRCA was conducted. Interestingly, many metastatic and
343 immune-related pathways were observed to be enriched in the highest quartile,
344 including EMT, TNF α signaling and Immune response regulating signaling pathways
345 (Figure 5A and S12A). Subsequent survival analysis of pan-cancer based on TCGA
346 cohorts was performed via Kaplan-Meier Plotter (<https://kmplot.com/analysis/>) (33),
347 indicating that APOA5 may serve as prognostic indicator in many cancers (Figure
348 S12B). The breast cancer patients with the highest APOA5 expression have a worse
349 survival outcome (Figure S12B). Overall, our analysis hinted that APOA5 may exert
350 its prognostic function to affect the immune activity in breast cancer, and it is likely to
351 be a potential target for the future research of breast cancer therapy.

352 Genomic mutations are mostly involved in the survival prognosis of various cancers
353 (34). Thus, we tested the associations between somatic mutations and MIRS in TCGA
354 BRCA data. According to the analysis in the study of Chen et al (35), only the genes
355 with somatic mutation frequencies more than 2.5% were included. By analyzing the
356 mutation annotation of TCGA BRCA cohort, we selected the top 10 genes by mutation
357 frequency. As provided in Figure 5B and C, MIRS^{low} group has increased frequency of
358 mutation events than MIRS^{high} group. Rizvi et al (36) and Capalbo et al' studies (37)
359 demonstrated that the patients with more mutations might have an increased number of
360 neoantigens that enhance response to immunotherapy. This result might explain, in the
361 present study, the reason that MIRS^{low} group has better prognostic outcomes than
362 MIRS^{high} group.

363 Recently, tumor mutation burden (TMB) is the paramount prognostic measure in

364 cancer survival (38). We further investigated the associations between MIRS and TMB.
365 As illustrated in Figure 5D, the patients in MIRS^{low} group exhibited markedly increased
366 TMB when compared with those with MIRS^{high} group. Lee et al (39) and Karn et al's
367 studies (40) showed that high TMB was associated with improved survival.
368 Additionally, Chen et al (35) reported that the increased TMB was correlated to
369 improved response to PD-1 blockades therapy. Correlation analysis between MIRS and
370 TMB demonstrated that MIRS score was negligibly correlated with TMB (Spearman
371 coefficient: $R = -0.1$, $p = 0.0011$; Figure 5E). These findings indicate that MIRS may
372 be related to immunotherapy response, and the patients with lower MIRS may have
373 probably response in immunotherapy.

374 **The role of MIRS in the prediction of therapeutic** 375 **benefits**

376 To explore predictive ability of MIRS in immunotherapy for each patient, T cell
377 inflamed score (TIS), IFN- γ signature, antigen presenting machinery genes
378 (APM) and Immunotherapyscore (IPS) (20, 41, 42), which are prevailing predictors of
379 clinical response to immunotherapy across different tumor types were compared.
380 Notably, the higher of TIS, IFN- γ score, APM and IPS mean that patients
381 receiving immunotherapy are more likely to response All patients in GSE20711 and
382 GSE58812 with MIRS^{low} showed significantly increased predictor scores than those
383 with MIRS^{high} (Figure 6A and S13A), which hints that MIRS^{low} group is more likely to
384 have immunotherapy response. To further appraise the prognostic capability of

385 MIRS^{low} group in immunotherapy, the differences in overall survival between MIRS^{high}
386 and MIRS^{low} groups were compared using KM survival analysis in breast cancer testing
387 cohort. Unfortunately, there are hitherto few public datasets of breast cancer patients
388 receiving immunotherapy. Instead, the data of melanoma from Liu et al (43) and
389 TCGA-SKCM dataset with patient receiving immunotherapy were used in present
390 analysis. As a result, compared with PD-1 and TMB biomarkers upon receiving anti-
391 PD-1 treatment, MIRS showed robust AUCs (Figure 6B-D). Furthermore, the patients
392 with MIRS^{high} have significantly shorter overall survival than their counterparts (Figure
393 6E and Figure S13B). MIRS significantly increases in patients with stable disease (SD)
394 or progressive disease (PD) when compared with those with complete response (CR)
395 or partial response (PR) (Figure 6F and Figure S13CB). Besides, the distributions of
396 CR/PR and SD/PD across MIRS^{high} and MIRS^{low} groups were also validated. We found
397 that patients in MIRS^{low} group had better response to immunotherapy than those in
398 MIRS^{high} group (Figure 6G and Figure S13DC).

399 Moreover, to assess therapeutic value of MIRS in chemotherapy, we examined its
400 predictive potential in GSE20685 with the breast cancer patients who receive adjuvant
401 chemotherapy. The optimal cutoffs of MIRS were determined by the median cutoff,
402 then the patients were stratified into MIRS^{high} and MIRS^{low} group. Survival analysis
403 displays that the breast cancer patients with MIRS^{low} had much better survival than
404 those with MIRS^{high} in adjuvant chemotherapy cases (Figure 6H). We also investigated
405 the prognosis of different MIRS subtypes with or without adjuvant chemotherapy. As
406 illustrated in Figure 6I, we found that MIRS^{high} group had statistically significant

407 differences between the patients who were treated with adjuvant chemotherapy and
408 those without adjuvant chemotherapy. However, a consistent result was not observed
409 in those patients with MIRS^{low} (Figure 6J). These results indicated that adjuvant
410 chemotherapy might be more beneficial to MIRS^{high} group. Based on the gene sets of
411 different drug treatments retrieved from MSigDB database, GSEA predicted that
412 MIRS^{high} was significantly correlated with drug sensitivity in TCGA cohort (Figure 6K).
413 Moreover, the R package pRRophetic was used to estimate the sensitivity of three
414 chemotherapeutic drugs, including cisplatin, tamoxifen and vincristine, which have
415 been commonly used in breast cancer treatment. The results showed that estimated IC50
416 values of cisplatin and vincristine significantly decrease in MIRS^{high} subtype (Figure
417 6L). We did not display IC50 boxplot of tamoxifen due to the R package ‘pRRophetic’
418 does not contain resistant information regarding tamoxifen.

419 These results suggest that MIRS holds massive potential for predicting the response
420 to chemotherapy and immunotherapy in breast cancer patients. In brief, the patients
421 with MIRS^{high} may benefit from the chemotherapy, and patients with MIRS^{low} are likely
422 to be more sensitive to the immunotherapy.

423

424 **Comparison of MIRS with the previously prognostic** 425 **models**

426 Before the creation of MIRS, Shimizu et al (26) demonstrated that 23-gene panel
427 (mPS) helps predict OS in breast cancer patients based on analogous neuron network

428 model; Cui's score (44) constructed 8-gene signature based on traditional Lasso Cox
429 model. We then comprehensively evaluate the prognostic power of our MIRS, mPS and
430 Cui's score by 0prognostic Cox analyses based on a variety of public datasets. Our
431 MIRS performed very well in different cohorts (Figure 7A). Although mPS showed to
432 be more robust than MIRS in many datasets, some of the HRs in mPS panel were not
433 significant (P value > 0.05) (Figure 7B). Cui's score performed the worst among these
434 models (Figure 7C).

435 Furthermore, we scrutinized the predictive potential of these three models in the
436 response to immunotherapy. The malignant melanoma cohort data (43) that receives
437 anti-PD-1 therapy was used. The optimal cutoffs of Cui's score and mPS value were
438 determined by the median. KM survival curves of MIRS show a significant difference
439 in OS between MIRS^{high} and MIRS^{low} group (Figure 7E). On the contrary, the survival
440 analysis of mPS and Cui's score revealed that patients with low mPS or Cui's score
441 showed no statistically significant difference when compared with those with high
442 mPS or Cui's score (Figure 7F-G). MIRS, mPS and Cui's score were also examined
443 with time-dependent ROC analysis in the testing cohort for prediction in
444 immunotherapeutic benefits. Notably, our MIRS exhibited much better predictive
445 ability than mPS and Cui's score for OS at 1 year, 1.5 years, and 2 years, respectively
446 (Figure 7D).

447 Discussion

448 With the development of transformative technologies, analyses of high throughput

449 sequencing data have significantly deepened the understanding of modern biology,
450 enabling the scientists to thoroughly explore key characteristics in a variety of cancers.
451 Metastasis and tumor-immune infiltration are two of the major characteristics, and have
452 been extensively proven to be associated with tumorigenesis, drug resistance and
453 prognosis in breast cancer (43). Quite a few studies have disclosed the roles of
454 metastasis and tumor-immune infiltration as prognostic factors in predicting the
455 survival outcomes for breast cancer (45). Unfortunately, breast tumors are highly
456 heterogeneous among individuals, and much current work has only considered organ-
457 specific metastasis or immune infiltration level and thus insufficient to achieve robust
458 predictive power on prognosis. To address this issue, in this study we developed a
459 comprehensive and efficient prognosis model, considering metastasis and immune
460 infiltration levels together, to aid clinicians in providing precise treatment strategies.

461 Given the promising predictive value of MIRS, we systematically investigated the
462 relationships between MIRS and clinical pathological characteristics. In different
463 sequencing platform data, MIRS demonstrated as an independent prognosis factor
464 compared with other conventional clinical features (Figure S14). As illustrated in
465 Figure S17A, we observed differences between MIRS and Age, Gender and Metastasis
466 variables. Subsequently, we used decision curve analysis (DCA) to decipher the effect
467 in combining MIRS with clinical indicators. In the DCA analysis, the net benefit of
468 clinical indicators combined with MIRS were better than of sole clinical indicator
469 (Figure S17B). Additionally, we employed TCGA and GSE96058 datasets to
470 investigate whether MIRS is suitable for all BRCA subtypes due to its complete subtype

471 information. However, we have not observed consistently predictive ability in both
472 datasets (Figure S18). This unsatisfied performance may come from the fact that we
473 built our MIRS model without tumor subtype information. Together, these results
474 demonstrate the validity and reliability of MIRS in clinical applications, but it no
475 suitable to all subtypes in breast cancer.

476 Next, we compared MIRS with the representative prognostic models, mPS and
477 Cui's score. Univariate cox regression analysis using nine public cohorts indicated that
478 MIRS and mPS performed well in most cohorts. These results indicated that,
479 constructing a prognostic system considering only metastatic features may be
480 insufficient. Compared with AI methods, traditional survival model showed weak
481 power. Nonetheless, mPS scoring system, based on an analogous AI approach, does not
482 work well in predicting immunotherapeutic. It might explain that the establishment of
483 mPS does not consider immunogenomic features, thus failing to achieve satisfactory
484 immunotherapeutic prediction.

485 Apart from being informative regarding prognosis, MIRS can also act as an
486 independent predictor to guide therapeutic strategies. Our analyses indicated that
487 MIRS^{high} group had lower TIS, IPS, IFN-gamma score and APM score, implying
488 MIRS^{high} group is more likely to escape from immunity in breast cancer. For further
489 validation, we tested if the OS between MIRS^{high} and MIRS^{low} groups was associated
490 with immunotherapy. We used two malignant melanoma cohorts with
491 immunotherapeutic information by conducting KM analysis. This survival analysis
492 showed that MIRS^{low} group exhibited improved survival and better response to

493 immunotherapy than MIRS^{high}. We speculate that the immunotherapy may achieve
494 beneficial treatment for MIRS^{low} patients.

495 Currently, chemotherapy is one of the main treatments for breast cancer. Hence, it
496 is necessary to identify patients who may potentially benefit from chemotherapy.
497 Through the analysis of breast cancer patients with chemotherapy clinical
498 information, we found that patients with MIRS^{high} respond better to chemotherapy than
499 patients with MIRS^{low}. Chemotherapy has been reported to be related to immune
500 infiltration (46). In the Ahn et al's study (47), they demonstrated that the high level of
501 the CD8+ TILs filtration is associated with chemotherapy resistance. This may be the
502 reason that high filtration MIRS^{low} subtype shows favorable chemotherapy. These
503 results emphasize the significance of MIRS^{high} patients who could benefit from
504 chemotherapy.

505 As a gene prognostic signature particularly designed for breast cancer patients,
506 MIRS is a novel and robust approach in risk stratification and personalized treatment.
507 However, there are still flaws in the current study. First, due to the remarkable intra-
508 tumor heterogeneity in breast cancer, we cannot cover all metastatic signatures despite
509 a large numbers of breast cancer patients used in this study. Second, only the median
510 cutoff of MIRS is used to classify the patients into high and low subtypes, the optimal
511 cutoff of MIRS would be needed to provide rational strategies. Lastly, all the
512 conclusions in this research are obtained from *in silico* studies, clinical experiments are
513 required to confirm our findings.

514 MIRS has the potential to assist oncologists to screen patients who are more likely

515 to benefit from immunotherapy or chemotherapy. It would be of great significance to
516 validate the value of MIRS in prospective clinical trials.

517 **Contributors**

518 CH and XDZ conceived the presented idea. CH, DLL and MD collected the public data,
519 MD and CH developed the methodology, MD and DLL analyzed the data under the
520 supervision of CH. CH and MD took the lead in drafting the manuscript with input from
521 all authors. CH, ELHL and XDZ revised the manuscript, PYZ and BQS interpreted
522 results from a clinical point of view. All authors read and approved the final manuscript.

523 **Declaration of competing interests**

524 The authors have declared that no competing interest exists

525 **Acknowledgements**

526 This work was supported by Dr. Neher's Biophysics Laboratory for Innovative Drug Discovery
527 (File no. 001/2020/ALC), by the Science and Technology Development Fund, Macau
528 Government (File no. 0020/2021/A), by the University of Macau (grant numbers: FHS-CRDA-
529 029-002-2017 and MYRG2018-00071-FHS), Zhongnanshan Medical Foundation of
530 Guangdong Province (grant number: ZNSA-2021016) and the Science and Technology
531 Development Fund, Macau SAR (File no. 0004/2019/AFJ and 0011/2019/AKP).

532 **Data sharing statement**

533 Data are available in a public, open access repository. All used data in the current study
534 are downloaded from Gene Expression Omnibus (GEO, <https://www.ncbi.nlm.nih.gov/geo/>) and The Cancer Genome Atlas (TCGA) database (<https://cancergenome.nih.gov/>).

537 **References**

- 538 1. Dumas A, Vaz Luis I, Bovagnet T, El Mouhebb M, Di Meglio A, Pinto S, et al.
539 Impact of Breast Cancer Treatment on Employment: Results of a Multicenter
540 Prospective Cohort Study (CANTO). *J Clin Oncol.* 2020;38(7):734-43.
- 541 2. Siegel RL, Miller KD, Jemal A. Cancer statistics, 2018. *CA Cancer J Clin.*
542 2018;68(1):7-30.
- 543 3. Afifi AM, Saad AM, Al-Husseini MJ, Elmehrath AO, Northfelt DW, Sonbol MB.
544 Causes of death after breast cancer diagnosis: A US population-based analysis.
545 *Cancer.* 2020;126(7):1559-67.
- 546 4. Gansler T, Ganz PA, Grant M, Greene FL, Johnstone P, Mahoney M, et al. Sixty
547 years of CA: a cancer journal for clinicians. *CA Cancer J Clin.* 2010;60(6):345-50.
- 548 5. Landemaine T, Jackson A, Bellahcene A, Rucci N, Sin S, Abad BM, et al. A six-gene
549 signature predicting breast cancer lung metastasis. *Cancer research.*
550 2008;68(15):6092-9.
- 551 6. Scott E. Androgen deprivation with or without radiation therapy for clinically
552 node-positive prostate cancer. Lin CC, Gray PJ, Jemal A, Efstathiou JA, Surveillance
553 and Health Services Research Program, Intramural Research, American Cancer
554 Society, Atlanta, GA (CCL, AJ); Department of Radiation Oncology, Massachusetts
555 General Hospital, Harvard Medical School, Boston, MA (PJG, JAE); e-mail:
556 jefstathiou@partners.org. *J Natl Cancer Inst.* 2015 May 9;107(7). pii: djv119. [Print
557 2015 Jul]. doi: 10.1093/jnci/djv119. *Urol Oncol.* 2017;35(3):122-3.
- 558 7. Early Breast Cancer Trialists' Collaborative G. Effects of chemotherapy and
559 hormonal therapy for early breast cancer on recurrence and 15-year survival: an
560 overview of the randomised trials. *Lancet.* 2005;365(9472):1687-717.
- 561 8. Greenberg PA, Hortobagyi GN, Smith TL, Ziegler LD, Frye DK, Buzdar AU. Long-
562 term follow-up of patients with complete remission following combination
563 chemotherapy for metastatic breast cancer. *J Clin Oncol.* 1996;14(8):2197-205.
- 564 9. Cremasco V, Astarita JL, Grauel AL, Keerthivasan S, Maclsaac K, Woodruff MC, et
565 al. FAP Delineates Heterogeneous and Functionally Divergent Stromal Cells in
566 Immune-Excluded Breast Tumors. *Cancer Immunol Res.* 2018;6(12):1472-85.

- 567 10. Landemaine T, Jackson A, Bellahcene A, Rucci N, Sin S, Abad BM, et al. A six-gene
568 signature predicting breast cancer lung metastasis. *Cancer Res.* 2008;68(15):6092-9.
- 569 11. Li J, Lenferink AE, Deng Y, Collins C, Cui Q, Purisima EO, et al. Identification of
570 high-quality cancer prognostic markers and metastasis network modules. *Nat*
571 *Commun.* 2010;1:34.
- 572 12. Xiao Z, Hu L, Yang L, Wang S, Gao Y, Zhu Q, et al. TGFbeta2 is a prognostic-
573 related biomarker and correlated with immune infiltrates in gastric cancer. *J Cell Mol*
574 *Med.* 2020;24(13):7151-62.
- 575 13. Shen Y, Peng X, Shen C. Identification and validation of immune-related lncRNA
576 prognostic signature for breast cancer. *Genomics.* 2020;112(3):2640-6.
- 577 14. Erdag G, Schaefer JT, Smolkin ME, Deacon DH, Shea SM, Dengel LT, et al.
578 Immunotype and immunohistologic characteristics of tumor-infiltrating immune cells
579 are associated with clinical outcome in metastatic melanoma. *Cancer Res.*
580 2012;72(5):1070-80.
- 581 15. Yang L, Wang S, Zhang Q, Pan Y, Lv Y, Chen X, et al. Clinical significance of the
582 immune microenvironment in ovarian cancer patients. *Mol Omics.* 2018;14(5):341-
583 51.
- 584 16. Barnes TA, Amir E. HYPE or HOPE: the prognostic value of infiltrating immune
585 cells in cancer. *British journal of cancer.* 2017;117(4):451-60.
- 586 17. Pérez-Romero K, Rodríguez RM, Amedei A, Barceló-Coblijn G, Lopez DH. Immune
587 Landscape in Tumor Microenvironment: Implications for Biomarker Development
588 and Immunotherapy. *International Journal of Molecular Sciences.* 2020;21(15):5521.
- 589 18. Liu R, Hu R, Zeng Y, Zhang W, Zhou H-H. Tumour immune cell infiltration and
590 survival after platinum-based chemotherapy in high-grade serous ovarian cancer
591 subtypes: A gene expression-based computational study. *EBioMedicine.*
592 2020;51:102602.
- 593 19. Davis S, Meltzer PS. GEOquery: a bridge between the Gene Expression Omnibus
594 (GEO) and BioConductor. *Bioinformatics.* 2007;23(14):1846-7.
- 595 20. Charoentong P, Finotello F, Angelova M, Mayer C, Efremova M, Rieder D, et al.
596 Pan-cancer immunogenomic analyses reveal genotype-immunophenotype
597 relationships and predictors of response to checkpoint blockade. *Cell reports.*
598 2017;18(1):248-62.
- 599 21. Charoentong P, Finotello F, Angelova M, Mayer C, Efremova M, Rieder D, et al.
600 Pan-cancer Immunogenomic Analyses Reveal Genotype-Immune Phenotype
601 Relationships and Predictors of Response to Checkpoint Blockade. *Cell Rep.*
602 2017;18(1):248-62.
- 603 22. Bhattacharya S, Dunn P, Thomas CG, Smith B, Schaefer H, Chen J, et al. ImmPort,
604 toward repurposing of open access immunological assay data for translational and
605 clinical research. *Sci Data.* 2018;5:180015.
- 606 23. Hänzelmann S, Castelo R, Guinney J. GSVA: gene set variation analysis for
607 microarray and RNA-seq data. *BMC bioinformatics.* 2013;14(1):7.
- 608 24. Gentleman R, Carey VJ. Unsupervised machine learning. *Bioconductor case*
609 *studies: Springer;* 2008. p. 137-57.
- 610 25. Benjamini Y, Hochberg Y. Controlling the false discovery rate: a practical and

- 611 powerful approach to multiple testing. *Journal of the Royal statistical society: series*
612 *B (Methodological)*. 1995;57(1):289-300.
- 613 26. Shimizu H, Nakayama KI. A 23 gene–based molecular prognostic score precisely
614 predicts overall survival of breast cancer patients. *EBioMedicine*. 2019;46:150-9.
- 615 27. Agarap AF. Deep learning using rectified linear units (relu). *arXiv preprint*
616 *arXiv:180308375*. 2018.
- 617 28. Yoshihara K, Shahmoradgoli M, Martínez E, Vegesna R, Kim H, Torres-Garcia W,
618 et al. Inferring tumour purity and stromal and immune cell admixture from
619 expression data. *Nature communications*. 2013;4(1):1-11.
- 620 29. Wakita K, Kohno N, Sakoda Y, Ishikawa Y, Sakaue M. Decreased expression of the
621 DCC gene in human breast carcinoma. *Surgery today*. 1996;26(11):900-3.
- 622 30. Muenst S, Soysal S, Gao F, Obermann E, Oertli D, Gillanders W. The presence of
623 programmed death 1 (PD-1)-positive tumor-infiltrating lymphocytes is associated
624 with poor prognosis in human breast cancer. *Breast cancer research and treatment*.
625 2013;139(3):667-76.
- 626 31. Dallongeville J, Cottel D, Montaye M, Codron V, Amouyel P, Helbecque N. Impact
627 of APOA5/A4/C3 genetic polymorphisms on lipid variables and cardiovascular
628 disease risk in French men. *International journal of cardiology*. 2006;106(2):152-6.
- 629 32. Lin Y-C, Nunez V, Johns R, Shiao SPK. APOA5 gene polymorphisms and
630 cardiovascular diseases: metaprediction in global populations. *Nursing research*.
631 2017;66(2):164-74.
- 632 33. Nagy Á, Munkácsy G, Gyórfy B. Pancancer survival analysis of cancer hallmark
633 genes. *Scientific reports*. 2021;11(1):1-10.
- 634 34. Gotea V, Gartner JJ, Qutob N, Elnitski L, Samuels Y. The functional relevance of
635 somatic synonymous mutations in melanoma and other cancers. *Pigment cell &*
636 *melanoma research*. 2015;28(6):673-84.
- 637 35. Chen Z, Yuan Y, Chen X, Chen J, Lin S, Li X, et al. Systematic comparison of
638 somatic variant calling performance among different sequencing depth and mutation
639 frequency. *Scientific reports*. 2020;10(1):1-9.
- 640 36. Rizvi NA, Hellmann MD, Snyder A, Kvistborg P, Makarov V, Havel JJ, et al.
641 Mutational landscape determines sensitivity to PD-1 blockade in non–small cell lung
642 cancer. *Science*. 2015;348(6230):124-8.
- 643 37. Capalbo C, Scafetta G, Filetti M, Marchetti P, Bartolazzi A. Predictive biomarkers
644 for checkpoint inhibitor-based immunotherapy: the galectin-3 signature in NSCLCs.
645 *International journal of molecular sciences*. 2019;20(7):1607.
- 646 38. Riviere P, Goodman AM, Okamura R, Barkauskas DA, Whitchurch TJ, Lee S, et al.
647 High tumor mutational burden correlates with longer survival in immunotherapy-
648 naïve patients with diverse cancers. *Molecular Cancer Therapeutics*.
649 2020;19(10):2139-45.
- 650 39. Lee D-W, Han S-W, Bae JM, Jang H, Han H, Kim H, et al. Tumor mutation burden
651 and prognosis in patients with colorectal cancer treated with adjuvant
652 fluoropyrimidine and oxaliplatin. *Clinical Cancer Research*. 2019;25(20):6141-7.
- 653 40. Karn T, Denkert C, Weber K, Holtrich U, Hanusch C, Sinn B, et al. Tumor
654 mutational burden and immune infiltration as independent predictors of response to

- 655 neoadjuvant immune checkpoint inhibition in early TNBC in GeparNuevo. *Annals of*
656 *Oncology*. 2020;31(9):1216-22.
- 657 41. Ayers M, Lunceford J, Nebozhyn M, Murphy E, Loboda A, Kaufman DR, et al. IFN-
658 γ -related mRNA profile predicts clinical response to PD-1 blockade. *The Journal of*
659 *clinical investigation*. 2017;127(8):2930-40.
- 660 42. Kamoun A, de Reyniès A, Allory Y, Sjö Dahl G, Robertson AG, Seiler R, et al. A
661 consensus molecular classification of muscle-invasive bladder cancer. *European*
662 *urology*. 2020;77(4):420-33.
- 663 43. Liu D, Schilling B, Liu D, Sucker A, Livingstone E, Jerby-Amon L, et al. Integrative
664 molecular and clinical modeling of clinical outcomes to PD1 blockade in patients with
665 metastatic melanoma. *Nature medicine*. 2019;25(12):1916-27.
- 666 44. Cui Q, Tang J, Zhang D, Kong D, Liao X, Ren J, et al. A prognostic eight - gene
667 expression signature for patients with breast cancer receiving adjuvant
668 chemotherapy. *Journal of cellular biochemistry*. 2020;121(8-9):3923-34.
- 669 45. Wang S, Zhang Q, Yu C, Cao Y, Zuo Y, Yang L. Immune cell infiltration-based
670 signature for prognosis and immunogenomic analysis in breast cancer. *Briefings in*
671 *bioinformatics*. 2020.
- 672 46. Mesnage S, Auguste A, Genestie C, Dunant A, Pain E, Drusch F, et al.
673 Neoadjuvant chemotherapy (NACT) increases immune infiltration and programmed
674 death-ligand 1 (PD-L1) expression in epithelial ovarian cancer (EOC). *Annals of*
675 *Oncology*. 2017;28(3):651-7.
- 676 47. Ahn S, Chung YR, Seo AN, Kim M, Woo JW, Park SY. Changes and prognostic
677 values of tumor-infiltrating lymphocyte subsets after primary systemic therapy in
678 breast cancer. *PloS one*. 2020;15(5):e0233037.
- 679

680 Tables

681 **Table 1. The 12 prognostic genes for calculating the risk score in TCGA**
682 **data**

Gene ID	Category	Gene expression (high)	Gene expression (low)	Weight
APOA5	Dangerous	1	0	0.4703
FAM9C	Dangerous	1	0	0.5585
IVL	Dangerous	1	0	0.4467
PAGE5	Dangerous	1	0	0.5637
CACNA1E	Protective	0	1	0.3596
CCL25	Protective	0	1	0.5013
CD1A	Protective	0	1	0.1782
CD1B	Protective	0	1	0.7733
GPR55	Protective	0	1	0.6999

LAX1	Protective	0	1	0.6383
TNFRSF8	Protective	0	1	0.6234
WNT10A	Protective	0	1	0.4189

683

684

685 **Figure Legends**

686

687 **Figure 1. Systematic bioinformatic analysis pipeline.**

688

689 **Figure 2. Exploration of the immune cell infiltration grouping, and 52 candidate**
690 **genes were expressed in BRCA samples from the TCGA, GSE10893, and**
691 **GSE3521 datasets.**

692 (A) Heatmap for the high and low immune-cell infiltration grouping from the TCGA cohort.

693 (B) Boxplots for the expression levels of the CD family gene between high and low infiltration groups.

694 (C) Comparison of Stromal score, Immunity score, ESTIMATE score and Tumor purity between the
695 high and low immune infiltration groups.

696 (D) Boxplots illustrate the 22 immune cell proportions between high and low immune infiltration
697 groups.

698 (E) Venn plot of the differentially expressed genes from the TCGA data and GEO datasets.

699 (F) Heatmap of the 52 candidate gene expression values between high and low immune infiltration
700 groups from the TCGA dataset. Cluste1 represents the low immune infiltration level group, cluster 2
701 represents the high immune infiltration level group.

702 (G) Volcano plot of the 52 candidate genes between the primary and metastasis tumor groups both
703 from the GSE3521 and GSE10893. The blue dots show the DE genes are down regulated in the
704 metastasis group. The red dots display the DE genes are up regulated in the metastasis group. The p-
705 values were calculated using Wilcox rank sum test.

706

707 **Figure 3. Construction and validation of the MIRS in the training and testing**
708 **cohorts.**

709 **A.** The square root of the variance inflation factor value for each candidate gene in the training
710 data.

711 **B.** Correlations between the candidate genes in the training TCGA data. Different correlations
712 between two genes are represented by different colors.

713 **C.** ROC curve for the patient's overall survival prediction in the training TCGA data.

714 **D.** Kaplan-Meier curves of overall survival according to the MIRS subtypes in the training TCGA
715 data.

716 **E.** ROC curve for the patient's overall survival prediction in GSE96058.

717 **F.** Kaplan-Meier curves of overall survival according to the MIRS subtypes in GSE96058.

718 **G.** ROC curve for the patient's overall survival prediction in GSE86166.

719 **H.** Kaplan-Meier curves of overall survival according to the MIRS subtypes in GSE86166.

720

721 **Figure 4. Correlation of MIRS with the metastatic and immunogenomic**
722 **landscape between the high and low MIRS subtypes.**

- 723 **A.** Comparison of the Stromal score, ESTIMATE score, Immune score, and Tumor purity between
724 high and low MIRS subtypes in GSE86166. The p-values were calculated using Wilcoxon rank
725 sum test. * $p < 0.05$; ** $p < 0.01$; *** $p < 0.0001$.
- 726 **B.** Function enrichment bar plot for the genes in GSE86166 which were highly correlated
727 (Spearman correlation coefficient ≥ 0.04) with 12 prognostic genes in GSE86166.
- 728 **C.** Boxplots of the ssGSEA score for 17 immune-related biological functions and pathways
729 between two MIRS subtypes in the GSE86166. The p-values were calculated using Wilcoxon
730 rank sum test. * $p < 0.05$; ** $p < 0.01$; *** $p < 0.0001$.
- 731 **D.** The spearman correlation between the gene expression levels of PD-1, PD-L1 and CTLA4 and
732 MIRS score in the GSE86166 data, respectively.
- 733 **E.** The boxplots of PD-1, PD-L1 and CTLA4 for two MIRS subtypes in the GSE86166 data. The
734 p-values were calculated using Wilcoxon rank sum test. * $p < 0.05$; ** $p < 0.01$; *** $p < 0.0001$.
- 735 **F.** The boxplots of DCC, MMP9 and ETS1 for two MIRS subtypes in GSE86166 dataset. The p-
736 values were calculated using Wilcoxon rank sum test. * $p < 0.05$; ** $p < 0.01$; *** $p < 0.0001$.
- 737 **G.** Boxplots of the ssGSEA score for 23 metastatic biological functions and pathways between two
738 MIRS subtypes in the GSE86166. The p-values were calculated using Wilcoxon rank sum test.
739 * $p < 0.05$; ** $p < 0.01$; *** $p < 0.0001$.
- 740 **H.** The spearman correlation between the gene expression levels of DCC, MMP9 and ETS1 and
741 MIRS score in the GSE86166, respectively.
- 742 **I.** Sankey diagram for the MIRS values with different intrinsic molecular subtypes in TCGA
743 patients.
- 744 **J.** Violin plots for the distribution of MIRS values in different intrinsic molecular subtypes at
745 TCGA BRCA cohort. The p-values were calculated using Kruskal-Wallis test. * $p < 0.05$;
746 ** $p < 0.01$; *** $p < 0.0001$.

747
748 **Figure 5. Identification of MIRS-related biological characteristics in prognosis of**
749 **breast cancer.**

- 750 **A.** GSEA enrichment plots in TCGA.
- 751 **B.** The Oncoplot of top 10 genes with the highest mutation frequency in high MIRS group (TCGA
752 data).
- 753 **C.** The Oncoplot of top 10 genes with the highest mutation frequency in low MIRS group (TCGA
754 data).
- 755 **D.** Boxplots of the MIRS score between the high and low TMB subtypes in TCGA data. The p-
756 values were calculated using Wilcoxon rank sum test.
- 757 **E.** The spearman correlation between the MIRS score and TMB values in TCGA data.

758
759 **Figure 6. The therapeutic benefit of the MIRS value.**

- 760 **A.** The boxplot of TIS, IPS, APM score and IFN gamma score between the high and low MIRS
761 in GSE20711.
- 762 **B.** ROC curves between the expression level of PD-1, TMB and MIRS of anti-PD1
763 immunotherapy response prediction in Liu et al data.
- 764 **C.** Time-dependent ROC curves of MIRS for anti-PD1 immunotherapy response prediction in the

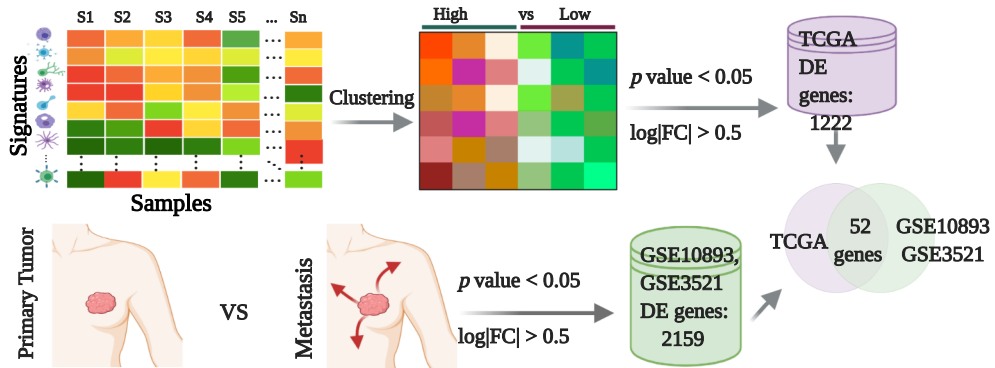
- 765 Liu et al data.
- 766 **D.** Time-dependent ROC curves of the expression level for anti-PD1 immunotherapy response
767 prediction in Liu et al data.
- 768 **E.** Kaplan-Meier curves of overall survival according to MIRS subtypes in the Liu et al data.
- 769 **F.** Violin plot illustrating the distribution of MIRS for patients with different immunotherapy
770 response in Liu et al data.
- 771 **G.** Bar graph showing the number of clinical responses to anti-PD-1 immunotherapy in the high
772 and low MIRS subtypes in Liu et al data.
- 773 **H.** Kaplan-Meier curves of overall survival according to MIRS subtypes with chemotherapy in
774 GSE20685.
- 775 **I.** Kaplan-Meier curves of overall survival according to the high MIRS subtype with or without
776 chemotherapy in GSE20685.
- 777 **J.** Kaplan-Meier curves of overall survival according to the low MIRS subtype with or without
778 chemotherapy in GSE20685.
- 779 **K.** GSEA predict that high MIRS group is negatively correlated with drug resistance in TCGA
780 cohort.
- 781 **L.** Chemotherapeutic sensitivity of two drugs (Cisplatin, Vincristine) were estimated and
782 compared in TCGA cohort.

783

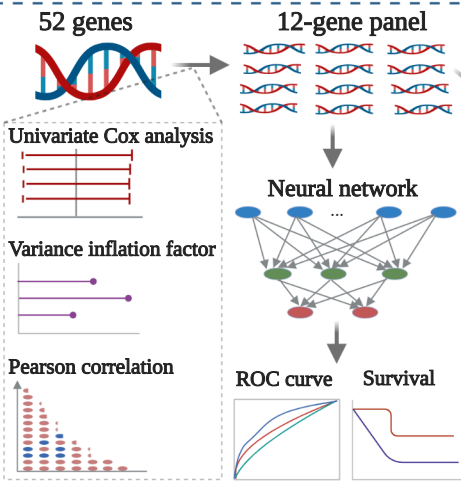
784 **Figure 7. Compare MIRS with previous prognosis signatures.**

- 785 **A.** A meta-analysis was performed using the prognosis results of MIRS in nine public datasets.
- 786 **B.** A meta-analysis was performed using the prognosis results of mPS in nine public datasets.
- 787 **C.** A meta-analysis was performed using the prognosis results of Cui's score in nine public
788 datasets.
- 789 **D.** Time-dependent ROC curves of anti-PD-1 immunotherapy on the 1-,1.5-,2-year survival rates
790 for Liu et al data.
- 791 **E.** Kaplan-Meier curves of overall survival according to MIRS subtype with immunotherapy in
792 Liu et al data.
- 793 **F.** Kaplan-Meier curves of overall survival according to mPS subtype with immunotherapy in Liu
794 et al data.
- 795 **G.** Kaplan-Meier curves of overall survival according to Cui's score subtype with immunotherapy
796 in Liu et al data.

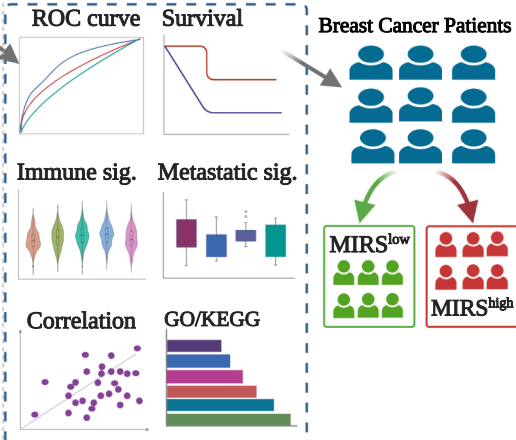
Step 1: Selection of candidates

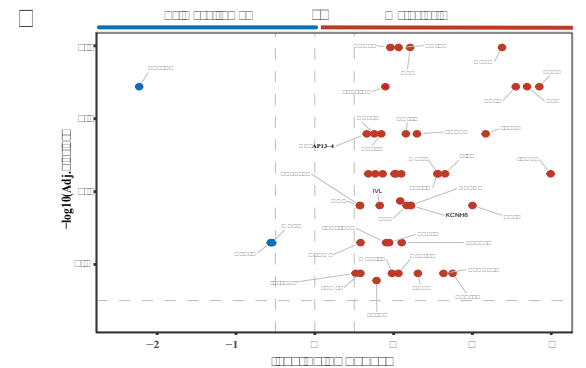
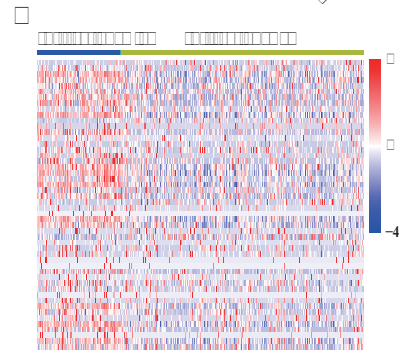
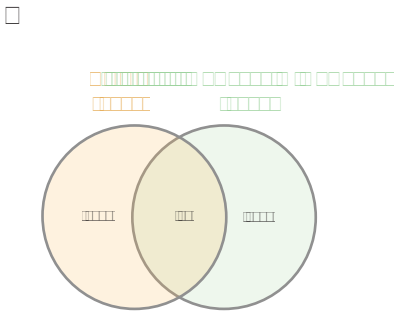
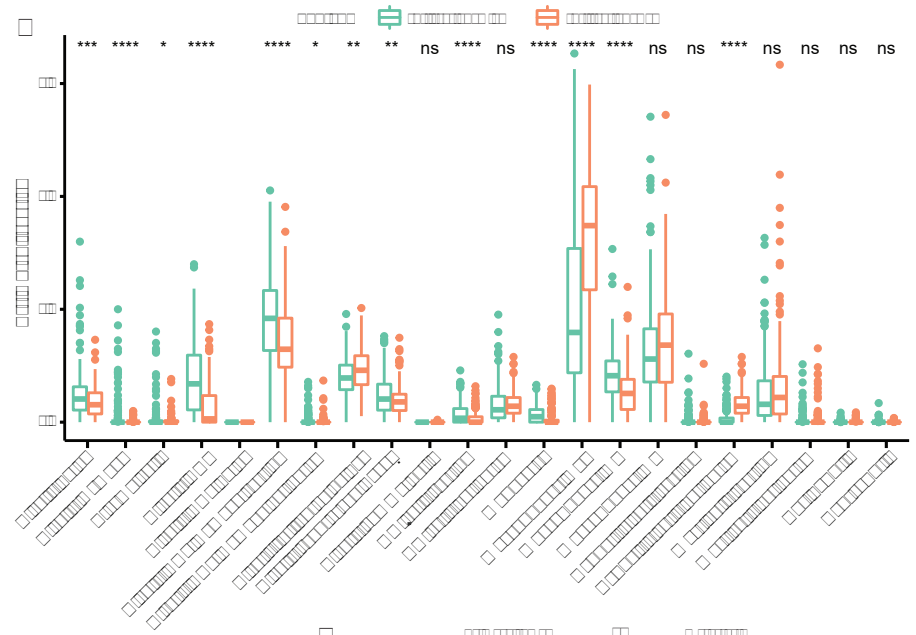
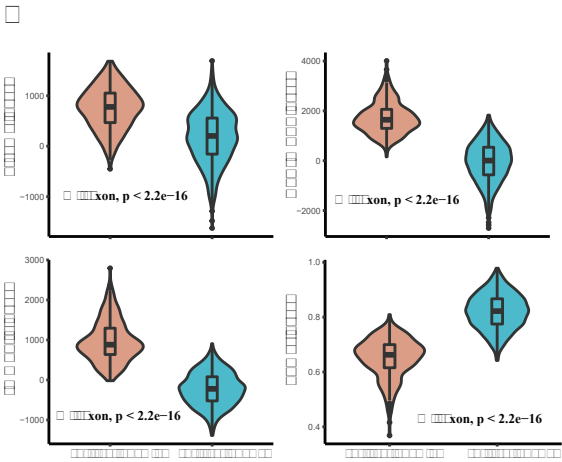
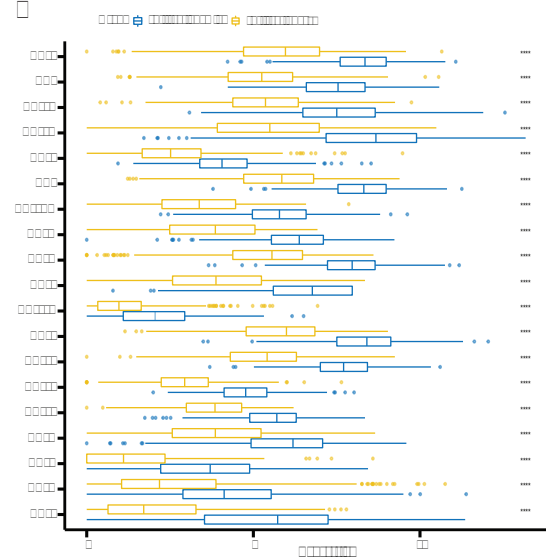
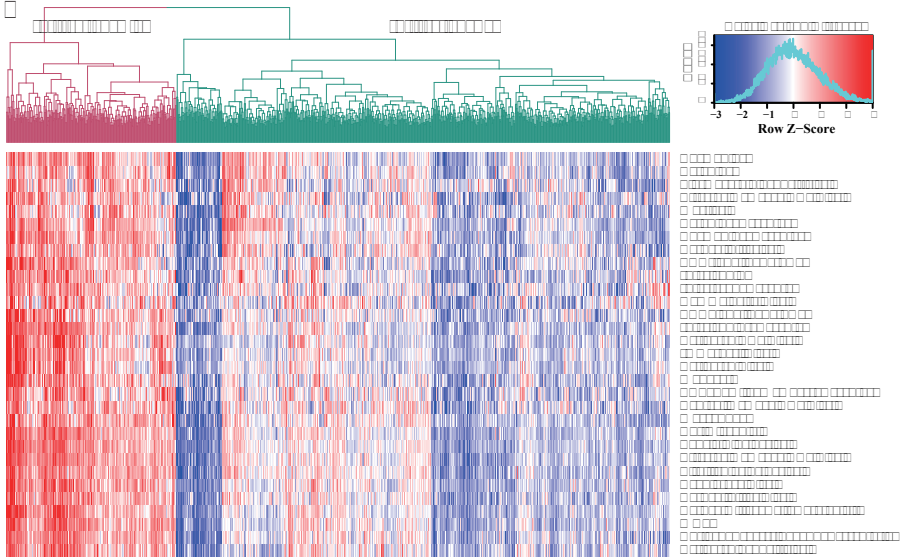


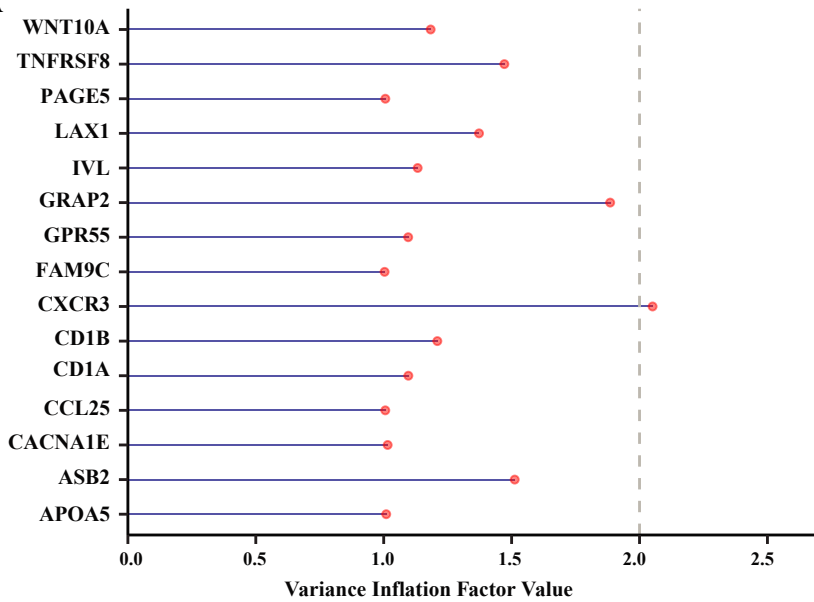
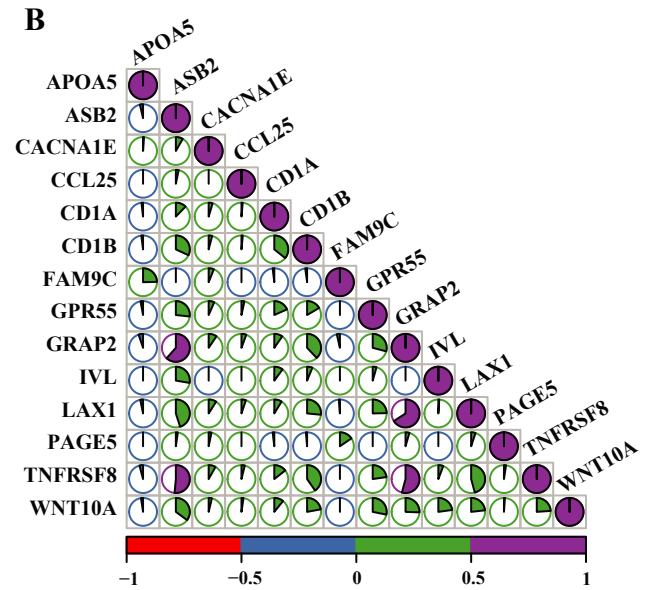
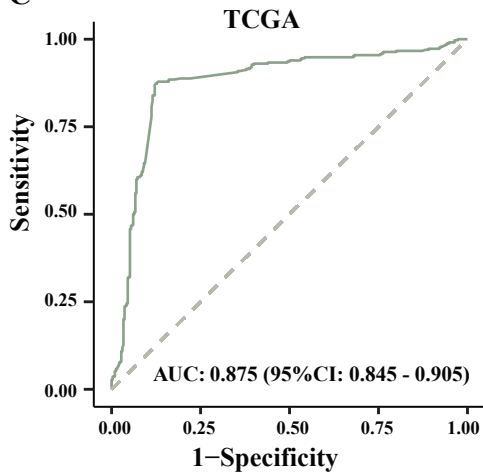
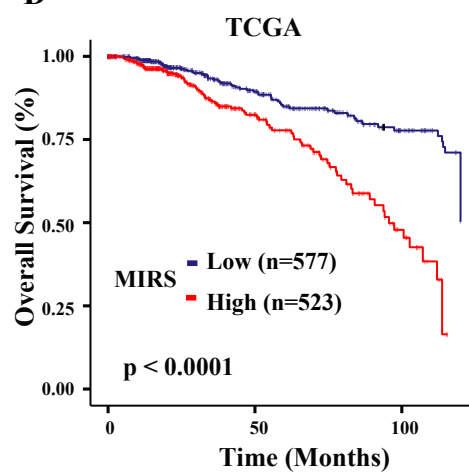
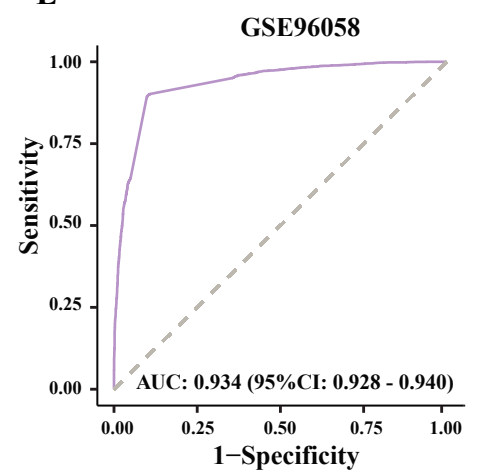
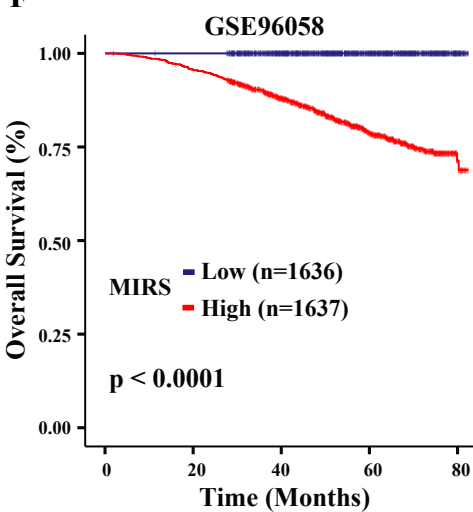
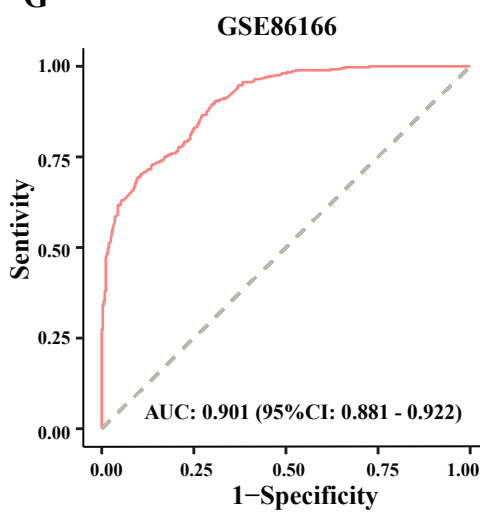
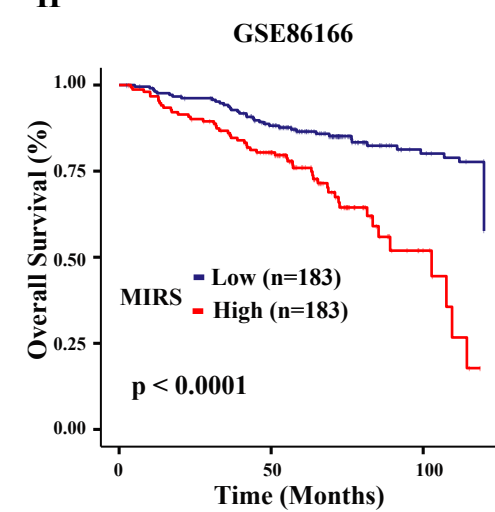
Step 2: Construction of MIRS

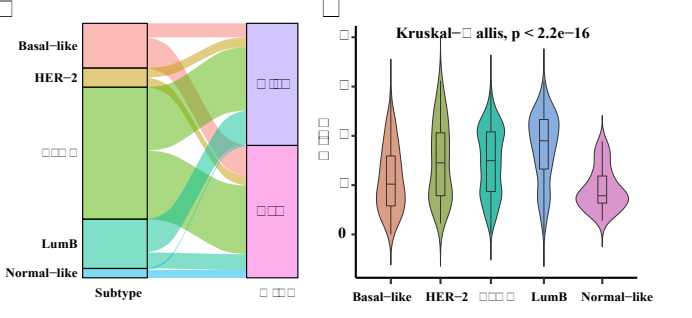
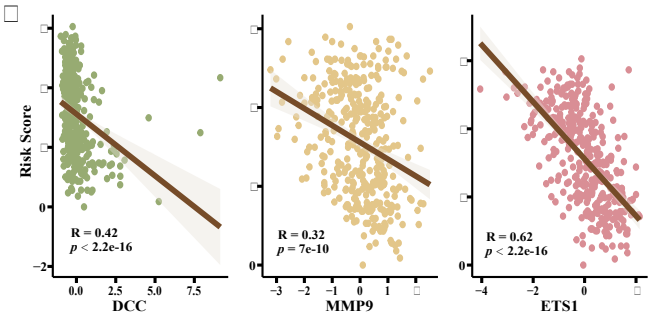
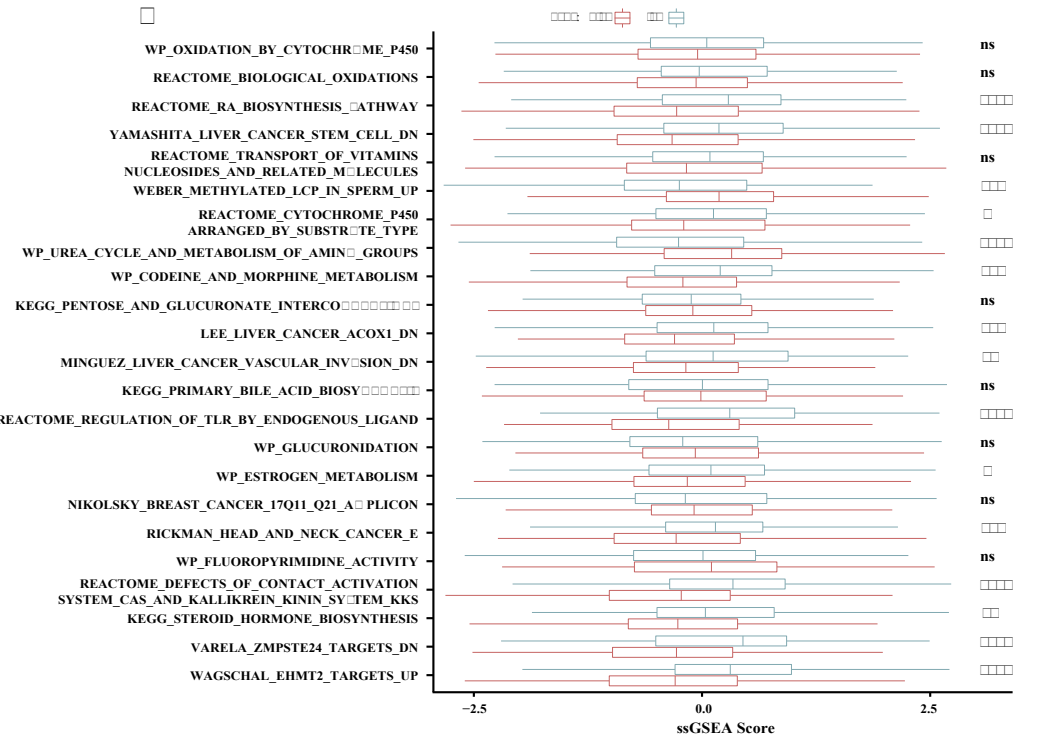
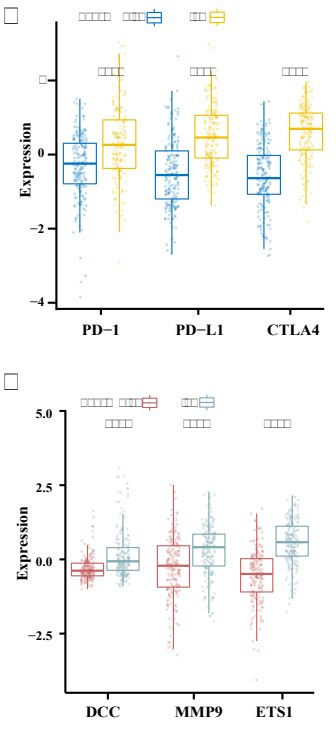
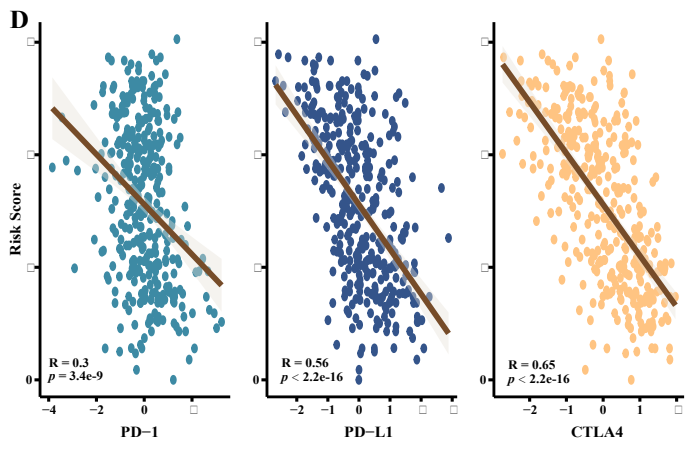
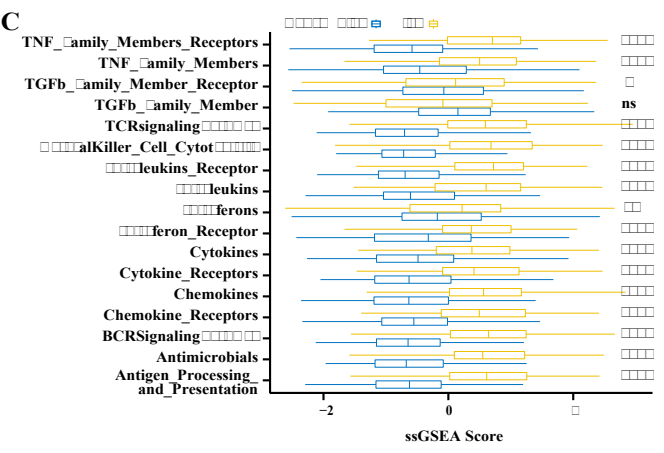
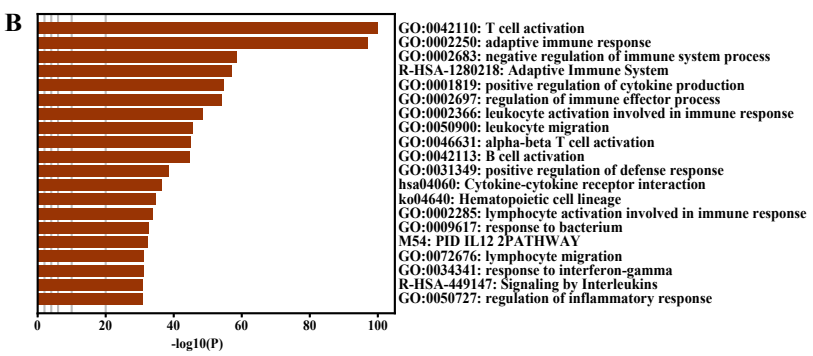
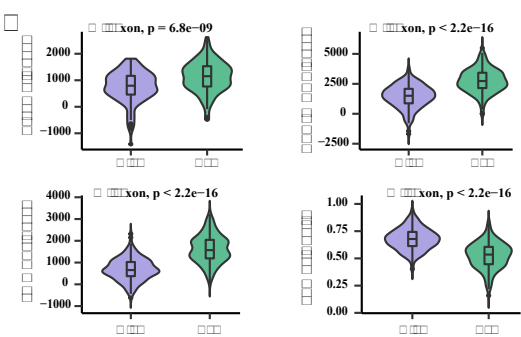


Step 3: Validation of MIRS

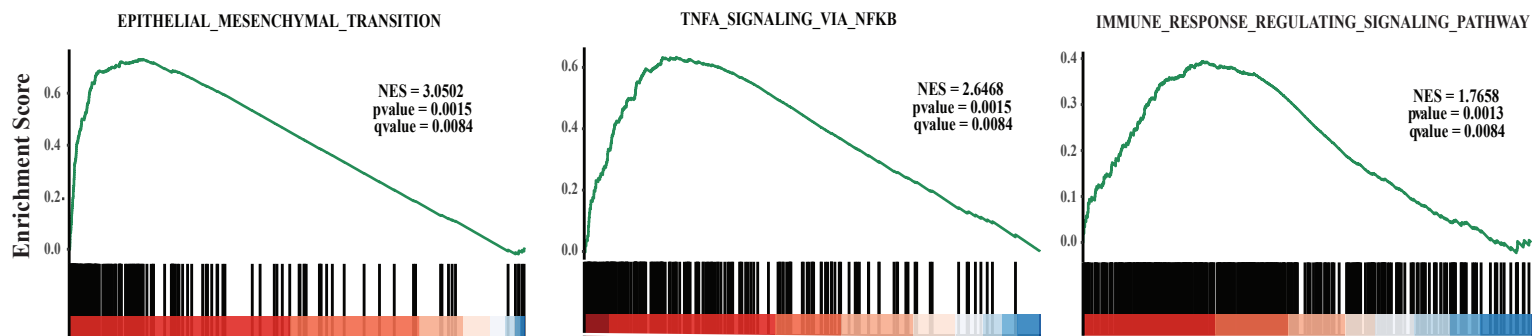




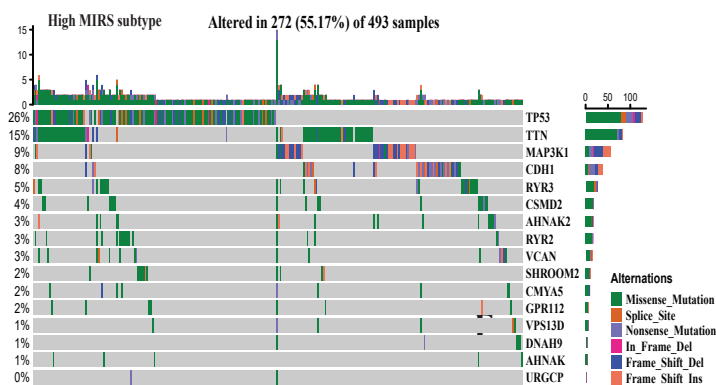
A**B****C****D****E****F****G****H**



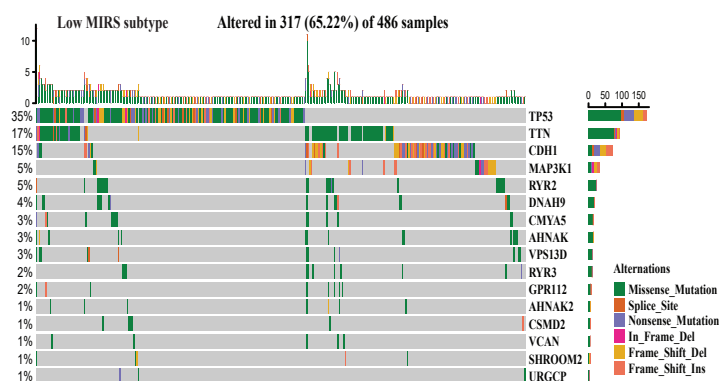
A



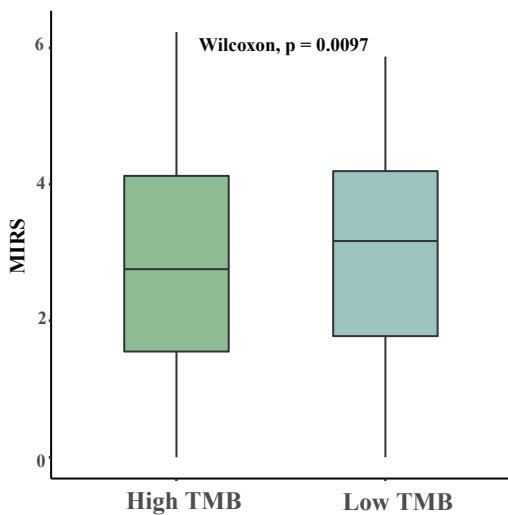
B



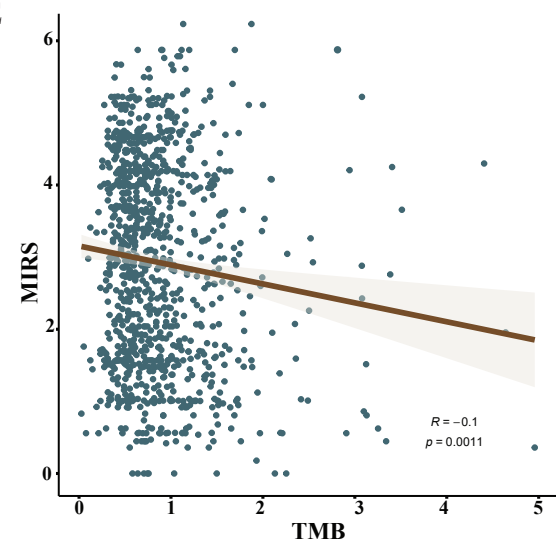
C

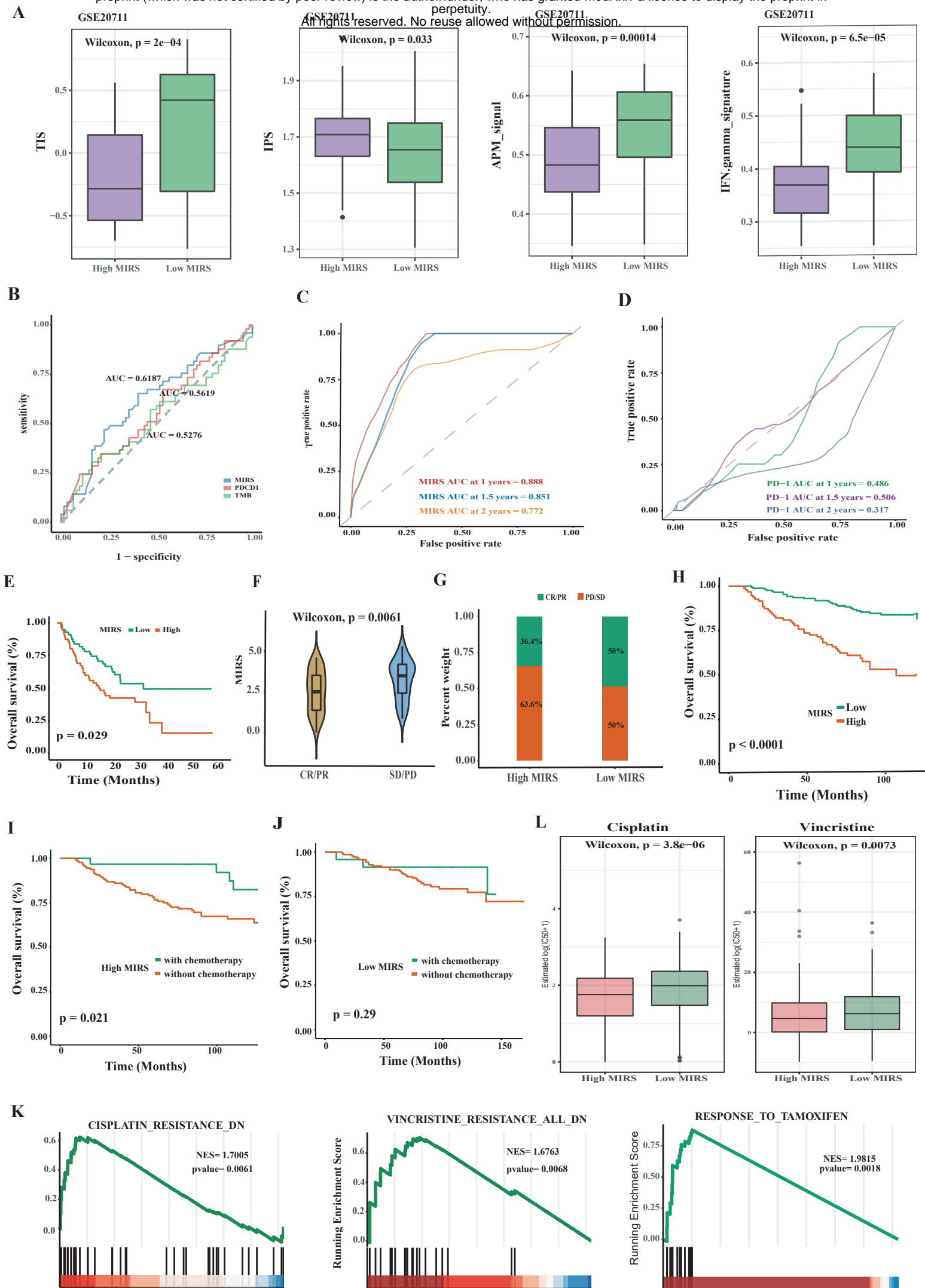


D

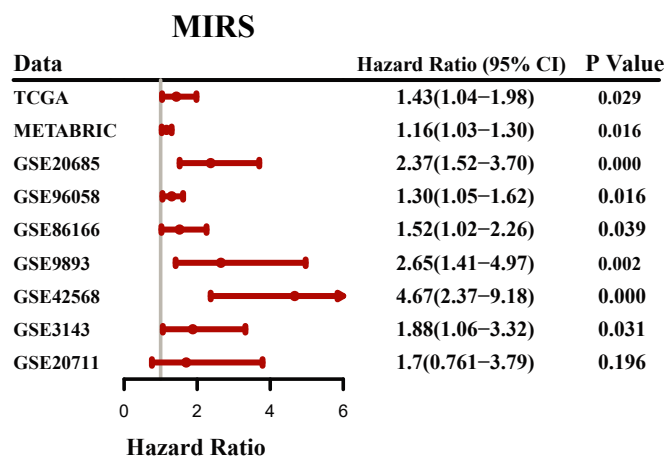


E

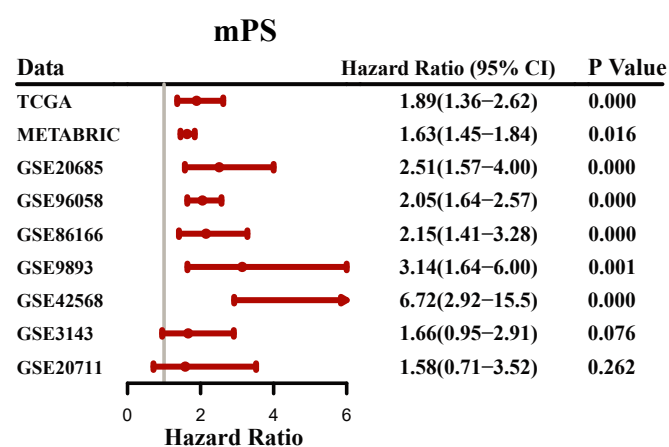




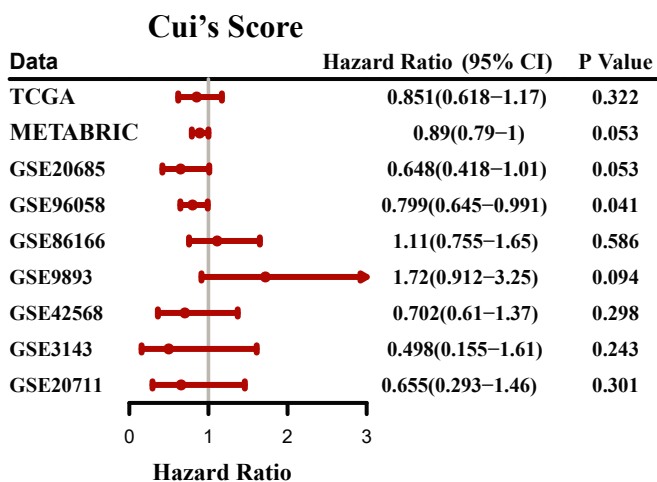
A



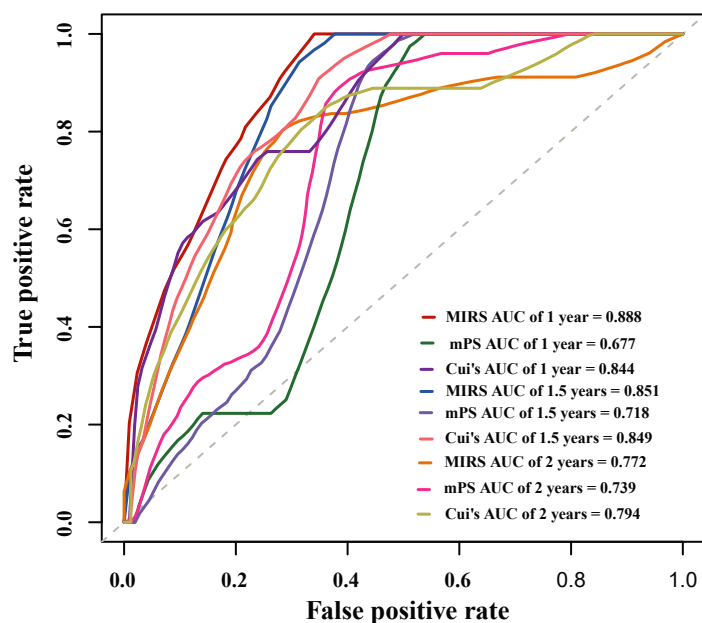
B



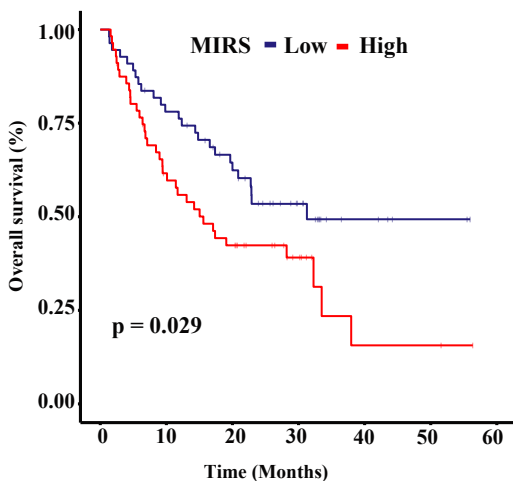
C



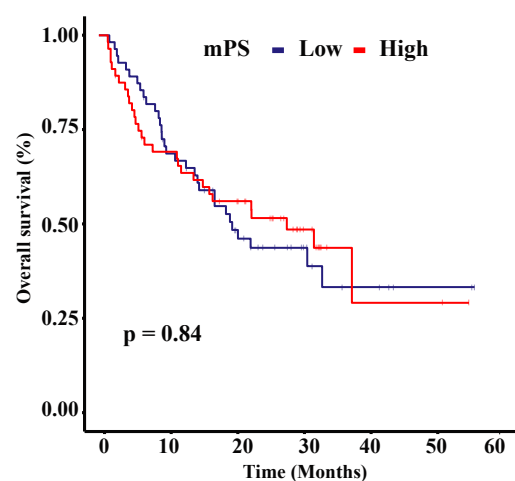
D



E



F



G

