

Gene expression-based identification of prognostic markers in lung adenocarcinoma

Annette Salomonsson¹, Daniel Ehinger^{1,2}, Mats Jönsson¹, Johan Botling^{3,4}, Patrick Micke⁴,
Hans Brunnström^{5,6}, Johan Staaf^{1,7}, Maria Planck^{1,8,9*}

¹Department of Clinical Sciences Lund, Division of Oncology, Lund University, Lund, Sweden

²Department of Genetics, Pathology, and Molecular Diagnostics, Skåne University Hospital, Helsingborg, Sweden

³Department of Laboratory Medicine, Institute of Biomedicine, Sahlgrenska Academy at University of Gothenburg, Gothenburg, Sweden

⁴Department of Immunology, Genetics, and Pathology, Uppsala University, Uppsala, Sweden

⁵Department of Clinical Sciences Lund, Division of Pathology, Lund University, Lund, Sweden

⁶Department of Genetics, Pathology, and Molecular Diagnostics, Skåne University Hospital, Lund, Sweden

⁷Department of Laboratory Medicine, Division of Translational Cancer Research, Lund University, Lund, Sweden

⁸Department of Clinical Sciences Lund, Division of Respiratory Medicine, Allergology, and Palliative Medicine, Lund University, Lund, Sweden

⁹Department of Respiratory Medicine and Allergology, Skåne University Hospital, Lund, Sweden

* Corresponding author: maria.planck@med.lu.se

Abstract

Introduction: Many studies have aimed at identifying additional prognostic tools to guide treatment choices and patient surveillance in lung cancer by assessing the expression of individual proteins through immunohistochemistry (IHC) or, more recently, through gene expression-based signatures. As a proof-of-concept, we used a multi-cohort, gene expression-based discovery and validation strategy to identify genes with prognostic potential in lung adenocarcinoma. The clinical applicability of this strategy was further assessed by evaluating a selection of the markers by IHC.

Materials and Methods: Publicly available gene expression data sets from six microarray-based studies were divided into four discovery and two validation data sets. First, genes associated with overall survival (OS) in all four discovery data sets were identified. The prognostic potential of each identified gene was then assessed in the two validation data sets, and genes associated with OS in both data sets were considered as potential prognostic markers. Finally, IHC for selected potential prognostic markers was performed in two independent and clinically well-characterized lung cancer cohorts.

Results and Conclusions: The gene expression-based strategy identified 19 genes with correlation to OS in all six data sets. Out of these genes, we selected Ki67, MCM4 and TYMS for further assessment with IHC. Although an independent prognostic ability of the selected markers could not be confirmed by IHC, this proof-of-concept study demonstrates that by employing a gene expression-based discovery and validation strategy, potential prognostic markers can be identified and further assessed by a technique universally applicable in the clinical practice. The concept of studying potential prognostic markers through gene expression-based strategies, with a subsequent evaluation of the clinical utility, warrants further exploration.

1 Introduction

2 Despite recent advancements in the understanding and treatment of lung cancer, the prognosis
3 is poor and lung cancer continues to be the leading cause of cancer-related mortality worldwide
4 [1]. Non-small cell lung cancer (NSCLC) accounts for the majority of cases, with
5 adenocarcinoma (AC) as the most frequent histological subtype [2]. Disease stage and patient's
6 performance status are the most well-established and clinically used prognostic factors.
7 Patients with localized disease can be candidates for curatively intended surgery. However,
8 also among these patients, there is a substantial mortality and a 5-year survival rate of only
9 around 60% [3]. For patients with tumors of TNM stage 1B or higher, post-operative adjuvant
10 chemotherapy leads to a decreased risk of recurrence and improved survival [3]. Since recently,
11 the addition of targeted therapy (for *EGFR*-mutated cases) or immunotherapy (for *EGFR*- and
12 *ALK*-negative tumors of stage 2 or higher that show high expression of PDL1), is also
13 recommended [4,5]. The varied outcome for surgically treated patients, also within the same
14 disease stage, illustrates a need for additional tools to guide treatment choices and patient
15 surveillance. With the emergence of yet more strategies involving immunotherapy or targeted
16 therapy in the preoperative and/or postoperative curative setting, treatment decisions will
17 become more and more complex [6-8]. Many studies have aimed at identifying prognostic
18 markers, often by assessing the expression of individual proteins through
19 immunohistochemistry (IHC). However, despite a plethora of IHC studies in lung cancer, no
20 such markers are in clinical use today. More recently, gene expression-based lung cancer
21 signatures turned out as promising prognosticators that deserve further validation for patient
22 benefit in clinical praxis, but the feasibility of such costly and labor-intensive analyses in a
23 clinical routine remain disputable [9]. In this proof-of-concept study, we hypothesized that by
24 utilizing a multi-cohort, gene expression-based discovery and validation strategy, we could
25 identify genes with prognostic potential in lung adenocarcinoma. Subsequently, to increase a

26 potential clinical applicability of this strategy for identifying prognostic markers, a selection of
27 the identified markers was further assessed by IHC.

28 **Materials and methods**

29 All analytical steps, and the public and in-house lung cancer cohorts that we used, are outlined
30 in Figure 1. In brief, we explored six different publicly available gene expression data sets, in
31 total comprising 1,167 lung adenocarcinomas, to identify and validate markers with consistent
32 correlation to overall survival (OS), and then evaluated a selection of these markers by IHC in
33 two independent cohorts.

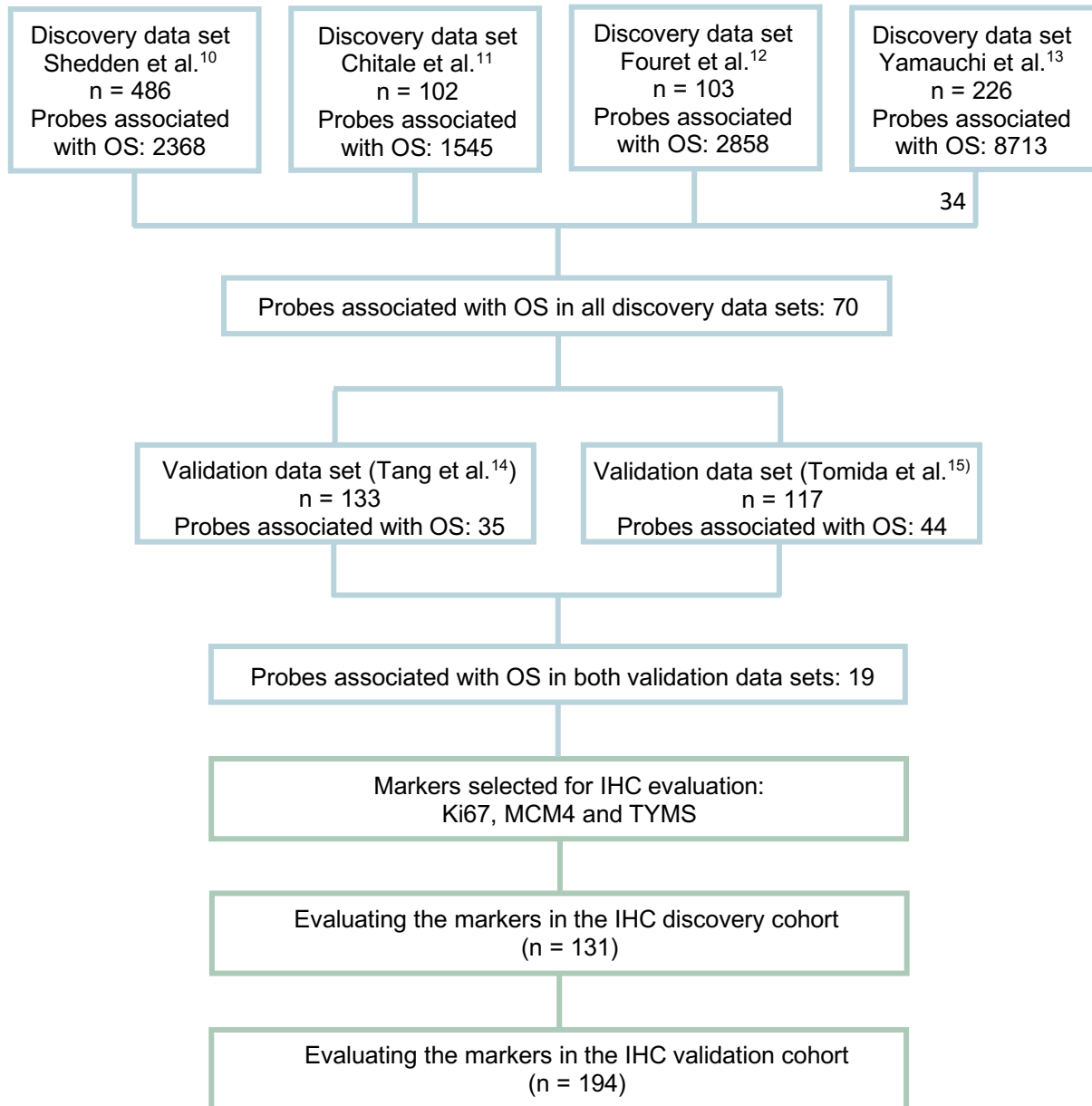


Figure 1. Schematic image of the gene expression-based strategy for identification of prognostic markers and subsequent IHC evaluation. For each probe (matching to a gene) in the four discovery data sets, the median gene expression value was used to divide the samples into two groups (high/low). The log-rank test was employed to identify probes significantly associated with OS (P -value < 0.05). Results from the four discovery data sets were then compared and probes that were significantly associated with OS in all four data sets were tested in the same manner in two validation data sets. The genes significantly associated with OS in both data sets were classified as potential prognostic markers. Out of these genes, three were selected for IHC evaluation in two patient cohorts. One of the cohorts was used as an IHC discovery cohort where optimal cut-offs for each marker were selected. These cut-offs were then applied to the cases in the IHC validation cohort.

Abbreviations: OS = overall survival, IHC = immunohistochemistry.

35 **Gene expression-based discovery and validation**

36 Publicly available transcriptomic profiles and matched survival data were obtained from six
37 microarray-based lung cancer studies [10-15]. Samples with AC histology (n = 1,167) were
38 chosen for further comparisons of the gene expression data, which were processed as
39 previously described [16]. Four of the data sets were used in the discovery step, all based on
40 the Affymetrix platform [10-13]. For each probe (matching to a gene) in the data sets, the
41 median gene expression value was used to divide the samples into two groups (high/low). Then,
42 the log-rank test was employed to identify probes significantly associated with OS (P -value <
43 0.05). Results from the four discovery data sets were then compared and probes that were
44 significantly associated with OS in all four data sets advanced to the validation step, in which
45 the probes generated from the discovery step were tested in the same manner in two validation
46 data sets, based on non-Affymetrix platforms [14-15]. The genes significantly associated with
47 OS in both data sets were then classified as potential prognostic markers.

48 **Immunohistochemical evaluation of potential prognosticators**

49 Among the potential prognostic markers obtained by our discovery and validation strategy, we
50 selected three genes (*Ki67*, *MCM4* and *TYMS*, as further discussed below) for further IHC
51 evaluation of the corresponding proteins. Immunohistochemical staining was performed in two
52 independent and clinically well-characterized lung cancer cohorts. The first cohort was used as
53 an IHC discovery cohort for identification of cut-offs for classifying samples as having a low
54 or high expression of each marker, and the second cohort was used as an IHC validation cohort,
55 where these cut-offs were then applied.

56 The IHC discovery cohort was based on the “Southern Swedish Lung Cancer Study” which
57 prospectively included patients with primary lung cancer who underwent surgical treatment at

58 the Skåne University Hospital, Lund, Sweden, in 2005 – 2011 [17]. The present study included
59 131 AC and 69 squamous cell carcinomas (SqCC). The IHC validation cohort, was based on
60 194 AC cases from the "Uppsala NSCLC II cohort" which included patients with primary lung
61 cancer who underwent surgical treatment at the University Hospital in Uppsala, Sweden, in
62 2006 – 2010 [18,19]. Patient characteristics and clinicopathological data in the two IHC cohorts
63 were described previously [20]. The studies were approved by the Regional Ethical Review
64 Board in Lund (Dnr 2004/762 and 2008/702) and Uppsala (Dnr 2012/532) and conducted in
65 adherence with the Declaration of Helsinki.

66 Only patients that were surgically treated for primary NSCLC tumors were included. Patients
67 receiving neoadjuvant treatment, or chemotherapy for another malignancy six months before
68 surgery, were excluded from the present study. All cases were previously reviewed by two
69 pathologists (HB and PM), who updated the diagnoses in accordance with the 2015 WHO
70 classification and TNM 7 and who confirmed all changes from the original diagnoses [17,19,
71 21, 22]. Furthermore, growth patterns were evaluated (HB) for stratification into three groups:
72 minimally invasive/predominant lepidic, predominant acinary/papillary, and mucinous or
73 predominant micropapillary/solid. Overall survival data were retrieved from the Swedish
74 Cancer Registry, to which reporting is mandatory by law. The registry was consulted on June
75 26, 2018 (the IHC discovery cohort), and on March 29, 2019 (the IHC validation cohort). For
76 one patient in the IHC validation cohort, survival data were unavailable. Analysis of
77 recurrence-free interval (RFI) was performed as previously described [20], and included 122
78 AC in the IHC discovery cohort, and 164 AC in the IHC validation cohort. Tissue microarrays
79 (TMA) were used for IHC analysis. The TMA-blocks had, for each case, three (the IHC
80 discovery cohort) or two (the IHC validation cohort) cores, 1 mm in diameter. For IHC
81 analysis, 4- μ m thick sections were stained according to Supplementary Table 1. The slides
82 were scanned and evaluated using the pathXL software (Philips, Amsterdam, The

83 Netherlands). For further analysis, we required a minimum of 200 assessable tumor cells on
84 the TMAs, with most cases having over 1000 evaluable cells.

85 All stainings were evaluated by three independent observers (AS, DE and MJ) who were
86 blinded to clinical data and patient outcome. Nuclear staining for Ki67 and MCM4 was
87 considered positive. Cytoplasmic or nuclear staining were considered positive for TYMS,
88 though only cells with visible nuclei were counted. Attention was paid to exclude stained non-
89 tumor cells. In case of varying expression of the marker between the cores within a sample, the
90 mean proportion of cells expressing the marker across all cores was assessed. For cases with
91 differences in the scoring between the evaluators, the cases were jointly reviewed, and
92 consensus was reached.

93 In the IHC discovery cohort, the fraction of viable tumor cells expressing the marker was scored
94 as 0 (0-1%), 1 (>1-10%), 2 (>10-25%), 3 (>25-50%), 4 (>50-75%) or 5 (>75%). For TYMS,
95 in addition to the recorded fraction of positive tumor cells, the staining intensity was scored as
96 0 (negative), 1 (mild), 2 (moderate), or 3 (strong) and a final score was constructed by
97 multiplying these two parameters. To establish the optimal cut-off for each marker for
98 categorizing samples into high or low expression groups, Kaplan-Meier plots with log-rank
99 tests were used and the cut-offs yielding the lowest p-values in the log-rank tests were selected.
100 Prognostic analyses were performed separately on AC and SqCC.

101 In the IHC validation cohort, tumors were scored as high or low expressors using the optimized
102 cut-offs selected in the IHC discovery cohort. Furthermore, in both cohorts, the combined
103 prognostic ability of the three markers was examined by each case receiving one point per
104 positive maker, thus resulting in a combined score ranging from 0 to 3 points.

105 **Gene expression of *Ki67*, *MCM4* and *TYMS* in the IHC validation cohort**

106 Gene expression data of *Ki67*, *MCM4*, and *TYMS* were available for 104 AC cases in the IHC
107 validation cohort. Gene expression data are available as GSE81089 and RNA sequencing
108 analysis was performed as previously described by Djureinovic et al. [23]. We tested two
109 different cut-offs for classifying samples as having low or high gene expression levels of *Ki67*,
110 *MCM4* or *TYMS* by dividing the samples into either two or three equally sized groups.

111 **Statistical analysis**

112 Kaplan-Meier plots with log-rank test were used for OS analyses and for analyses of RFI.
113 Univariable and multivariable Cox proportional hazards regression models were used for
114 further comparisons between groups. Multivariable models were adjusted for stage (I, II, III,
115 and IV), age, smoking status (current, past, or never), gender, adjuvant therapy, growth pattern,
116 and patients' performance status (the latter available for the IHC validation cohort only).
117 Spearman's rank correlation was used to assess the correlations between gene expression levels
118 of the potential prognostic genes. The Mann-Whitney U test/Wilcoxon rank-sum test and
119 Fisher's exact test were used to compare data between groups. A *P*-value < 0.05 was considered
120 statistically significant. All statistical analyses were performed using R (version 3.6.1) [24].

121

122 **Results**

123 **Gene expression-based identification of genes with prognostic** 124 **potential**

125 For 70 probes (genes), the gene expression levels were associated with OS in all four discovery
126 data sets, as schematically illustrated in Figure 1. Of these, 19 genes (listed in Table 1) were

127 associated with OS in the two gene expression data sets used for validation and were thus
128 considered as having prognostic potential in lung adenocarcinoma.

Table 1. Genes with prognostic potential identified in the gene expression-based discovery and validation step.

Gene Symbol	Gene Name
KI67	Marker of proliferation Kiel 67
MCM4	Minichromosome maintenance complex component 4
TYMS	Thymidylate synthetase
CCNA2	Cyclin A2
CCNE1	Cyclin E1
BUB1B	Budding uninhibited by benzimidazoles 1 homolog beta
DLGAP5	Discs large homolog associated protein 5
KIF14	Kinesin family member 14
NUSAP1	Nucleolar and spindle-associated protein 1
RACGAP1	Rac GTPase activating protein 1
ECT2	Epithelial cell transforming sequence 2 oncogene
ASPM	Abnormal spindle-like microcephaly-associated protein
PRC1	Protein regulator of cytokinesis 1
BTG2	B-cell translocation gene 2
HLF	Hepatic leukemia factor
GDF10	Growth differentiation factor 10
CTTN	Cortactin
COL4A3	Collagen, type IV, alpha 3
CIRBP	Cold inducible RNA binding protein

129 In the two validation data sets, correlation plots (Figure 2) showed a generally strong
130 correlation in gene expression levels between the 19 genes. Broadly, the 19 genes could be
131 divided into two groups that were inversely correlated to each other. By using Kaplan-Meier
132 plots, it was demonstrated that high gene expression levels were associated with worse outcome
133 for one group of genes, while low expression levels were associated with worse outcome for
134 the other group, as exemplified in Supplementary Figure 1.

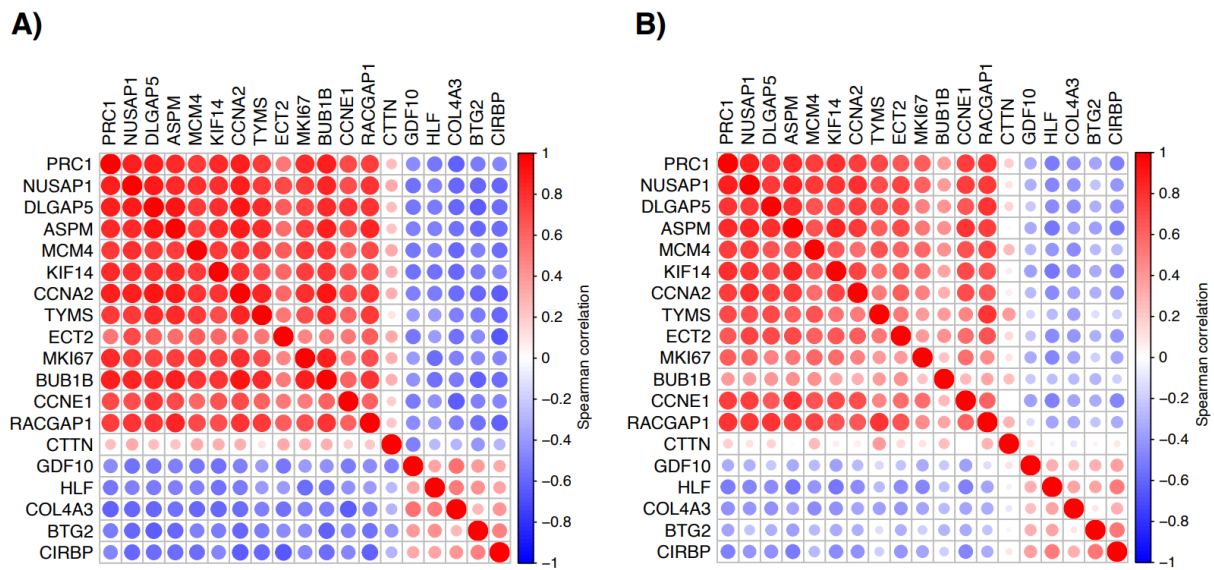


Figure 2. Spearman correlation of gene expression levels of the 19 candidate genes in the two validation data sets. (A) Tomida et al.¹⁵, (B) Tang et al.¹⁴. If multiple probes were available for a gene then the probes with largest standard deviation was chosen to represent the gene.

135 Evaluation of the clinical utility

136 To further explore the potential clinical utility of our gene-expression based strategy for
 137 identifying prognostic markers, a selection of the identified markers was further assessed by
 138 IHC. Considering clinical practicability, we preferred markers where high expression was
 139 associated with worse outcome. Also, availability of reliable antibodies was considered in
 140 selecting candidate genes for further analyses with IHC. In the two validation data sets, the
 141 patients classified as having high gene expression levels of *Ki67*, *MCM4*, and *TYMS* did not
 142 fully overlap (Supplementary Figure 2), thereby suggesting that the three markers could
 143 possibly complement each other. Based on these considerations, we chose these three genes
 144 for further analyses with IHC in two independent lung cancer cohorts. Representative
 145 microscopic images of the stainings for *Ki67*, *MCM4* and *TYMS* are shown in Supplementary
 146 Figure 3.

147 **Protein expression of Ki67, MCM4 and TYMS in the IHC discovery cohort**

148 For AC in the IHC discovery cohort, the protein expression could be evaluated for Ki67 in all
149 131 cases, for MCM4 in 129 cases and for TYMS in 120 cases. For SqCC, 68 cases could be
150 evaluated for Ki67 and MCM4, and 60 cases could be evaluated for TYMS. When comparing
151 the IHC scores of the respective markers with regards to histology, SqCC had significantly
152 higher expression of Ki67, MCM4 and TYMS compared to AC (Wilcoxon test, P -value <0.05
153 in all three tests). Therefore, further analyses were performed on AC and SqCC separately.

154 For Ki67, a cut-off of $>10\%$ positive tumor cells most clearly identified prognostic groups in
155 the OS analysis among the AC cases and was therefore chosen for identification of samples
156 with a low or high expression (Figure 3A). By applying this cut-off, 74 AC cases (56%) were
157 classified as having a high Ki67 protein expression. A cut-off of $>75\%$ positive tumor cells
158 was selected for MCM4 among the AC cases in the OS analysis, which resulted in 15 cases
159 (12%) identified as having a high MCM4 expression (Figure 3B). For TYMS, a score (obtained
160 by multiplying fraction and intensity) of >2 p was chosen for identification of AC samples with
161 a high TYMS expression in the OS analysis, which resulted in 19 cases (16%) classified as
162 having a high expression of TYMS (Figure 3C). The prognostic value of these cut-offs in the
163 RFI-analysis for Ki67, MCM4, and TYMS are shown in Supplementary Figure 4.

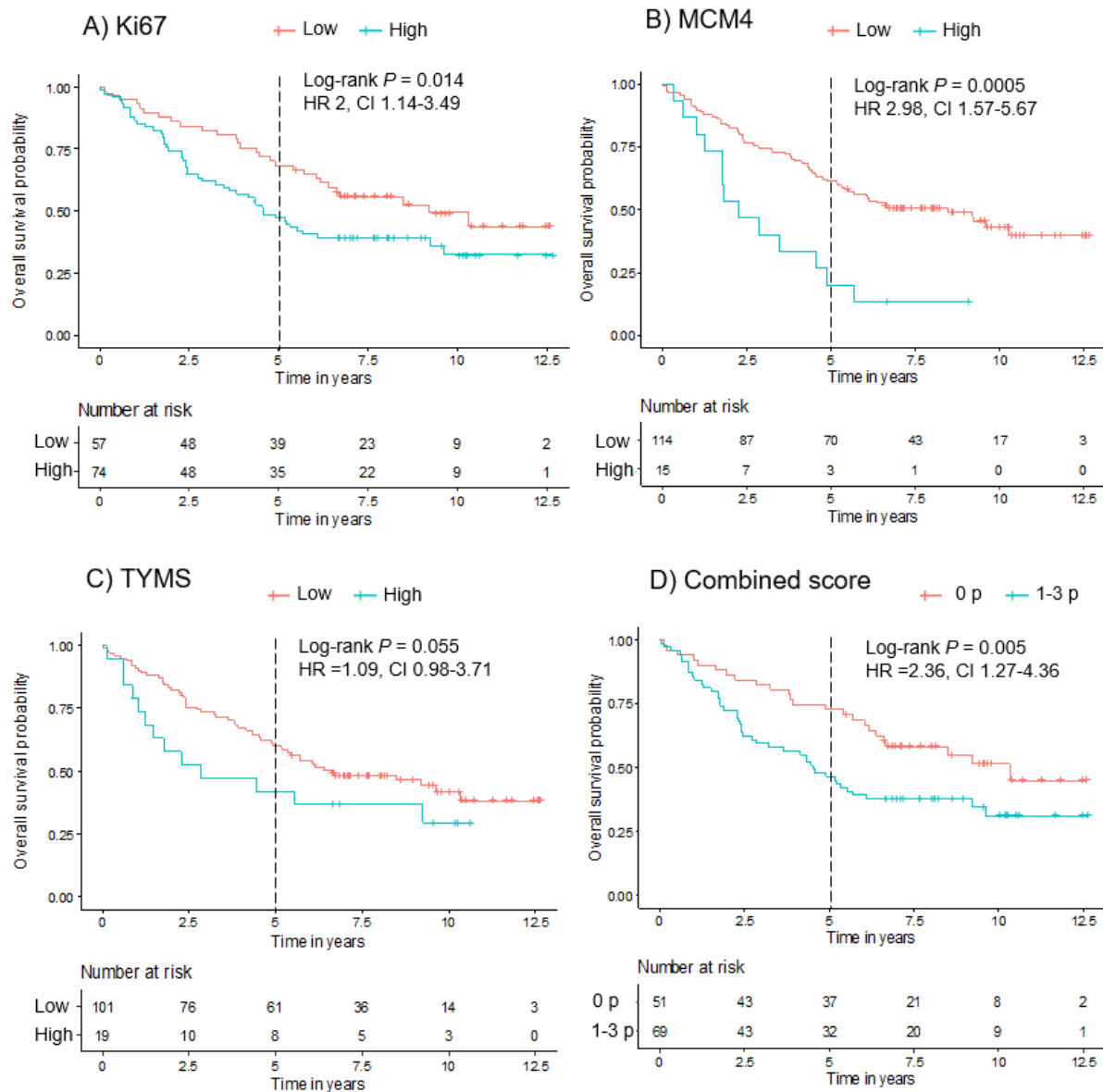


Figure 3. Prognostic value of Ki67 (A), MCM4 (B), TYMS (C), and combined score (D), on overall survival in the IHC discovery cohort.

164 For SqCC, no prognostically meaningful subgroups could be identified for any of the markers
165 in the OS analysis, and therefore the SqCC cases were omitted from further analyses and were
166 not evaluated in the IHC validation cohort.

167 For the selected cut-offs in the IHC discovery cohort, smoking status was significantly
168 associated with Ki67 expression levels, with low expression often observed in never smokers
169 and high expression in current smokers (Fisher's exact test, $P = 0.007$). For MCM4 and TYMS,
170 no associations with smoking were found. Furthermore, expression of Ki67 was associated
171 with AC growth pattern, as samples with a high expression were more frequently found in the
172 group with mucinous or predominant micropapillary/solid pattern (Fisher's test, $P = 0.03$). For
173 MCM4 and TYMS, no associations with growth patterns were found. For age, gender, stage,
174 and number of cases receiving adjuvant treatment, no associations between these parameters
175 and patients with a high or low expression of Ki67, MCM4 or TYMS, respectively, were found.

176 **Protein expression of Ki67, MCM4 and TYMS in the IHC validation** 177 **cohort**

178 In the IHC validation cohort, the protein expression of Ki67, MCM4, and TYMS could be
179 assessed in 159, 178, and 146 AC cases, respectively. By applying the identified cut-offs from
180 the IHC discovery cohort for respective gene, high expression was found in 91 cases (57%) for
181 Ki67, in 17 cases (10%) for MCM4, in and 17 cases (12%) for TYMS. High expression of
182 Ki67, MCM4, and TYMS was associated with male gender (Fisher's test, $P < 0.05$ in all three
183 tests) and high expression of Ki67 was associated with more advanced stages (stage III)
184 (Fisher's test, $P = 0.007$). Furthermore, the expression of Ki67 and MCM4 was associated with
185 growth pattern, as proportionally more cases with a high expression were found in the group

186 with mucinous or predominant micropapillary/solid pattern compared to cases with a low
187 expression, where proportionally more cases were minimally invasive/lepidic or
188 acinary/papillary (Fisher's test, $P < 0.05$). The expression of Ki67 was associated with smoking
189 as there were proportionally more never smokers among cases with a low expression compared
190 to cases with a high expression (Fisher's test, $P < 0.001$). Apart from these findings, no other
191 associations between age, gender, stage, smoking status, growth pattern, WHO performance
192 status, and number of cases receiving adjuvant treatment and patients with a high or low
193 expression of Ki67, MCM4 or TYMS, respectively, were found.

194 High protein expression of Ki67 was associated with a worse prognosis in the 5-year OS
195 analysis (log-rank test, $P = 0.0002$, Figure 4A). In the univariable Cox proportional hazards
196 regression model, Ki67 expression was significantly associated with prognosis (HR 2.54, 95%
197 CI 1.54-4.21). However, these results did not remain statistically significant in the
198 multivariable model adjusted for stage, growth pattern, age, gender, smoking, WHO
199 performance status, and adjuvant treatment (HR 1.29, 95% CI 0.68-2.45). In the RFI analysis,
200 patients with a high expression of Ki67 had a higher rate of recurrence (log-rank test, $P =$
201 0.0003, Supplementary Figure 5A). For MCM4 and TYMS, no statistically significant
202 associations between protein expression and survival or RFI could be demonstrated (Figure
203 4B/C and Supplementary Figure 5B/C).

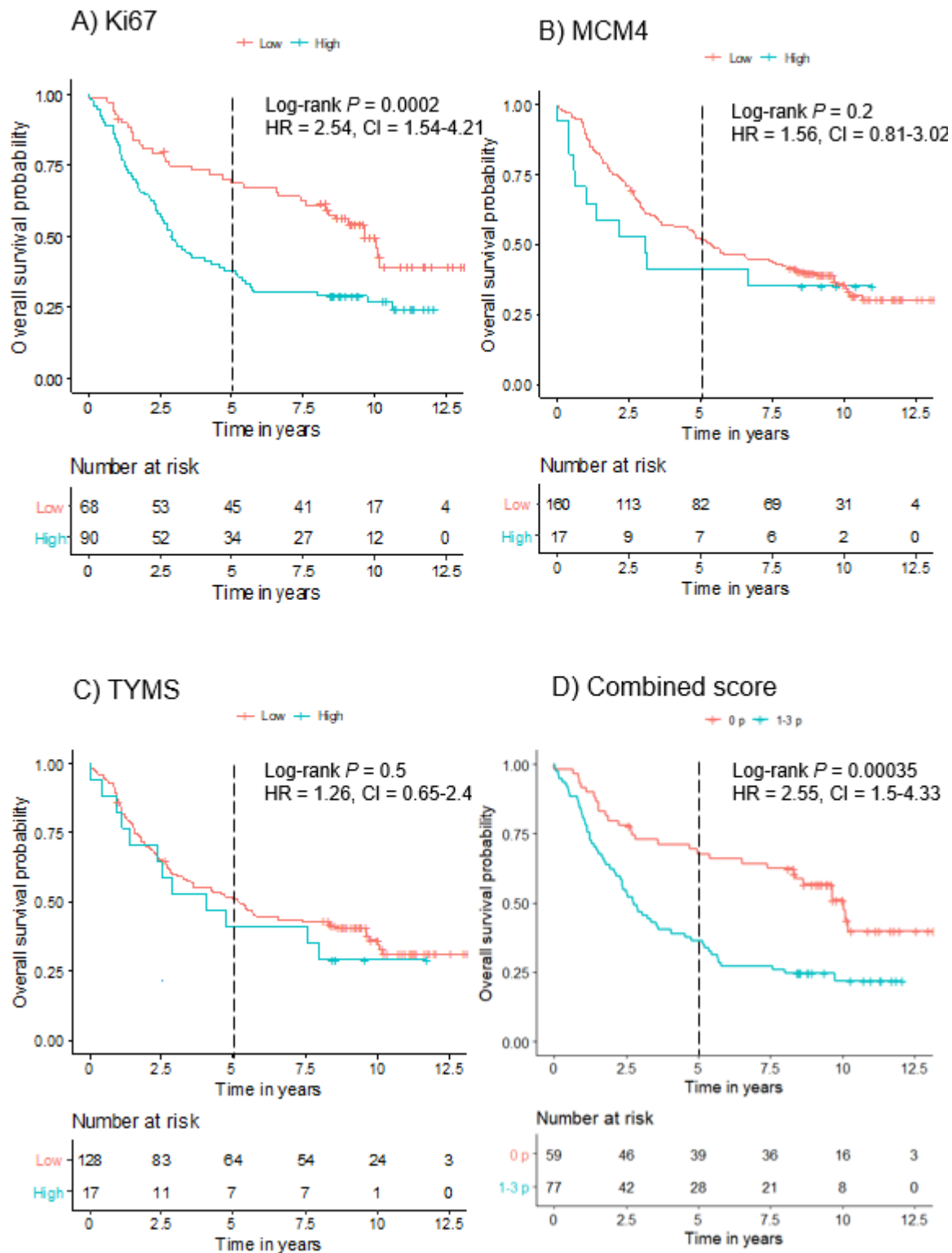


Figure 4. Prognostic value of Ki67 (A), MCM4 (B), TYMS (C), and combined score (D), on overall survival in the IHC validation cohort.

204 **Combining markers for improved prognostication**

205 Considering the gene expression findings, we also examined the combined prognostic ability
206 of the three markers. The number of patients positive for each marker, and the overlap between
207 these, are presented in Supplementary Figure 6. All three markers did not independently add
208 prognostic information in the combined score as there were no patients that were positive for
209 only MCM4 in both cohorts, and TYMS only added one (IHC validation cohort) or two (IHC
210 discovery cohort) cases to the high-risk group (more than one point) compared to Ki67 alone.
211 In both IHC cohorts, cases that were positive for one or more markers had a worse prognosis
212 in the 5-year OS analysis (Figure 3D and 4D) and a higher rate of recurrence (Supplementary
213 Figure 4D and 5D) compared to cases that were negative for all three markers. However, in the
214 IHC validation cohort, these associations did not remain statistically significant in the
215 multivariate model.

216 **Gene expression of *Ki67*, *MCM4* and *TYMS* in the IHC validation cohort**

217 Gene expression levels of *Ki67*, *MCM4*, and *TYMS* were available for 104 AC cases in the IHC
218 validation cohort. Out of these 104 cases, IHC data were missing for 20 cases for Ki67, four
219 cases for MCM4, and 28 cases for TYMS. For all three markers, a correlation between gene
220 expression levels and IHC classification (low or high expression) could be observed (Wilcoxon
221 test, $P < 0.01$ all three tests, Supplementary Figure 7).

222 The prognostic value of *Ki67*, *MCM4*, and *TYMS* gene expression levels in the IHC validation
223 cohort were evaluated by dividing the samples into two or three equally sized groups. For *Ki67*
224 and *TYMS*, no statistically significant differences between the groups could be identified,
225 although potentially prognostic subgroups could be visualized in the Kaplan-Meier plots
226 (Supplementary Figure 8). For *MCM4*, patients with high expression levels had a worse

227 prognosis compared to patients with low expression levels in the 5-year OS analysis (log-rank
228 test, $P = 0.004$ when dividing the samples into two groups, and $P = 0.02$ when dividing the
229 samples into three groups, Supplementary Figure 8).

230 **Discussion**

231 Studies of potential prognosticators in lung cancer are often based on the immunohistochemical
232 expression of protein markers or on gene expression-based prognostic signatures. There are no
233 prognostic IHC markers in clinical use for lung cancer today, and the reproducibility and
234 clinical benefit of gene expression-based prognostic signatures needs to be thoroughly
235 validated before being implemented in a clinical setting [9]. However, gene expression can
236 already now be employed as a research tool for identifying potential prognosticators. As a
237 proof-of-concept, we identified markers with prognostic impact in lung adenocarcinoma
238 through a gene expression-based, multi-cohort discovery and validation strategy, where the
239 expression of 19 genes was correlated to survival in six independent studies based on global
240 gene expression profiling arrays, published in reputable journals [10-15]. We also selected
241 markers identified by this strategy for further evaluation by IHC, a method more adapted to the
242 current clinical setting, thus underlining a potential future clinical utility.

243 Several of the 19 potential prognostic markers that we identified in our expression-based
244 discovery and validation strategy (Table 1) are linked to proliferation, and either higher
245 expression levels (e.g. for *Ki67*) or lower expression levels (e.g. *BTG2*) are correlated to poor
246 outcome [25,26]. The prognostic impact of proliferation has long been recognized in many
247 types of cancer, and many IHC-based markers target proliferation [27]. Furthermore, it has
248 been suggested that proliferation-associated genes are key components in gene-expression-
249 derived adenocarcinoma prognostic phenotypes [28]. Accordingly, genes linked to

250 proliferation proved important in our current multi-cohort approach to associate gene
251 expression with patient overall survival. As illustrated in correlation plots for the two validation
252 data sets (Figure 2), the 19 candidate genes could broadly be divided into two groups that were
253 inversely correlated to each other. In the larger of these two groups, all genes are directly
254 implicated in proliferation.

255 To further explore the potential clinical applicability of our gene-expression based strategy for
256 identifying prognostic markers, we selected three of the 19 markers for further assessment. The
257 gene expression levels of the three selected markers (Ki67, MCM4, and TYMS) were
258 correlated to each other in the two validation data sets, and high expression was associated with
259 worse prognosis. Furthermore, as illustrated in Supplementary Figure 2, the markers could
260 possibly complement each other in identifying high-risk patients. Based on the gene expression
261 correlation analyses, it could be hypothesized that alternative gene selections would have
262 resulted in similar results. IHC has the advantage of being an accessible and applicable method
263 in the clinical routine and, for all three markers, there are reliable antibodies available and the
264 genes have a recognized prognostic potential in lung adenocarcinoma [26, 29-32]. Difficulties
265 in standardization and reproducibility across IHC studies remain challenges, but may improve
266 with the emergence of digital image analysis [33]. However, other methods suitable for
267 formalin-fixed paraffin embedded tissue, such as RNA-based NanoString technology or
268 quantitative reverse transcription polymerase chain reaction (qRT-PCR), might as well have
269 been considered.

270 The robustness of the three selected markers was demonstrated through a clear correlation
271 between gene expression levels and IHC classification in both IHC cohorts. Furthermore, when
272 assessing the prognostic value of gene expression levels of the three markers in the IHC
273 validation cohort, potentially prognostic subgroups could be visualized in the Kaplan-Meier
274 plots, although not statistically significant for *Ki67* and *TYMS*. For the IHC stainings, we were

275 able to test the consistency of the chosen cut-offs by using two independent lung cancer cohorts.
276 The first cohort was used to establish cut-off values for categorizing samples into groups (high
277 and low expression). Subsequently, these cut-offs were applied when evaluating the cases in
278 the IHC validation cohort. However, based on these cut-offs, we could only confirm the
279 prognostic ability of Ki67, although the association did not remain prognostic in the
280 multivariable model. For MCM4 and TYMS, the cut-offs chosen in the IHC discovery cohort
281 identified only a small proportion of the samples with a worse prognosis. It is possible that a
282 lower cut-off for these two markers, identifying more patients and more resembling the cut-off
283 chosen for Ki67, would have performed better in the validation cohort.

284 Only AC cases were included in the gene expression-based discovery and validation step for
285 detection of potential prognostic markers. Therefore, it was not unexpected that we were unable
286 to define prognostic subgroups among the SqCC cases for any of the three prognostic markers
287 selected for our IHC validation. The impact of histological subtyping for accurate choice of
288 prognosticators has been demonstrated also for other potential prognostic markers in lung
289 cancer and, indeed, prognostic gene expression signatures developed in lung cancer have often
290 been derived from specific histological subgroups [9,20]. The SqCC samples had significantly
291 higher protein expression of Ki67, MCM4 and TYMS compared to AC. As the three markers
292 are related to proliferation, these results imply that most of the SqCC cases were highly
293 proliferative, and thereby it becomes more challenging to find meaningful subgroups based on
294 these markers.

295 As illustrated in Supplementary Figure 2, more high-risk patients could be identified by
296 combining the three markers compared to using one single marker for gene expression levels
297 of *Ki67*, *MCM4* and *TYMS* in the two validation data sets. However, for the IHC evaluations
298 in our study, the combined score was in both cohorts dependent on the prognostic ability of
299 *Ki67* alone, and the patients identified by MCM4 and TYMS overlapped with the patients

300 identified by Ki67 (Supplementary Figure 6). These results possibly reflect that we set the cut-
301 offs for MCM4 and TYMS at a level where these markers identified a too small proportion of
302 the patients. Also, all three selected markers are associated with proliferation and, as such, may
303 be redundant as they assess the same cancer characteristic. Possibly, a combination of markers
304 that assess different biological processes could better identify additional high-risk patients.
305 However, in a previous study by Grinberg et al., a prognostic model based on a biomarker
306 panel consisting of five protein markers with diverse biological functions was developed,
307 where each marker also was associated with prognosis in gene expression data sets [34]. When
308 the model was applied to the validation cohort, it failed to improve survival prediction beyond
309 clinical parameters alone, thus questioning the prognostic impact of protein biomarkers and
310 further stresses the difficulties of implementing additional prognosticators into clinical
311 practice.

312 Our study has several limitations. Out of the 19 potential prognostic markers generated in the
313 gene expression-based step, we selected three for further evaluation with IHC in this proof-of-
314 concept study. It is possible that choosing other, or more, markers would have yielded a
315 different result. Ideally, more AC cases in the IHC-cohorts would have permitted more
316 extensive evaluations and subgroup analyses of the prognostic value of the markers. Also, in
317 the IHC-cohorts, the SqCC cases were few and results should be interpreted with care. The
318 markers were further evaluated by using TMAs instead of whole tumor sections, which could
319 have an impact on the validity of the results, particularly when assessing markers with unknown
320 intra-tumoral heterogeneity. The cut-off values were determined by using log-rank tests in the
321 IHC discovery cohort, and the threshold with the lowest p-value from these log-rank tests was
322 considered the optimal cut-off which were then applied to the IHC validation cohort. It is
323 conceivable that another method for identifying an optimal cut-off would have resulted in a
324 different, and perhaps better balanced, cut-off for the markers. The lack of consensus in how

325 to set cut-off values highlight some of the challenges with conducting IHC-based prognostic
326 studies.

327 To summarize, through our gene expression-based discovery and validation strategy, we
328 identified 19 genes with prognostic potential in lung adenocarcinoma and assessed three of
329 these markers further by IHC. In conclusion, this proof-of-concept study demonstrates that a
330 gene-expression based strategy for identifying prognostic markers, combined with a
331 subsequent evaluation of the clinical utility, is a justified approach that warrants further
332 exploration.

References

1. Bray F, Laversanne M, Sung H, Ferlay J, Siegel RL, Soerjomataram I, et al. Global cancer statistics 2022: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries. *CA Cancer J Clin*. 2024 May-Jun;74(3):229-263.
2. Lewis DR, Check DP, Caporaso NE, Travis WD, Devesa SS. US lung cancer trends by histologic type. *Cancer*, 2014. 120(18): p. 2883-92.
3. Burdett S, Pignon JP, Tierney J, Tribodet H, Stewart L, Le Pechoux C, et al. Adjuvant chemotherapy for resected early-stage non-small cell lung cancer. *Cochrane Database Syst Rev*, 2015(3): p. Cd011430.
4. Felip E, Altorki N, Zhou C, Csöszi T, Vynnychenko I, Goloborodko O, et al. Adjuvant atezolizumab after adjuvant chemotherapy in resected stage IB-IIIa non-small-cell lung cancer (IMpower010): a randomised, multicentre, open-label, phase 3 trial. *Lancet*, 2021. 398(10308): p. 1344-1357.
5. Tsuboi M, Herbst RS, John T, Kato T, Majem M, Grohé C, et al., Overall Survival with Osimertinib in Resected EGFR-Mutated NSCLC. *N Engl J Med*, 2023. 389(2): p. 137-147.
6. Wu YL, Dziadziuszko R, Ahn JS, Barlesi F, Nishio M, Lee DH, et al. Alectinib in Resected ALK-Positive Non-Small-Cell Lung Cancer. *N Engl J Med*, 2024. 390(14): p. 1265-1276.
7. Provencio M, Nadal E, González-Larriba JL, Martínez-Martí A, Bernabé R, Bosch-Barrera J, et al. Perioperative Nivolumab and Chemotherapy in Stage III Non-Small-Cell Lung Cancer. *N Engl J Med*, 2023. 389(6): p. 504-513.
8. Wakelee H, Liberman M, Kato T, Tsuboi M, Lee SH, Gao S, et al. Perioperative Pembrolizumab for Early-Stage Non-Small-Cell Lung Cancer. *N Engl J Med*, 2023. 389(6): p. 491-503.
9. Tang H, Wang S, Xiao G, Schiller J, Papadimitrakopoulou V, Minna J, et al. Comprehensive evaluation of published gene expression prognostic signatures for biomarker-based lung cancer clinical studies. *Ann Oncol*, 2017. 28(4): p. 733-740.
10. Shedden K, Taylor JM, Enkemann SA, Tsao MS, Yeatman TJ, Gerald WL, et al. Gene expression-based survival prediction in lung adenocarcinoma: a multi-site, blinded validation study. *Nat Med*, 2008. 14(8): p. 822-7.
11. Chitale D, Gong Y, Taylor BS, Broderick S, Brennan C, Somwar R, et al. An integrated genomic analysis of lung cancer reveals loss of DUSP4 in EGFR-mutant tumors. *Oncogene*, 2009. 28(31): p. 2773-83.

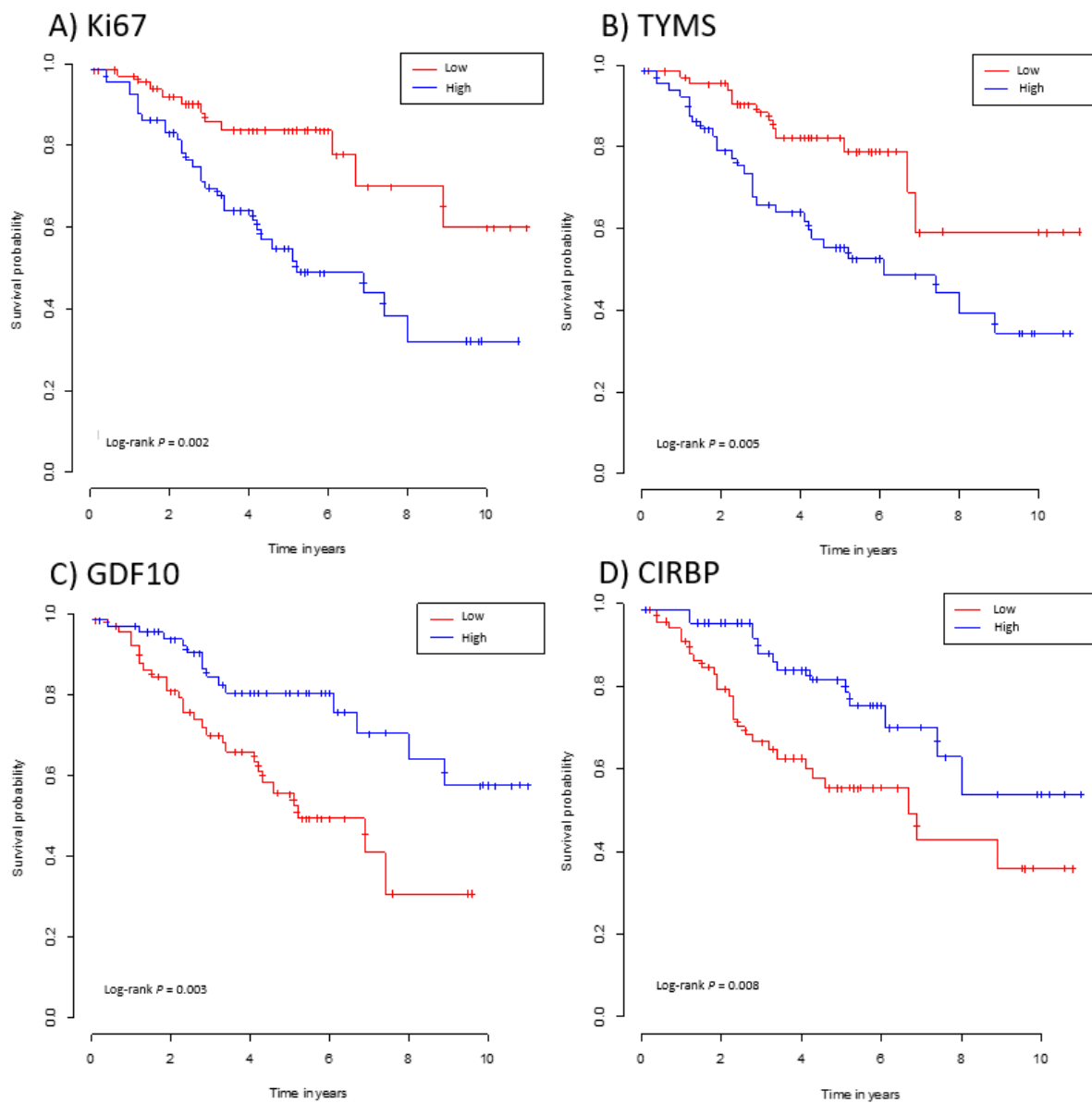
12. Fouret R, Laffaire J, Hofman P, Beau-Faller M, Mazieres J, Validire P, et al. A comparative and integrative approach identifies ATPase family, AAA domain containing 2 as a likely driver of cell proliferation in lung adenocarcinoma. *Clin Cancer Res*, 2012. 18(20): p. 5606-16.
13. Yamauchi M, Yamaguchi R, Nakata A, Kohno T, Nagasaki M, Shimamura T, et al. Epidermal growth factor receptor tyrosine kinase defines critical prognostic genes of stage I lung adenocarcinoma. *PLoS One*, 2012. 7(9): p. e43923.
14. Tang H, Xiao G, Behrens C, Schiller J, Allen J, Chow CW, et al. A 12-gene set predicts survival benefits from adjuvant chemotherapy in non-small cell lung cancer patients. *Clin Cancer Res*, 2013. 19(6): p. 1577-86.
15. Tomida S, Takeuchi T, Shimada Y, Arima C, Matsuo K, Mitsudomi T, et al. Relapse-related molecular signature in lung adenocarcinomas identifies patients with dismal prognosis. *J Clin Oncol*, 2009. 27(17): p. 2793-9.
16. Karlsson A, Ringnér M, Lauss M, Botling J, Micke P, Planck M, et al. Genomic and transcriptional alterations in lung adenocarcinoma in relation to smoking history. *Clin Cancer Res*, 2014. 20(18): p. 4912-24.
17. Brunnström H, Johansson L, Jirstrom K, Jönsson M, Jönsson P, Planck M. Immunohistochemistry in the differential diagnostics of primary lung cancer: an investigation within the Southern Swedish Lung Cancer Study. *Am J Clin Pathol*, 2013. 140(1): p. 37-46.
18. La Fleur L, Falk-Sörqvist E, Smeds P, Berglund A, Sundström M, Mattsson JS, et al. Mutation patterns in a population-based non-small cell lung cancer cohort and prognostic impact of concomitant mutations in KRAS and TP53 or STK11. *Lung Cancer*, 2019. 130: p. 50-58.
19. Tran L, Mattsson JS, Nodin B, Jönsson P, Planck M, Jirstrom K, et al. Various Antibody Clones of Napsin A, Thyroid Transcription Factor 1, and p40 and Comparisons With Cytokeratin 5 and p63 in Histopathologic Diagnostics of Non-Small Cell Lung Carcinoma. *Appl Immunohistochem Mol Morphol*, 2016. 24(9): p. 648-659.
20. Salomonsson A, Micke P, Mattsson JSM, La Fleur L, Isaksson J, Jönsson M, et al. Comprehensive analysis of RNA binding motif protein 3 (RBM3) in non-small cell lung cancer. *Cancer Med*, 2020. 9(15): p. 5609-5619.
21. Travis WD, B.E., Burke AP, Marx A, Nicholson AG (ed). *WHO Classification of Tumours of the Lung, Pleura, Thymus and Heart*. 4th ed. Lyon, France: IARC Press; 2015.
22. Sobin, L.H., Gospodarowicz, M.K. and Wittekind, C. (2009) *International Union against Cancer (UICC): TNM Classification of Malignant Tumours*. 7th Edition, Wiley-Blackwell, Chicester.
23. Djureinovic D, Hallström BM, Horie M, Mattsson JSM, La Fleur L, Fagerberg L, et al. Profiling cancer testis antigens in non-small-cell lung cancer. *JCI Insight*, 2016. 1(10): p. e86837.
24. R Core Team (2019). *R: A language and environment for statistical computing*. R Foundation for Statistical Computing, V., Austria. URL <https://www.R-project.org/>.
25. Zhang XZ, Chen MJ, Fan PM, Jiang W, Liang SX. BTG2 Serves as a Potential Prognostic Marker and Correlates with Immune Infiltration in Lung Adenocarcinoma. *Int J Gen Med*, 2022. 15: p. 2727-2745.
26. Jakobsen JN, Sørensen JB. Clinical impact of ki-67 labeling index in non-small cell lung cancer. *Lung Cancer*, 2013. 79(1): p. 1-7.
27. Zhu CQ, Shih W, Ling CH, Tsao MS. Immunohistochemical markers of prognosis in non-small cell lung cancer: a review and proposal for a multiphase approach to marker evaluation. *J Clin Pathol*, 2006. 59(8): p. 790-800.
28. Ringnér M, Jönsson G, Staaf J. Prognostic and Chemotherapy Predictive Value of Gene-Expression Phenotypes in Primary Lung Adenocarcinoma. *Clin Cancer Res*, 2016. 22(1): p. 218-29.
29. Kikuchi J, Kinoshita I, Shimizu Y, Kikuchi E, Takeda K, Aburatani H, et al. Minichromosome maintenance (MCM) protein 4 as a marker for proliferation and its clinical

- and clinicopathological significance in non-small cell lung cancer. *Lung Cancer*, 2011. 72(2): p. 229-37.
30. Liu Q, Yu Z, Xiang Y, Wu N, Wu L, Xu B, et al. Prognostic and predictive significance of thymidylate synthase protein expression in non-small cell lung cancer: a systematic review and meta-analysis. *Cancer Biomark*, 2015. 15(1): p. 65-78.
 31. Huang C, Lei C, Pan B, Fang S, Chen Y, Cao W, et al. Potential Prospective Biomarkers for Non-small Cell Lung Cancer: Mini-Chromosome Maintenance Proteins. *Front Genet*, 2021. 12: p. 587017.
 32. Lin CS, Liu TC, Lai JC, Yang SF, Tsao TC. Evaluating the Prognostic Value of ERCC1 and Thymidylate Synthase Expression and the Epidermal Growth Factor Receptor Mutation Status in Adenocarcinoma Non-Small-Cell Lung Cancer. *Int J Med Sci*, 2017. 14(13): p. 1410-1417.
 33. Acs B, Pelekanou V, Bai Y, Martinez-Morilla S, Toki M, Leung SCY, et al. Ki67 reproducibility using digital image analysis: an inter-platform and inter-operator study. *Lab Invest*, 2019. 99(1): p. 107-117.
 34. Grinberg M, Djureinovic D, Brunnström HR, Mattsson JS, Edlund K, Hengstler JG, et al. Reaching the limits of prognostication in non-small cell lung cancer: an optimized biomarker panel fails to outperform clinical parameters. *Mod Pathol*, 2017. 30(7): p. 964-977.

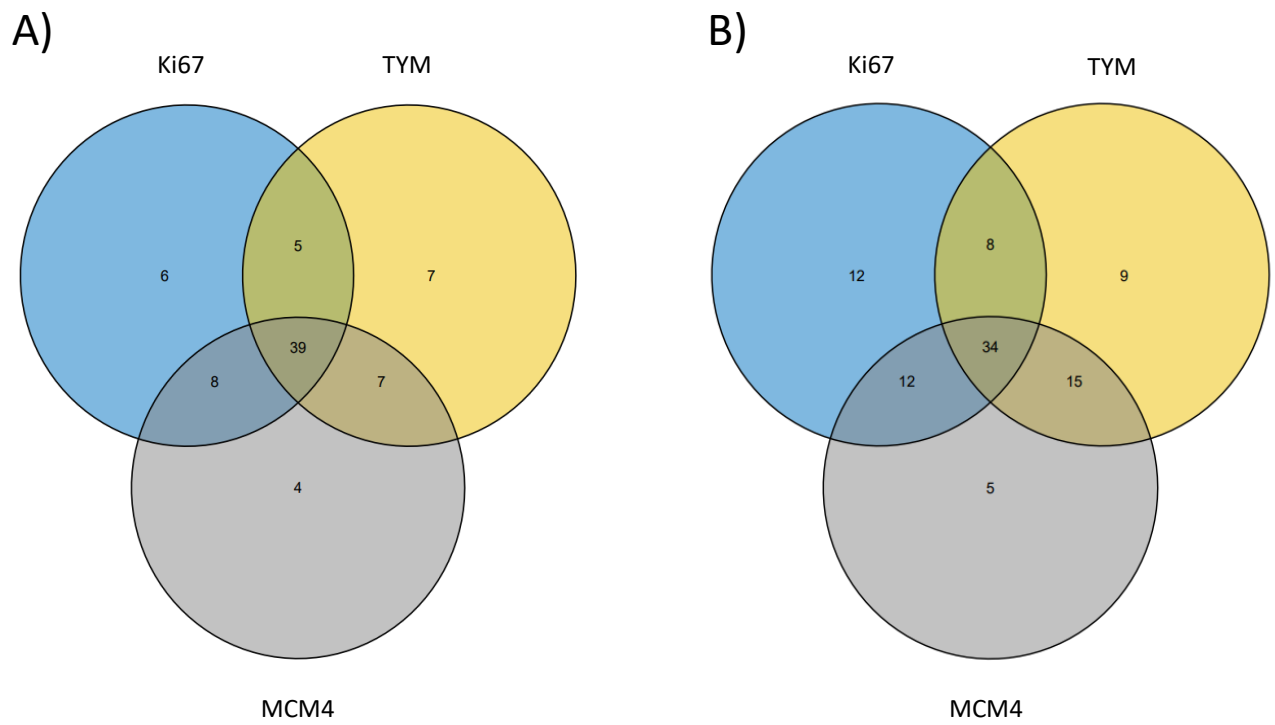
SUPPORTING INFORMATION

Supplementary Table 1. Immunohistochemical stainings for Ki67, MCM4 and TYMS.

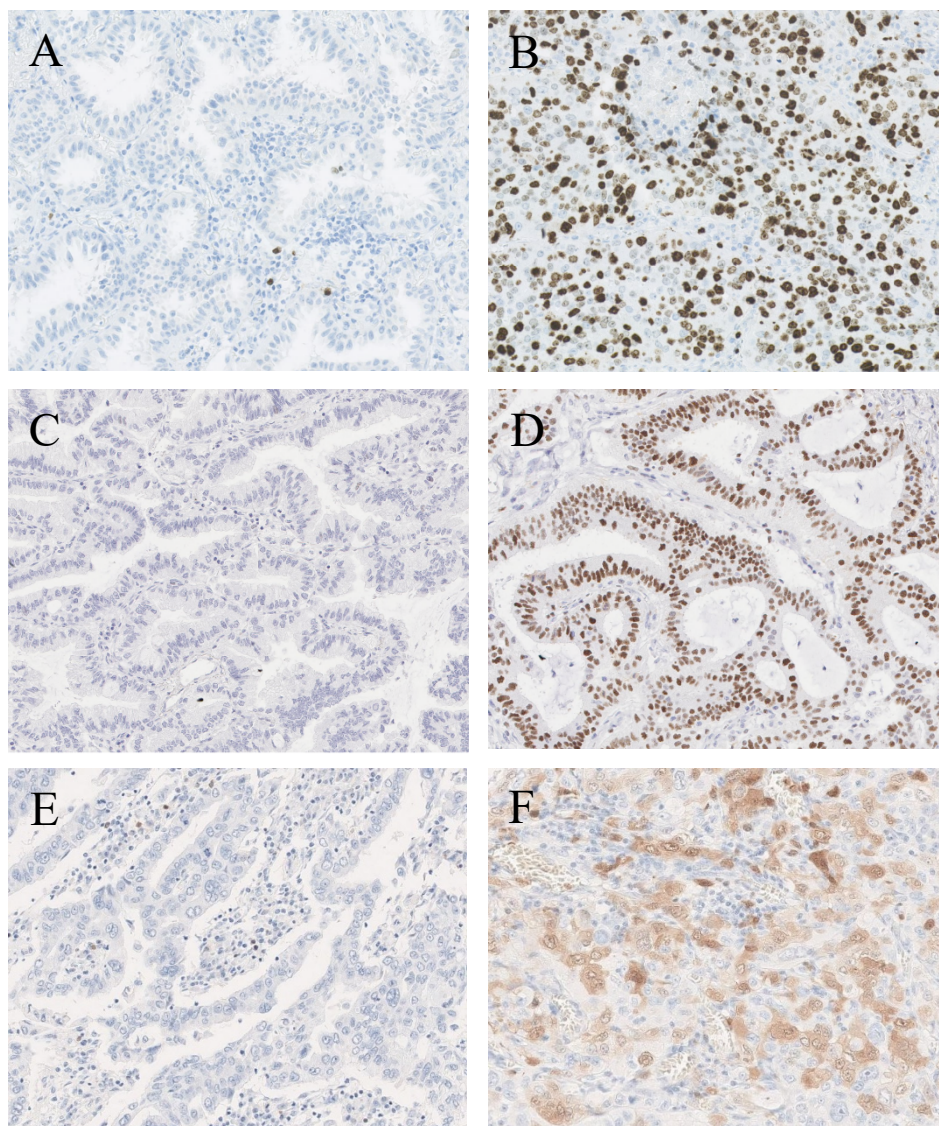
Antigen	Clone	Dilution	Pretreatment	Staining System
Ki67	MIB1	1:200	Dako high pH	Dako Autostainer+, EnVision
MCM4	D3H6N	1:200	Dako high pH	Dako Autostainer+, EnVision
TYMS	EPR4545	1:50	CC1	Ventana Discovery Ultra



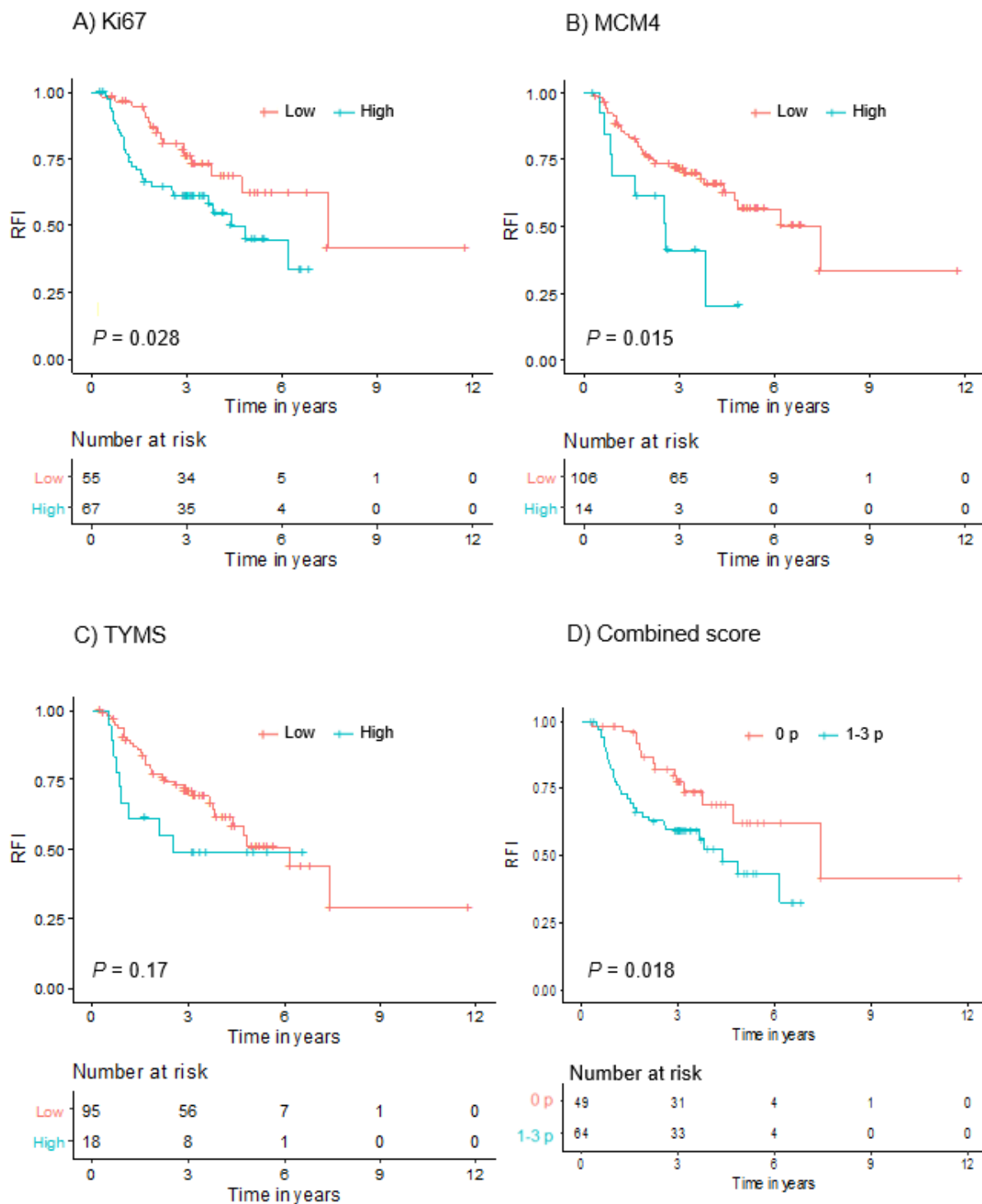
Supplementary Figure 1. Prognostic value of *Ki67* (A), *TYMS* (B), *GDF10* (C), and *CIRBP* (D) gene expression levels in one of the validation data set (Tang et al.¹⁴).



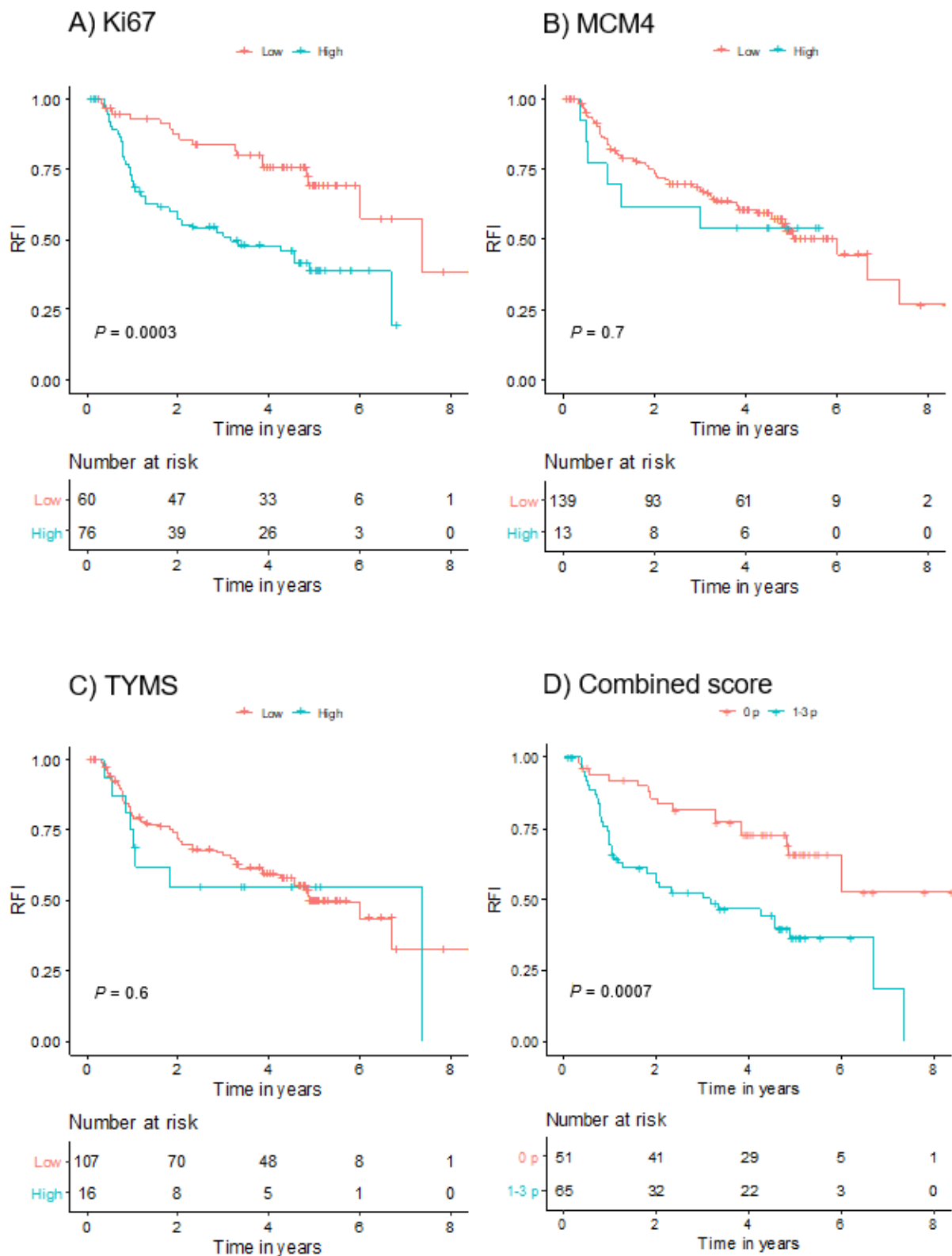
Supplementary Figure 2. The overlap between cases with high gene expression levels (cut-off based on the median gene expression values for each gene) of *Ki67*, *TYMS*, and *MCM4* in the two validation data sets. (A) Tomida et al.¹⁵, (B) Tang et al.¹⁴.



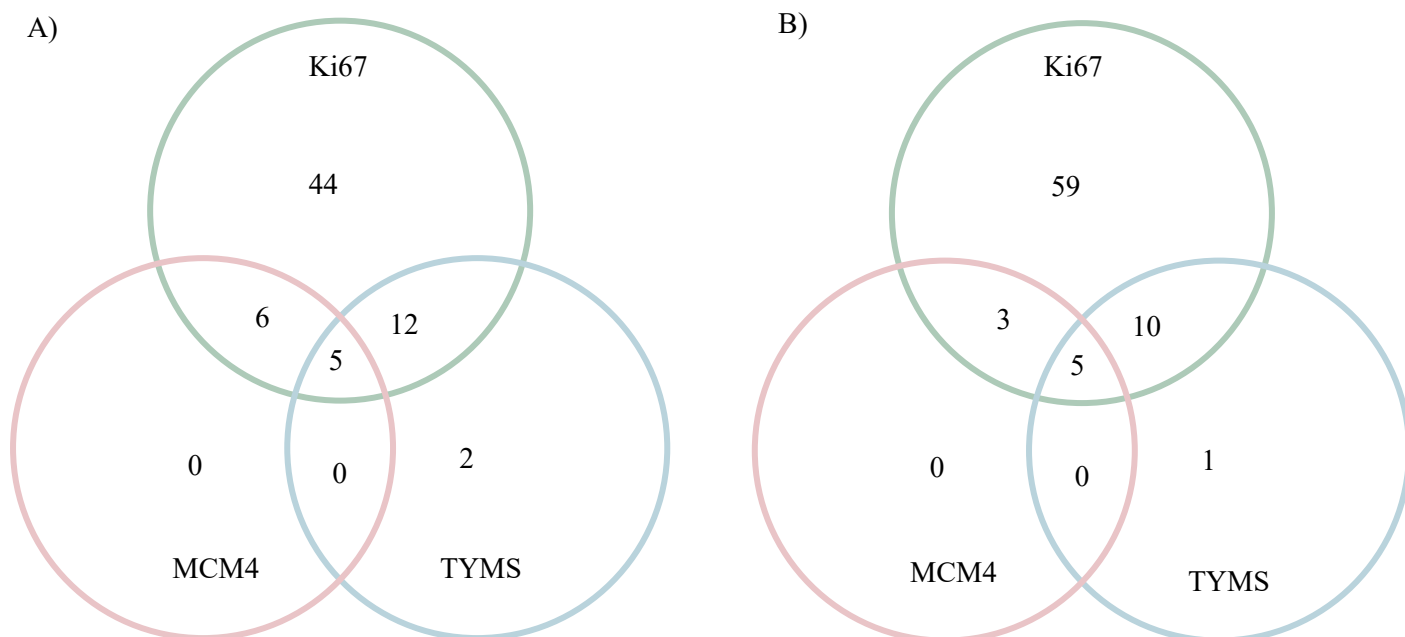
Supplementary Figure 3. Representative microscopic images of the stainings for Ki67 (A: low expression, B: high expression), MCM4 (C: low expression, D: high expression), and TYMS (E: low expression, F: high expression).



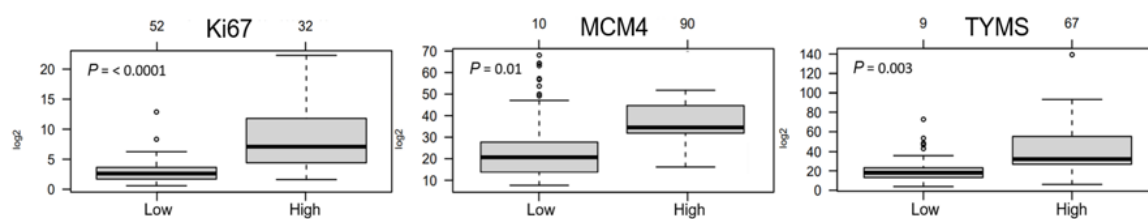
Supplementary Figure 4. Prognostic value of Ki67 (A), MCM4 (B), TYMS (C), and combined score (D), on recurrence-free interval (RFI) in the IHC discovery cohort.



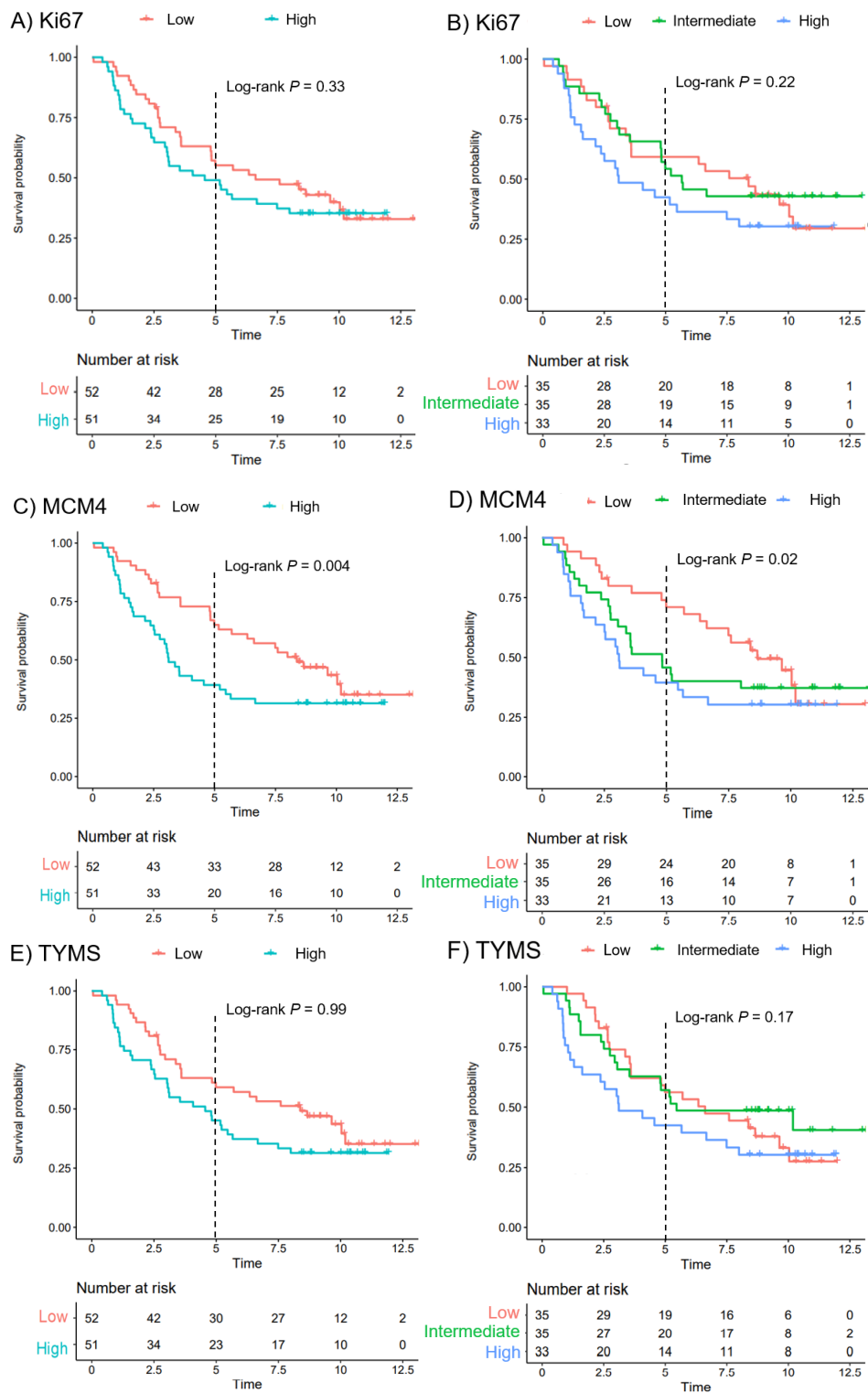
Supplementary Figure 5. Prognostic value of Ki67 (A), MCM4 (B), TYMS (C), and combined score (D), on recurrence-free interval (RFI) in the IHC validation cohort.



Supplementary Figure 6. The overlap between cases that were positive for the three markers in the IHC discovery cohort (A) and the IHC validation cohort (B).



Supplementary Figure 7. The association between gene expression levels and immunohistochemical classification (low or high expression) for the three markers in the IHC validation cohort.



Supplementary Figure 8. The prognostic value of Ki67 (A and B), MCM4 (C and D), and TYMS (E and F) gene expression levels in the IHC validation cohort, when dividing the samples into two or three groups.