

1 Article

# 2 A Novel Machine Learning 13-Gene Signature: Improving Risk 3 Analysis and Survival Prediction for Clear Cell Renal Cell 4 Carcinoma Patients

5 Patrick Terrematte<sup>1,2</sup>, Dhiego Souto Andrade<sup>1</sup>, Josivan Justino<sup>1,3</sup>, Beatriz Stransky<sup>1,4</sup>, Daniel Sabino A. de Araújo<sup>1</sup>  
6 and Adrião D. Dória Neto<sup>1,5</sup>

- 7
- 8     <sup>1</sup> Bioinformatics Multidisciplinary Environment (BioME), Metropole Digital Institute (IMD), Federal  
9 University of Rio Grande do Norte (UFRN), Natal, 59078-400, Brazil;  
10    <sup>2</sup> Department of Engineering and Technology (DETEC), Pau dos Ferros Multidisciplinary Center, Federal  
11 Rural University of Semi-arid (UFERSA), Pau dos Ferros, 59900-000, Brazil;  
12    <sup>3</sup> Department of Mathematics and Statistics (DME), Federal University of Rondônia (UNIR), Ji-Paraná, 76900-  
13 726, Brazil;  
14    <sup>4</sup> Biomedical Engineering Department, Center of Technology, UFRN, Natal, 59078-970, Brazil;  
15    <sup>5</sup> [Department of Computer Engineering and Automation, UFRN, Natal, 59078-970, Brazil](#);  
16 \* Correspondence: [patrick.terrematte@ufersa.edu.br](mailto:patrick.terrematte@ufersa.edu.br)

17 **Simple Summary:** Clear cell Renal cell carcinoma is a type of kidney cancer which comprises the  
18 majority of all renal cell carcinomas. Many efforts have been made to identify biomarkers which  
19 could help healthcare professionals better treat this kind of cancer. With extensive public data  
20 available, we conducted a machine learning study to determine a gene signature that could  
21 indicate patient survival with high accuracy. Through the min-Redundancy and Max-Relevance  
22 algorithm we generated a signature of 13 genes highly correlated to patient outcomes. These  
23 findings reveal potential strategies for personalized medicine in the clinical practice.

24  
25 **Citation:** Terrematte, P.; Andrade, D.  
26 S.; Justino, J.; Stransky, B.; Araújo,  
27 D.; Dória Neto, A. Title A Novel  
28 Machine Learning 13-Gene  
29 Signature: Improving Risk Analysis  
30 and Survival Prediction for Clear  
31 Cell Renal Cell Carcinoma Patients.  
32 *Cancers* **2022**, *14*, x.  
33 <https://doi.org/10.3390/xxxxx>

34 Academic Editor: Firstname Last-  
35 name

36 Received: date

37 Accepted: date

38 Published: date

39 **Publisher's Note:** MDPI stays  
40 neutral with regard to jurisdictional  
41 claims in published maps and  
42 institutional affiliations.



43 **Copyright:** © 2022 by the authors.  
44 Submitted for possible open access  
45 publication under the terms and  
46 conditions of the Creative Commons

2 Attribution (CC BY) license (<https://creativecommons.org/licenses/by/>)  
3 Cancers 2022, 14, x. <https://doi.org/10.3390/xxxxx>

4.0).

47 | **Keywords:** Kidney cancer; clear cell Renal Cell Carcinoma (ccRCC); Gene signature;  
48 | Prognosis; [Survival analysis](#); Feature selection; Mutual Information; Machine Learning  
49 |

50 |

## 1. Introduction

51 | Renal cell carcinoma (RCC) occurs in the renal cortex or the renal tubular  
52 | epithelial cell. The molecular subtypes of renal cancers are clear cell RCC (ccRCC),  
53 | papillary RCC (pRCC), and chromophobe RCC (ChRCC). RCC accounts for more  
54 | than 90% of cancers in the kidney [1], of which 80–90% are ccRCC [2], and more  
55 | than 30% of patients with ccRCC experience metastasis [3]. In 2020, the worldwide  
56 | mortality rate from kidney cancer was an estimated 179,368 cases International  
57 | Agency for Research on Cancer (IARC). The American Cancer Society estimated a  
58 | prevalence of 76,080 new cases of kidney cancer for 2021 in the United States (48,780  
59 | in men and 27,300 in women), and an estimated mortality rate of 13,780 people  
60 | (8,790 men and 4,990 women) [4]. Depending on the stage at diagnosis, the five-year  
61 | survival rates of RCC in the US are the following: 93% for localized disease (stage I),  
62 | 72.5% for regional disease (stage II/III, local lymph node involvement), and only 12%  
63 | for late-stage (stage IV metastatic) [5]. The poor survival outcomes of metastatic  
64 | patients with ccRCC reveal the importance of seeking new and robust biomarkers of  
65 | prognosis, and of preventing the progression of non-metastatic tumors.

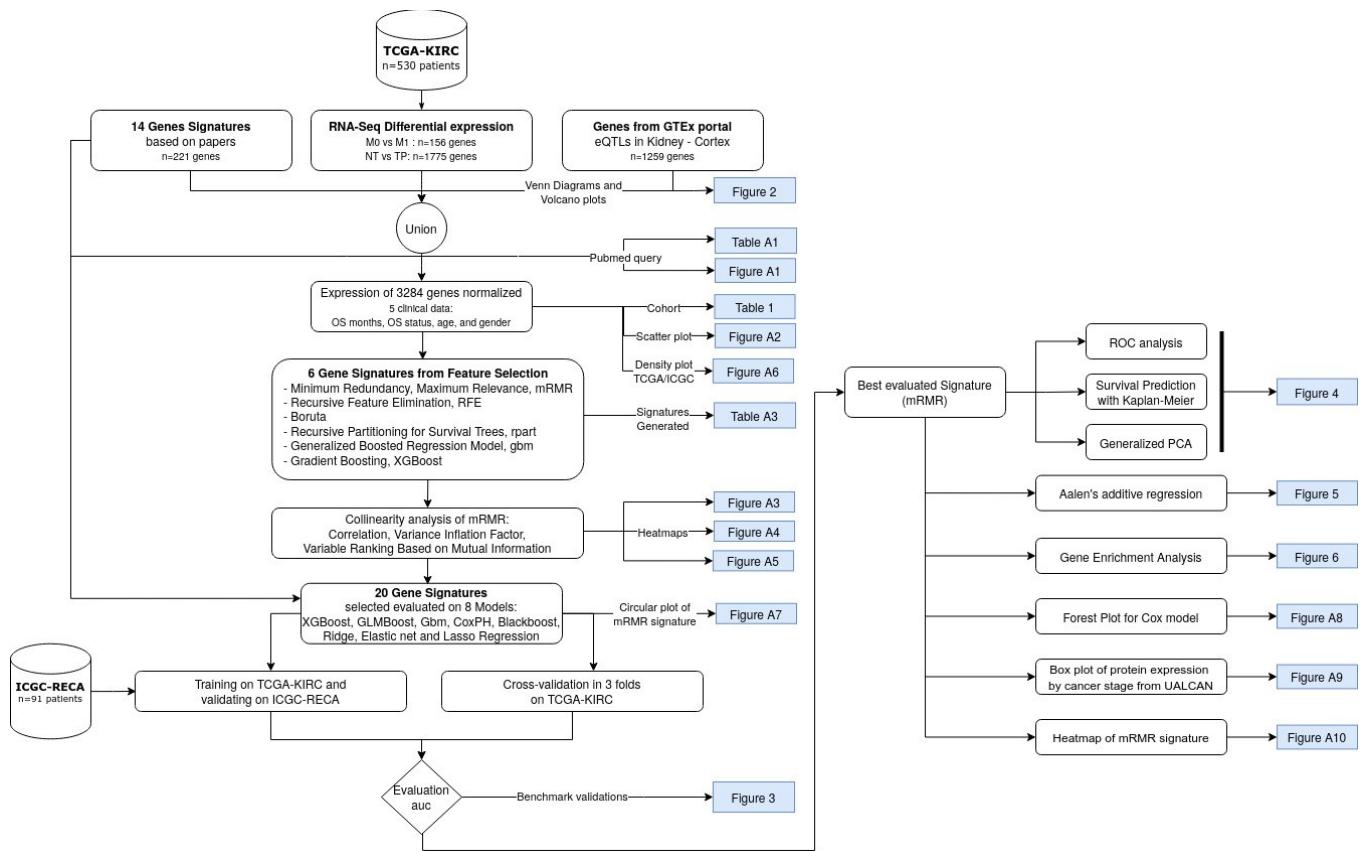
66 | The challenges of artificial intelligence (AI) applications to cancer care are  
67 | driven by the translation of models with clinical validity, utility, and usability into  
68 | feasible clinical treatment [6]. In the field of precision medicine applied to cancer,  
69 | feature selection is useful in detecting the most important traits and molecular  
70 | profiles for predicting the survival risks of a patient's outcome through a given gene  
71 | set. A gene signature is a set of genes whose expression pattern in a specific cell  
72 | type and condition can provide a biomarker for diagnosis, prognosis, or therapeutic  
73 | responses in cancer patients [7]. The gene signatures can be defined by the pattern of  
74 | the Single Nucleotide Variant (SNV) mutational profile; the copy number of  
75 | alterations (CNA); the methylation levels; or the expression of messenger or other  
76 | RNA types. Genes involved in the biological processes of many tumors might be  
77 | overexpressed or inhibited, signaling a better or worse prognosis for the patient [8].  
78 | While most of the studies used only mRNA data to build their signatures,  
79 | microRNA and/or clinical data can be explored as relevant features to build a  
80 | predictive signature [9–14].

81 | [Nowadays, the scientific community is still searching for new biomarkers for](#)  
82 | [ccRCC, and feature selection methods using survival analysis provide a robust](#)  
83 | [exploratory methodology before experimental validations. Survival analysis is a](#)  
84 | [field of statistics that predicts the time until an event of interest happens in many](#)  
85 | [domains \[15\]. The most commonly used method for survival analysis is the Cox](#)  
86 | [Regression model \[16\]. The Cox model is semi-parametric, that is, the distribution of](#)  
87 | [the event of interest is unknown. In addition, Cox models are widely used for](#)  
88 | [censored data, i.e., when the event is not observed during the study period due to](#)  
89 | [loss to follow-up, study termination, or the patient's death by other causes.](#)  
90 | [Regularized Cox models provide suitable predictions for high-dimensional data](#)  
91 | [using penalty functions with the main regularizers Lasso-Cox, Ridge-Cox, and](#)  
92 | [Elastic net-Cox \[15\]. Ensemble learning methods are committees of machine learning](#)  
93 | [models, in other words, they combine the majority of the votes for each model in an](#)  
94 | [ensemble or they adjust the weighted vote of each model. Moreover, this approach](#)  
95 | [results in a more robust, efficient, and stable model compared to singular models. In](#)

96  
97  
98  
99  
100  
101  
102  
103  
104  
105  
106  
107  
108  
109  
110  
this work, we applied Cox models and ensemble methods using gene expression to predict the overall survival (OS) after diagnosis of ccRCC.

111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
Lasso-Cox regression generated most of the reviewed gene signatures for ccRCC [9,10,13,17–19]. All the studies reviewed in this work use TCGA-KIRC dataset to train and validate the results. Fewer studies validated their results with other datasets such as GEO database [10,13,2], ICGC-RECA [2,11], and data from Fudan University Shanghai Cancer Center (FUSCC). The most common methodologies used to discover and validate gene signatures were differentially expression analysis (DEA), and gene set enrichment analysis (GSEA). Only one study compared its methodology to 3 other biomarker signatures from our literature selection [9]. In addition, there was a lack of comparisons between the gene signatures. As far we know, our study presents a most comprehensive comparison between gene signatures, including ensemble methods, machine learning, and feature selection.

This study aims to specify a gene signature based on the state-of-the-art algorithms of feature selection methods, and to be able to predict the survival risk of ccRCC patients. Moreover, this study compares the novel signatures obtained by these feature selection methods, and other previously published gene signatures. The best-performing gene signature was achieved using the mutual-information-based ensemble method of min-Redundancy and Max-Relevance (mRMR) [16]. Specifically, the mRMR is an ensemble-based method to select a minimal set of features with a maximum prediction performance. The flowchart shown in Figure 1 displays a summarized view of the discovery process for the novel mRMR gene signature of ccRCC.



122  
123  
124  
125  
126

**Figure 1.** Flowchart of the current study to obtain a gene signature based on mutual information, Minimum Redundancy Maximum Relevance (mRMR). The datasets are indicated by the cylinder, white rectangles represent a step of the analysis, and the blue rectangles indicate the resulting figures and tables. TCGA-KIRC and ICGC-RECA are datasets of ccRCC.

## 127 2. Materials and Methods

### 128 2.1. Literature Search Using PubMed

129 A literature search for gene signatures was conducted using PubMed to select  
130 studies from 2015 to 2020 (Figure A1);—given that the search was carried out in  
131 January 2021, from this date, we performed the analyses and wrote the manuscript.  
132 The majority of papers were published in the last five years since 2020, therefore we  
133 excluded the period of 2008 to 2014. The PubMed query of terms comprised the  
134 following: (renal OR kidney) AND (clear cell) AND (cancer) AND (prognosis OR  
135 survival OR outcomes) AND (regression) AND (gene signature). The search was  
136 conducted in January 2021.

137 The search query resulted in 770 papers, and we adopted the following as  
138 inclusion criteria: original articles on human ccRCC about survival prognosis or  
139 tumor staging classification. The exclusion criteria consisted of the following:  
140 reviews, editorials, conferences, or abstracts; studies about other RCC subtypes,  
141 such as pRCC, ChRCC, or Sarcomatoid renal cell carcinoma; and studies that  
142 evaluated genes based on their corresponding patient prognoses depending on  
143 chosen treatment, on biomarkers predicting treatment resistance, or on tolerance to  
144 renal allograft. Ultimately, we adopted 14 gene signatures with a total of 221 unique  
145 genes (Table A1).

### 146 2.2. Data

147 From a bottom-up perspective, this work is data-driven by the gene expression  
148 and survival data of the larger public dataset of ccRCC (n=530), The Cancer Genome  
149 Atlas Consortium of Kidney Renal Clear Cell Carcinoma (TCGA-KIRC) [17,18]  
150 [16,17]. For external data validation, in order to corroborate the findings within our  
151 novel gene signature, we used the dataset of ccRCC samples (n=91) from the  
152 International Cancer Genome Consortium (ICGC-RECA) [19,20][18,19].

### 153 2.3. Data Pre-Processing

154 Data pre-processing was undertaken to select the genes from both TCGA-KIRC  
155 (n=60,489) and ICGC-RECA (n=49,221) cohorts to obtain a consensus nomenclature  
156 for the genes in the signatures, and to map the latter with the HGNC Symbol and the  
157 Ensembl identifiers. The reference genomes used for the TCGA-KIRC and ICGC-  
158 RECA databases are the GRCh38 and the GRCh37 genomes, respectively. Despite  
159 this distinction, both reference genome versions are highly concordant [21][20].

160 For both datasets, we used unprocessed raw count data, and to reduce the  
161 batch-effect of datasets, we evaluated the following normalization methods: (1)  
162 scaling to the range interval between zero and one; (2) Variance Stabilizing  
163 Transformation with DESeq2 [22][21]; and (3) Box-Cox normalization with Caret R  
164 package [23][22] (v. 6.0-90). The chosen method was Box-Cox transformation, with  
165 the higher correlation  $R = 0.97$  between the median of each gene expression of  
166 datasets (Figure A42).

### 167 2.4. Feature Selection with Bioinformatics Analyses and Machine Learning

168 From a top-down perspective, in order to guide our feature selection, we  
169 performed two differential expression analyses of RNA-Seq with DESeq2 [22][21].

The first analysis was to compare solid normal tissue (NT) samples ( $n=71$ ) versus primary solid tumor (TP) samples ( $n=530$ ) using the absolute log<sub>2</sub> fold-change (LFC)  $> 3$  and p-value adjusted (FDR)  $< 0.01$ , from which we obtained 1,775 genes under- and over-expressed. The LFC of each gene expression is the ratio of the mean normalized by log<sub>2</sub> in the two groups of samples. The second analysis was to compare the non-metastatic (M0) samples ( $n=422$ ) against metastatic (M1) samples ( $n=78$ ); then using absolute LFC over 2 and p-value adjusted (FDR)  $< 0.01$ , we obtained 156 altered genes.

To optimize the right candidates to their ideal biomarker genes, we also included 221 genes from the literature, and selected 1259 tissue-specific genes of Kidney Cortex tissue with significant expression Quantitative trait locus (eQTL) obtained from the Genotype-Tissue Expression (GTEx) Project [24,25][23,24]. The feature selection methods and supervised Cox regression models were then trained by gene expressions of 33043,284 pre-selected genes, the overall survival (OS) days since the diagnosis, and the OS status (deceased or living) of TCGA-KIRC patients.

Inspired by the methodology of [26][25], we produced the new gene signatures using 6 feature selection methods divided into two main categories:

1. Filtering methods of feature importance: eXtreme Gradient Boosting (XGBoost), Generalized Boosted Regression Model (GBM), and Recursive Partitioning for Survival Trees (Rpart).
  2. Wrapper methods: Minimum Redundancy Maximum Relevance (mRMR); Recursive Feature Elimination (RFE); and Boruta.

For the filtering methods, we selected the 30 most important genes for patient survival. We chose the number of 30 genes based on this being the average number of genes in signatures referenced in the literature. The wrapper methods selected the most important genes based on the best performing metrics of models without predefining the number of genes on signatures. The new gene signatures generated by each feature selection are available in Table A2.

We evaluated the signatures using Machine Learning analyses of eight linear survival models with optimized auto-tuning hyper-parameters: Extreme Gradient Boosting (XGBoost), Cox Model with gradient boosting (GLMBoost), Generalized Boosted Regression Model (GBM), Cox proportional hazards Regression model (CoxPH), Gradient Boosting with Regression Trees (Blackboost), and three models of Penalised Cox Regression (glmnet) [27][26] – LASSO, ElasticNet and Ridge regression. Each model calculates the fitted coefficient for each gene.

The mRMR method applies mutual information to select features that maximize the statistical dependency on the joint distribution of the target variable of supervised learning [16,28][15,27]. The maximum relevance for the feature set  $S$ , given the mutual information of gene  $g_i$  in  $k$ -classes, is:

$$maxD(S, k), D = \frac{1}{|S|} \sum_{g_i \in S} I(g_i, k) \quad (1)$$

The minimum redundancy in the feature subset condition is given by the sample vectors of all genes  $g_i, g_j$ :

$$minR(S, k), D = \frac{1}{|\Sigma^2|} \sum_{g_i, g_j \in S} I(g_i, k) \quad (2)$$

This work uses the implementation of the R package mRMRe [29][28] (v. 2.1.2) available in CRAN on expression Data. The target features consisted of the Overall Survival days and Overall Survival status. We set an ensemble of 5 executions filtering 20 genes per run, resulting in a set of 64 unique genes as relevant features. Finally, we performed a forward search feature selection with Variable Ranking

216 | Based on Mutual Information Difference of the most representative genes with  
217 | respect to AJCC Staging, resulting in a 13-gene signature (Figure A34).  
218 |  
219 |  
220 |  
221 |  
222 |  
223 |  
224 |  
225 |  
226 |

The framework of Tidyverse in R (v. 4.1.1) was used for pre-processing, and the framework mlr3 (Machine Learning in R) [30][29] carried out the evaluation of the metrics of feature selection and model benchmark. All of the code for the experiments was written in R. For the multicollinearity analysis, we built the visualization with corrrplot [31][30] (v. 0.92), and we assessed the degree of collinearity among independent variables. None of the genes had Variance Inflation Factors > 5 (Figure A24). Also, no correlations greater than or equal to 0.7 were found between the genes (Figure A35). For the Variable Ranking Based on Mutual Information Difference, we used the R package varrank [32][34] (v. 0.4).

227 | *2.5. Model Evaluation and Statistical Analysis*

228 | The concordance C-index is a commonly used metric, but is not a proper  
229 | strategy to predict the t-year risk of an event [33][32]. Therefore, to evaluate the  
230 | performance of each survival model, we applied the measure of the area under the  
231 | time-dependent ROC curve (AUC Uno) [34][33]. For internal validation, we used  
232 | AUC Uno of 10-years on 3-fold cross-validation of TCGA-KIRC in 100 repetitions.  
233 | For external validation, we used AUC Uno of 7-years by training with TCGA-KIRC  
234 | and predicting the ICGC-RECA dataset using 100 repetitions through censored  
235 | regression models. We restrict the 10-years prediction for TCGA-KIRC to exclude  
236 | outliers in the long tail of the density plot of the patient's overall survival. For the  
237 | ICGC-RECA dataset, we decided to maintain a 7-years prediction in order to include  
238 | all samples, and limit the time prediction to the range of distribution of this dataset  
239 | for external validation (Figure A6). The sensitivity (SE) and the specificity (SP)  
240 | describe the distinguishing risk of patients to be deceased by time  $t$  from those who  
241 | will be alive, with values ranging from 0 to 1, where 1 corresponds to the best model  
242 | performance, and 0.5 represents a random prediction. The evaluation was  
243 | performed with the R package survAUC [35][34] (v. 1.0-5).

244 | The Kaplan-Meier analysis is the main visualization graph used to distinguish  
245 | between high-risk, moderate, and low-risk patients. The p-value was calculated by  
246 | the log-rank test using the survminer [36][35] (v. 0.4.9) R package and by comparing  
247 | the predicted survival distributions of groups' high, moderate, and low risk.

248 | The enrichment analysis was performed using the 13-gene signature on the  
249 | curated database of DisGeNET [37][36] (v7.0) with gene-disease associations (GDAs)  
250 | filtering by FDR (<0.05).

251 | The Flowchart was created using diagrams.net. The figures were implemented  
252 | in R 4.1.1 using the following packages: VennDiagram [38][37] (v. 1.7.1); the ggplot2  
253 | (v. 3.3.5) for Volcano plots, Heatmap and Boxplots; GOplopt [39][38] (v. 1.0.2) for the  
254 | circular visualization of mRMR genes and sets of genes; FactoMineR [40][39] (v. 2.4)  
255 | and factoextra [41][40] (v. 1.0.7) for the Principal Component Analysis (PCA);  
256 | survival [42][41] (v. 3.2-11) and ggstatsplot [43][42] (v. 0.9.0) for the Aalen's additive  
257 | cox regression; clusterProfiler [44][43] (v. 4.2.1) and disgenet2r [37][36] (v. 0.99) for  
258 | the enrichment analysis with a Heatmap-like functional classification; survminer  
259 | [36][35] (v. 0.4.9) and finalfit [45][44] (v. 1.0.4) for the survival curves and the Forest  
260 | plot for Cox proportional hazards model; and pheatmap [46][45] (v. 1.0.12) for the  
261 | Heatmap with Hierarchical clustering of RNA-seq expression and clinical  
262 | annotation with dendrograms (Figure A6).

263 | **3. Results**

264 | *3.1. Clinical Characteristics of the ccRCC Cohorts*

In order to produce our gene signature, we used the TCGA-KIRC ( $n = 530$ ) and ICGC-RECA ( $n = 91$ ) samples of RNASeq data of ccRCC. The characteristics of both cohorts for training and validation datasets are summarized in Table 1. The clinical characteristics with their respective p-value tests indicate that there is no significant distinction in the distributions between both datasets, except for Neoplasm.

**Table 1.** Study Characteristics of TCGA-KIRC and ICGC-RECA cohort [with the clinical characteristics](#):—Age, gender, tumor grade, metastasis, and staging by the American Joint Committee on Cancer (AJCC).

Clinical characteristics		Training cohort TCGA-KIRC ( $n=530$ ) <sup>1</sup>	Validation cohort ICGC-RECA ( $n=91$ )	p value <sup>2</sup>
Overall survival (days)	Mean (SD)	1343.2 (976.6)	1511.6 (634.6)	0.113
Overall survival status, N./total N. (%)	Alive Deceased	359/530 (67.7) 171/530 (32.3)	61/91 (67.0) 30/91 (33.0)	0.991
Age, years	Mean (SD)	60.5 (12.0)	60.5 (10.0)	0.99
Gender, N./total N. (%)	Female Male	183/530 (34.5) 347/530 (65.5)	39/91 (42.9) 52/91 (57.1)	0.158
AJCC stage, N. / Total (%)	T1 T2 T3 T4	270/530 (50.9) 70/530 (13.2) 179/530 (33.8) 11/530 (2.1)	54/91 (59.3) 13/91 (14.3) 22/91 (24.2) 2/91 (2.2)	0.343
Neoplasm, N. (%)	N0 N1 NX	79 (86.8) 2 (2.2) 10 (11.0)	239 (45.1) 16 (3.0) 275 (51.9)	<0.001
Metastasis, N. (%)	M0 M1 MX	422/528 (79.9) 78/528 (14.8) 28/528 (5.3)	81/91 (89.0) 9/91 (9.9) 1/91 (1.1)	0.081

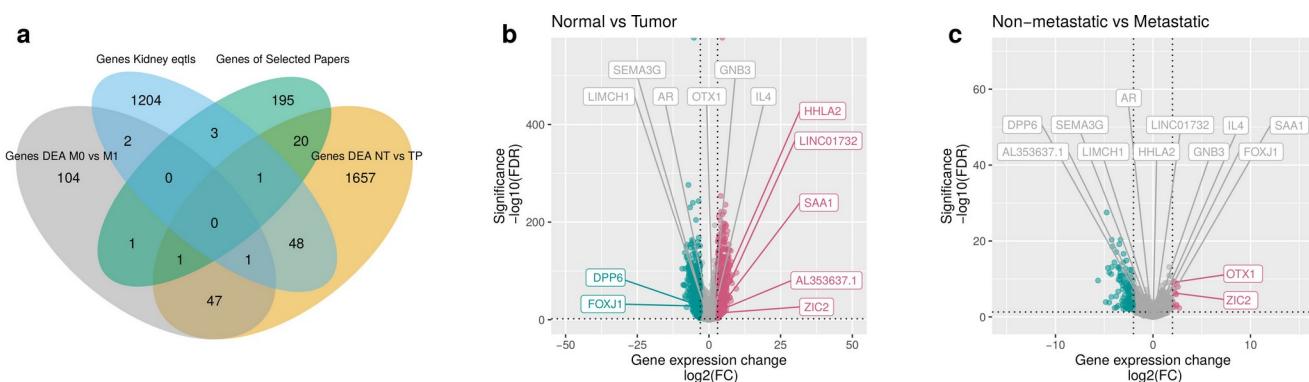
<sup>1</sup> The metastasis values do not sum up to heading totals because of missing data.

<sup>2</sup> The statistical tests for age and overall survival days are performed by Wilcoxon rank-sum test, and all other comparisons are by Fisher exact test.

### 3.2. mMRM Gene Selection

The mRMR executed a supervised gene selection of 3304 genes with ~~three~~four clinical features: overall survival (OS) days, OS status, age and sex. [In order to identify the most representative genes of the signature related to Stage AJCC, we performed a forward search feature selection Variable Ranking Based on Mutual Information Difference, resulting in a 13-gene signature \(AR, AL353637.1, DPP6, FOXJ1, GNB3, HHLA2, IL4, LIMCH1, LINC01732, OTX1, SAA1, SEMA3G, ZIC2 – Figure A3\) able to predict distinct outcomes \(high, moderate and low survival risk\) of patients with ccRCC.](#) In order to select the best independent predictors genes for survival risk, it is important to avoid multicollinearity, therefore we assessed the degree of collinearity amongst independent variables. None of the genes had Variance Inflation Factors >5 (Figure A42). Also, no correlations greater than 0.70 were found between the genes (Figure A53). [In order to identify the most representative genes of the signature related to Stage AJCC, we performed a forward search feature selection Variable Ranking Based on Mutual Information Difference, resulting in a 13-gene signature \(AR, AL353637.1, DPP6, FOXJ1, GNB3, HHLA2, IL4, LIMCH1, LINC01732, OTX1, SAA1, SEMA3G, ZIC2 – Figure A4\) able to predict distinct outcomes \(high, moderate and low survival risk\) of patients with ccRCC.](#)

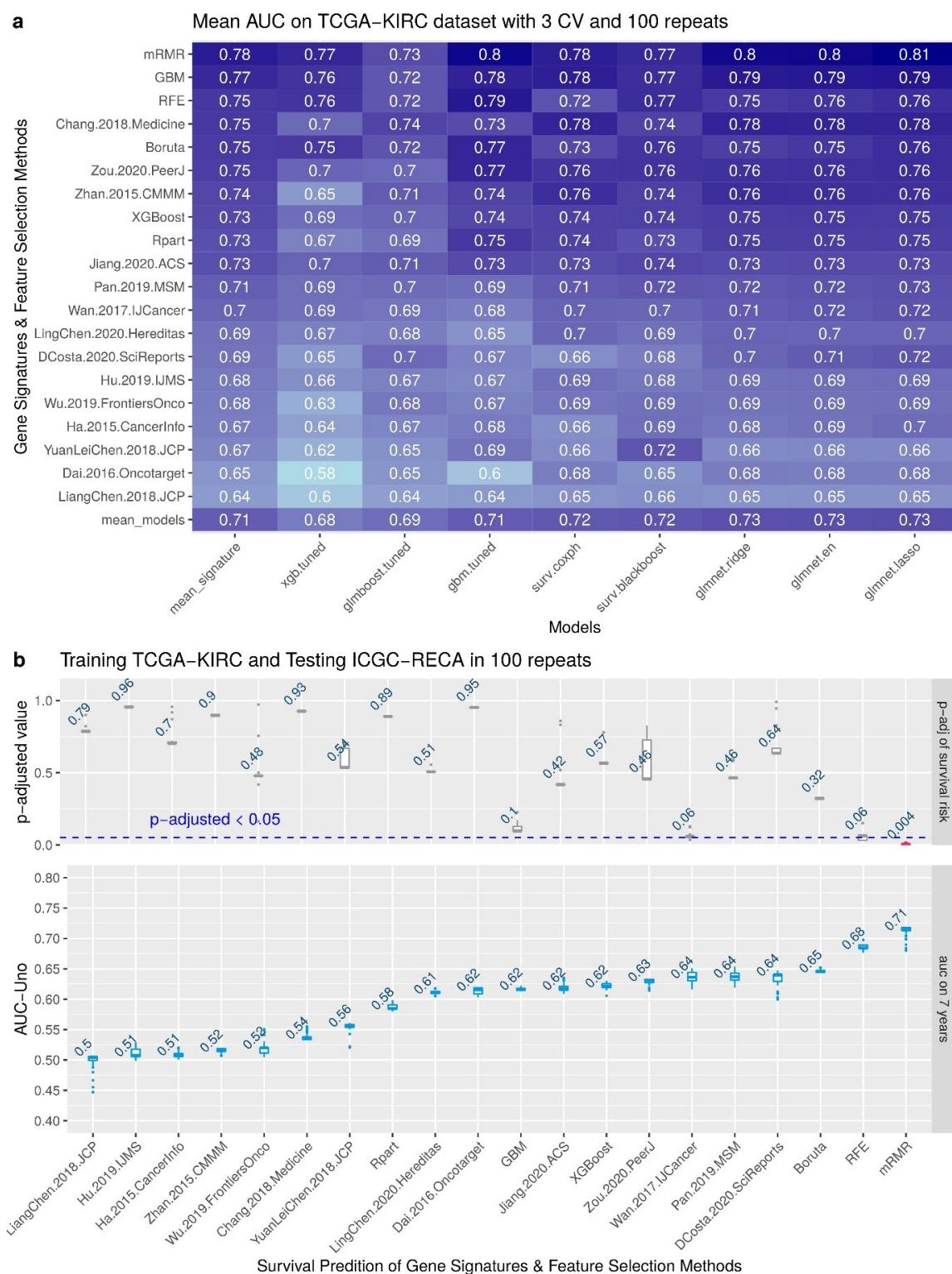
We visualized the composition of filtered genes with a Venn diagram (Figure 2a) with the intersection sizes of genes and the original sets of genes. In particular, most of the mRMR genes ( $n=7$ ) were obtained from the differential gene expression analysis (DEA) comparing normal tissues versus primary tumor samples (Figure 2ab), with a larger number of upregulated genes, including the mRMR genes HHLA2, LINC01732, SAA1, AL353637.1, and ZIC2. The downregulated mRMR genes for normal versus tumor samples are DPP6 and FOXJ1. The DEA of comparing non-metastatic versus metastatic samples (Figure 2c) has identified less differentiated genes ( $n=2$ ), with the upregulated genes OTX1 and ZIC2. The genes selected with mRMR on TCGA-KIRC samples are presented in Figure A72d with a circular visualization of the relationship between genes and their original sets of DEA, genes from GTEx portal of expression quantitative trait loci (eQTLs) in Kidney Cortex, and gene signatures from the literature (Table A1).



**Figure 2. Selected genes through mRMR.** (a) prefiltered gene sets and the mRMR gene signature. A total of 3284 prefiltered genes is given by the sets of DEA between non-metastatic versus metastatic (156), normal tissues versus primary tumor (1775), genes from literature (221), significant eQTLs genes (1259), and 124 genes overlapping in two or three intersections of sets. (b) Volcano Plot of Differential Expression Analysis EA comparing normal tissues versus primary tumor samples of TCGA-KIRC. In green, we see the downregulated genes of normal tissues versus primary tumors (DPP6 and FOXJ1). In red, we see the upregulated genes (HHLA2, LINC01732, SAA1, AL353637.1, and ZIC2). In gray, we see the not significant genes with low Fold-Change. (c) Volcano Plot of Differential Expression Analysis DEA comparing non-metastatic versus metastatic samples. In red, we see the upregulated genes (OTX1 and ZIC2). (d) Circular diagram of mRMR gene signature and the source of genes.

### 3.3. Performance of the Feature Selection Models for Internal and External Validations

In order to compare our mRMR signature with six feature selection methods (Recursive Feature Elimination, Boruta, Rpart, GBM and XGBoost for Survival) and 14 signatures published, we performed a benchmark using eight survival models of cox survival regressions (XGBoost, GLMBoost, Gbm, CoxPH, Blackboost, Ridge, Elastic Net, and Lasso). The benchmark results are shown in Figure 3a with the performance of 100 repetitions of predictions with Area Under the Curve (AUC) Receiving Operator Characteristics (ROC) Uno evaluating the 20 gene signatures using 3-fold cross-validation of TCGA-KIRC dataset. We can observe that the model Lasso Cox regression of glmnet had the best mean of AUC, 0.81, in internal validation for mRMR. The minimal set of genes with best performance to predict TCGA-KIRC as internal validation is: AL353637.1, DPP6, FOXJ1, GNB3, HHLA2, IL4, LIMCH1, OTX1, SAA1, and ZIC2.



336

337  
338  
339  
340341  
342  
343

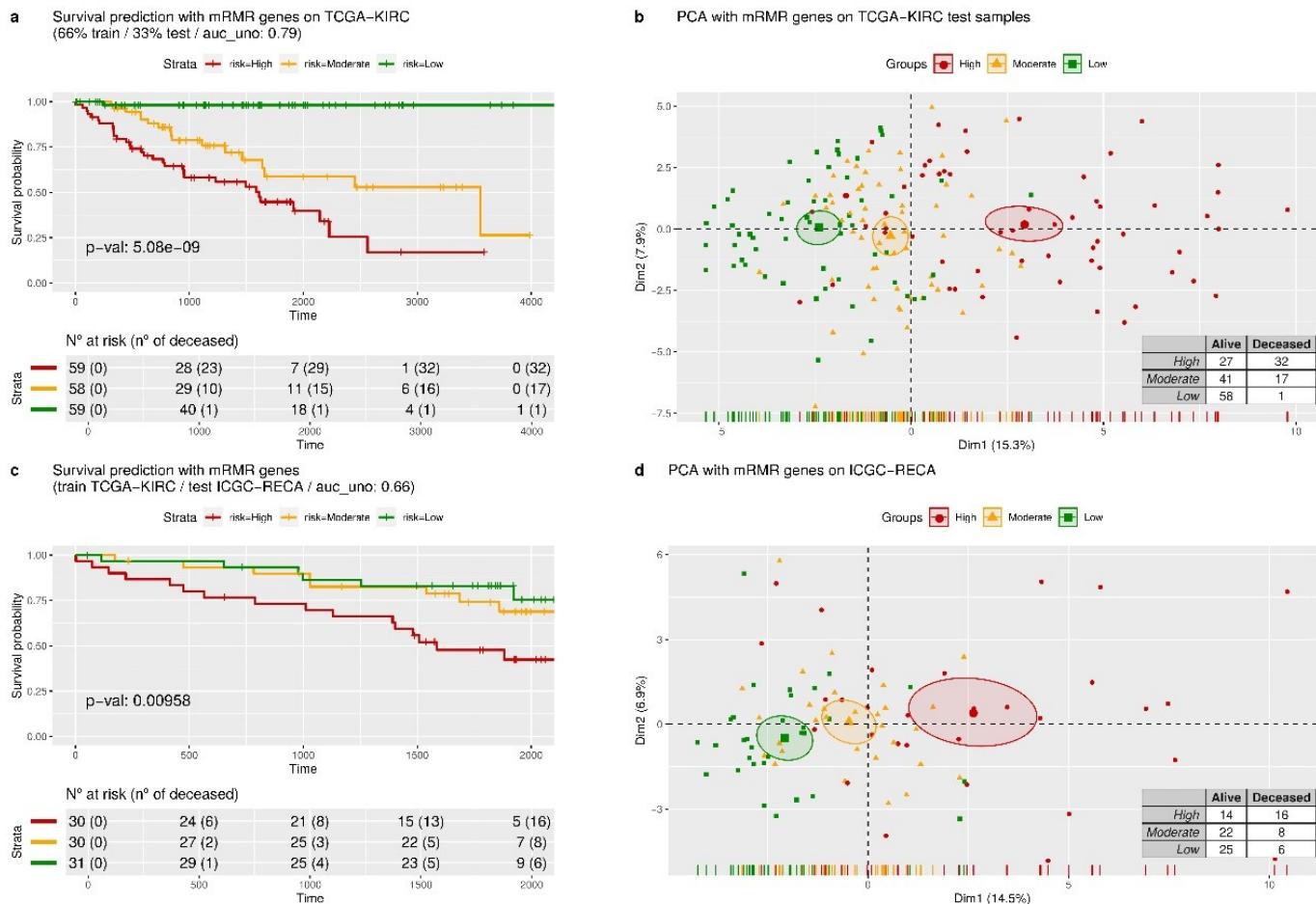
**Figure 3. Benchmark with internal and external validation.** (a) Comparing 14 gene signatures from literature and 6 feature selection on 8 models for survival risk, showing the predicted AUC of survival outcome in 10-years [prediction](#). (b) Box-plots of results of each gene signature and feature selection [in 7-years prediction](#).

The Figure 3b shows the box-plot of this result of external validation in 100 random repeats. The upper plot also displays the mean of the adjusted p-value of the log-rank test of survival risk. Note that the unique signature that has a

344  
345  
346  
347  
348  
349  
350  
351

significant adjusted p-value ( $p < 0.05$ ) is the mRMR. The lower plot displays the AUC metric of each survival prediction, and the number displayed on box-plots is the average value of all repeats. Note that the best mean of AUC is 0.71 for mRMR signature. The minimal set of genes for training with samples of TCGA-KIRC and predicting the survival risk of samples of ICGC-RECA is: AR, AL353637.1, FOXJ1, HHLA2, SEMA3G, and LINC01732.

In Figure 4, we display the Kaplan-Meier curves and a principal component analysis (PCA) of two random predictions of internal and external validations.



**Figure 4. Survival risk predictions with mRMR signature and dimensionality reduction.** (a) The survival curves are predicted in three equal-size strata of risk groups of the TCGA-KIRC dataset: higher risk (red), lower risk (green), and moderate risk (orange). (b) A dimension reduction of genes from the mRMR signature, using principal components analysis. (c) The survival curves were predicted by validating the ICGC-RECA dataset. (d) The principal components analysis of the ICGC-RECA dataset with genes of mRMR signature.

For internal validation of the mRMR gene signature, we performed cross-validation with 3-folds measured with AUC assessed on TCGA-KIRC with seven years of time-dependent intervals. Figure 4a shows a prediction of a random 33% sampling from TCGA-KIRC after training the regression model with 66% of samples. The Kaplan-Meier curves (Figure 4a) are evaluated by the p-values of the log-rank test, indicating the separation between patients with high, moderate and low risk. Figure 4a displays a PCA with the same predicted samples using only the expression of mRMR genes. Note that only one patient is deceased in the low-risk group, and there is a visible separation between the low-risk and high-risk groups of patients on the x-axis of PCA.

352  
353  
354  
355  
356  
357  
358

359  
360  
361  
362  
363  
364  
365  
366  
367  
368

369  
370  
371  
372  
373  
374

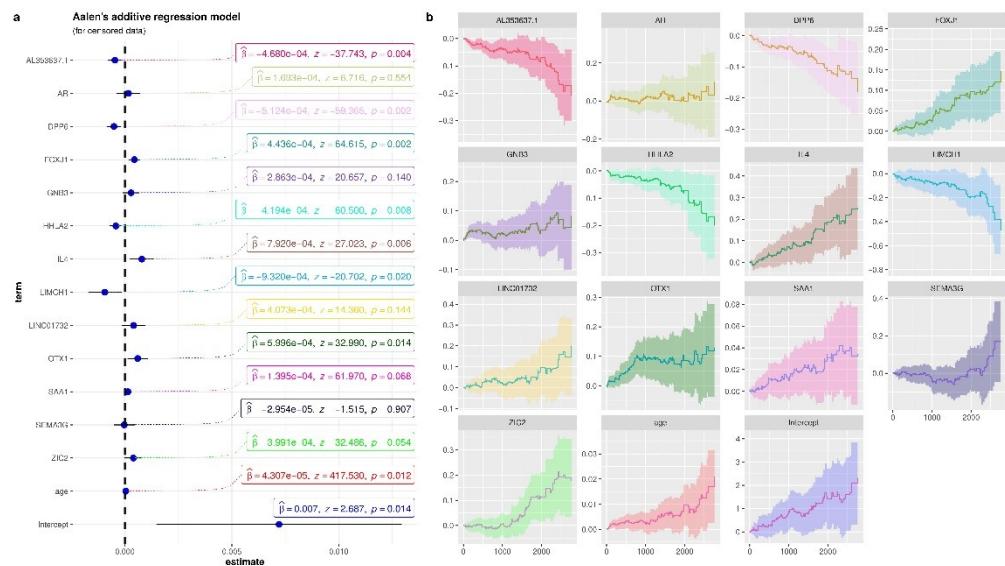
For external validation, we trained the model with the TCGA-KIRC dataset and predicted all the samples of ICGC-RECA. Analogously, in Figure 4c, a model trained with TCGA-KIRC data predicts ICGC-RECA samples in separated survival curves risks ( $p < 0.05$ ) and AUC of 0.66. In Figure 4d, we performed a PCA with mRMR gene on the same previously predicted samples of ICGC-RECA, and the x-axis also separates the centroids of the risk clusters.

375  
376

### 3.4. Biological Interpretation: Gene Contributions for Survival Risk and Enrichment Analysis

377  
378  
379  
380  
381  
382  
383  
384  
385  
386

In order to shed light on ability of each gene to predict ccRCC risk, we performed an additive regression, plotting the genes' coefficients with time-varying and covariate effects. [Similarly to the forest plot of hazard ratio regression \(Figure A8\)](#), Figure 5 shows the estimated coefficients of the increasing curves for the following significant high expression genes with a high risk of death: FOXJ1, OTX1, and IL4. On the other hand, the decreasing curves indicate that the high expression of the following genes is related to the low risk of death: AL353637.1, DPP6, HHLA2, and LIMCH1. A classical common representation of these covariate effects is the Hazard ratio in Figure 5 of the Forest Plot for Cox Proportional Hazards Model.

387  
388  
389  
390  
391

**Figure 5.** Aalen's additive cox regression model for censored data of mRMR signature, and the clinical features age and metastasis. (a) The dot-and-whisker plots with the estimated coefficients ( $\beta$ ), z-score, their confidence intervals (95%), and the p-values. (b) Curves of each term for the censored data in relation to time (days).

392  
393  
394  
395  
396  
397  
398  
399  
400  
401

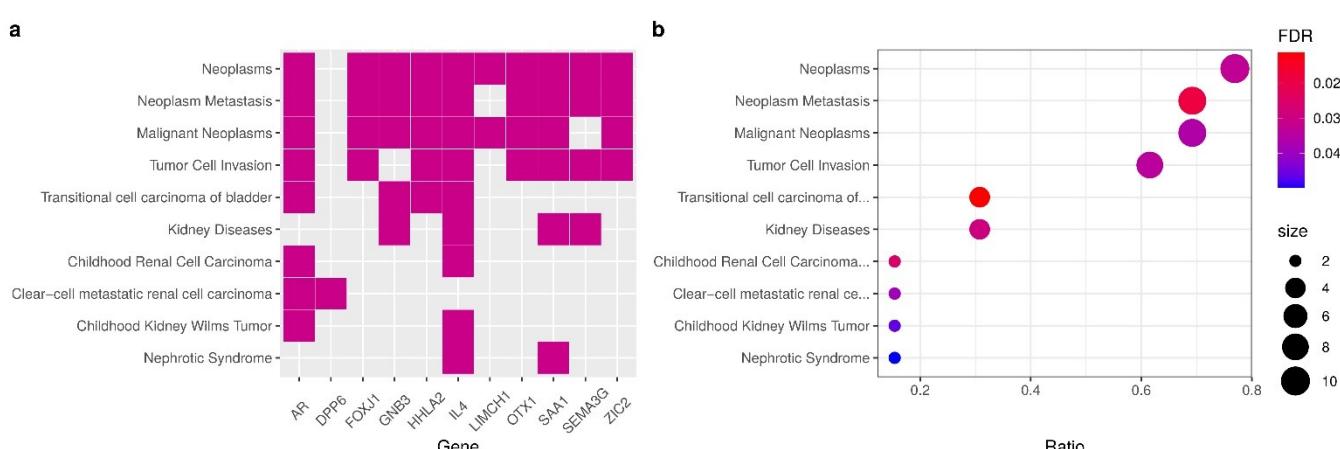
We confirmed the genes contributions to survival risk by checking protein expression according to cancer stage using the UALCAN dataset of ccRCC from Clinical Proteomic Tumor Analysis Consortium (CPTAC). As a result the gene expression by Overall Survival is corroborated with the levels of protein expression and the cancer stage. In CPTAC-ccRCC, the protein expression of genes AR, GNB3, HHLA2, LIMCH1, and SAA1 had statistical significance in some comparisons of normal samples and cancer stage (Figure A<sup>79</sup>a-e). In particular, HHLA2 protein expression in samples of Stage 1 is higher than in stage 4, but normal tissue has a lower protein expression than any tumor stage (Figure A<sup>79</sup>c). This protein expression shift is compatible with our results in Figure 5 of the decreasing curve for

402 HHLA2, and is in accordance with the TCGA-KIRC RNASeq data, since the higher  
403 expression of HHLA2 demonstrates a better prognosis (Figure A79c).

404 We verified patient survival curves by comparing the low/medium versus high  
405 expression of TCGA-KIRC data on UALCAN portal [47][46]. The above results  
406 correspond with the OS patients with low/medium versus high expression, available  
407 on the effect of expression level of patients survival. We noticed that patients with a  
408 poor prognosis had low expression of AR, DPP6, HHLA2, LIMCH1 and SEMA3G.  
409 Additionally, poor prognoses of patients can be identified with high expressions of  
410 FOXJ1, GNB3, OTX1, SAA1, and ZIC2.

411 Furthermore, in accordance with the above results, performing a Heatmap with  
412 hierarchical clustering combining RNA-Seq of patients from TCGA-KIRC and ICGC-  
413 RECA (Figure A106), we verified that the cluster of genes SAA1, OTX1, ZIC2,  
414 LINC01732, GNB3, and IL4, with high expression, is correlated with Stage T3 AJCC,  
415 Metastasis and poor prognoses. Likewise, the cluster of genes AL353637.1, AR,  
416 HHLA2, LIMCH1, SEMA3G, DPP6, and FOXJ1, with low expression, is correlated  
417 with poor prognoses.

418 To clarify the relationship between the genes and other kidney pathologies, we  
419 checked the statistical significance of multiple diseases associated with the enriched  
420 genes in the signature. Figure 6 shows a subset of 11 genes from within the  
421 signature, and most genes are related to Neoplasms, except for AL353637.1,  
422 LINC01732, and DPP6. Nevertheless, genes DPP6 and AR are enriched to clear-cell  
423 metastatic RCC diseases. We identified six genes enriched to Kidney Diseases and  
424 ccRCC (AR, DPP6, GNB3, IL4, SAA1, SEMA3G). Other enriched genes we found  
425 (AR, GNB3, HHLA2 and IL4) are related to Transitional cell carcinoma of the  
426 bladder (also known as Urothelial carcinoma). GNB3 and IL4 are enriched in both  
427 Kidney Diseases, Transitional cell carcinoma, and Neoplasm Metastasis. This  
428 enrichment analysis also confirms the results of benchmark and comparisons to the  
429 literature, indicating the importance of the selected mRMR genes in predicting  
430 ccRCC survival risk.



431  
432 **Figure 6. Gene Enrichment Analysis.** (a) Heatmap of enriched terms and relationships of  
433 genes, displaying the fold change of differential analysis of Normal tissues versus primary  
434 tumors of TCGA-KIRC samples. (b) Enrichment Analysis of gene-disease associations (GDAs)  
435 from DisGeNET (v7.0) of expert curated databases.

#### 436 4. Discussion

437 From our 13-gene signature, a subset of 8 genes were reported previously in  
438 distinct signatures for ccRCC, including other recent signatures that were not  
439 compared in our benchmark: AR [48,49][47,48], SEMA3G [19,52,53][47-49], LIMCH1

440 | [9], DPP6 [51,52][50,51], FOXJ1 [53,54][52,53], ZIC2 [11], IL4 [48,54–56], and OTX1  
441 | [12]. The concordance of this work with published signatures strengthens the  
442 | validity of our methodology to obtain a ccRCC survival signature.

443 | FOXJ1, IL4, HHLA2, and SEMA3G are immune-related genes [19,52,53][47–49],  
444 | corroborating with the high immunogenicity of ccRCC. Forkhead Box J1 (FOXJ1) is a  
445 | transcription factor, member of the FOX family, involved in ciliogenesis. Its  
446 | defective expression is associated with some inflammatory [58][57] and autoimmune  
447 | [59,60][58,59] diseases. FOXJ1 was already identified as a prognostic marker of RCC  
448 | where its expression was reported to be upregulated [54][53]. Moreover, it was  
449 | reported to be upregulated in bladder cancer [61][60], hepatocellular carcinoma [62]  
450 | [61] and colorectal cancer [63][62]. Conversely, its low expression was reported to be  
451 | correlated with gastric cancer [64][63], ependymoma and choroid plexus tumours  
452 | [65][64]. AL353637.1 is a pseudogene nearby the gene FOXB2, also belongs to the  
453 | FOX family of FOXJ1 [53][52], and contains a variant (rs115747230) associated with  
454 | chronic kidney disease [66][65]. Interleukin 4 (IL4) is a cytokine that induces  
455 | differentiation of T cells and is present in the tumor environment of many cancers.  
456 | The expression of IL4 in the tumor microenvironment can improve tumor growth  
457 | and the blockade of IL4 can delay the growth [67][66] and also can improve  
458 | immunotherapies (in mice models) such as CpG ODN or anti-OX40 AB [68][67].  
459 | Polymorphisms of the IL4 gene were associated with many cancers [69][68]. HERV-  
460 | H LTR-Associating 2 (HHLA2, also known as B7-H7) is a member of the B7-family  
461 | of immune checkpoint molecules, known to perform an inhibitory activity in human  
462 | CD4+ and CD8+ T cells by binding to their receptors [70,71][69,70]. It is known to  
463 | have a limited expression in normal tissues and to be highly expressed in cervical  
464 | adenocarcinoma [72][71], pancreatic and ampullary cancers [73][72], also widely  
465 | expressed in different subtypes of human lung cancer [70,74][69,73]. The 5-year  
466 | survival rate of patients with gastric cancer was significantly higher in patients with  
467 | HHLA2 highly expressed [75][74]. In particular, the overexpression of HHLA2 in  
468 | patients after surgery was identified to promote ccRCC progression when compared  
469 | to normal adjacent tissue [76][75], which corresponds with our results regarding  
470 | HHLA2 expression. The knockdown of HHLA2 decreased the expression of genes  
471 | related to the cell cycle, as well as the ability of the cells to migrate and invade [76]  
472 | [75]. SEMA3G belongs to the family of class-3 semaphorins, and studies indicate  
473 | that this gene is linked to kidney diseases [77,78][76,77], suggesting important roles  
474 | with neuropilin and plexin families in the etiology of cancer [79][78], and it is also an  
475 | inhibitor of glioma progression by competing with VEGF for receptor NRP1 [80][79].  
476 | In single-cell RNA-seq study of kidney with transplant biopsy, SEMA3G activates  
477 | an angiogenic program [81][80]. Patients with high expression of SEMA3G and AR  
478 | have better prognoses according the survival analysis of UALCAN RNASeq data  
479 | [47][46]. The presence of immune-related genes in our signature strengthens the  
480 | approach of focusing on the genes from the immune system to build a prognostic  
481 | signature [49,82][48,81]. Our findings reinforce that HHLA2 is an important  
482 | immune-related biomarker of ccRCC.

483 | The genes AR, OTX1, and ZIC2 are transcription factors. In particular,  
484 | Androgen Receptor (AR) is a transcription factor whose activity is highly critical to  
485 | prostate cancer evolution [83][82]. The expression of AR-V7, its isoform, which is  
486 | encoded by splice variant 7 in circulating tumor cells of prostate cancer, was  
487 | reported to be associated with drug resistance [84][83]. AR interacts with VHL to  
488 | modulate the metastasis of ccRCC [85][84], and AR inhibition can attenuate RCC  
489 | progression [86][85]. The epigenetic control of AR co-regulates lysine-specific  
490 | histone demethylase 1 (LSD1) in Kidney Cancer development, and the LSD1  
491 | inhibitor can reduce growth of kidney cancer cells [87][86]. Also, AR could suppress  
492 | ccRCC cell progression by increasing the expression of circRNA circHIAT1 [88][87].

In addition, in vitro research and in vivo mouse model studies indicate that AR mediates lncRNA-TANAR signals that might play a crucial role in ccRCC progression and metastasis [89][88]. The studies above indicate that AR might be a promising drug target for treatment of ccRCC. OTX1 is a protein-coding gene of the bicoid sub-family of homeodomain-containing transcription factors, involved in differentiation of young neurons of the deeper cortical layers, and in proliferative zones of the neocortex [90][89]. OTX1 is related to breast cancer, medulloblastomas, colorectal cancer, hepatocellular carcinoma and bladder cancer [12]. The zinc finger of the cerebellum 2 (ZIC2) is a transcription factor with an important role in neural development and mutations of ZIC2, which could lead to brain malformations [91,92][90,91]. ZIC2 is an oncogenic with overexpression correlated to progression of epithelial ovarian tumors [93][92]. In breast cancer, low expression of ZIC2 was correlated with poor outcomes and acts as a tumor suppressor by regulating STAT3 [94][93]. ZIC2 also upregulates gene RUNX2 and promotes ccRCC progression through inhibition of tumor suppressor NOLC1 [95][94].

Lim and Calponin Homology Domains 1 (LIMCH1) is an actin-stress-fibers-associated protein, a gene encoding zinc-binding protein, and is known to negatively regulate cell-spreading and migration [96][95]. It was reported to be down-regulated in malignant lung tissue [97][96] and upregulated for breast cancer [98][97]. LIMCH1 is upregulated with a strong association to poor prognoses, representing a potential biomarker for cervical cancer treatment [99][98]. According to survival analysis of the Human Protein Atlas [100][99], LIMCH1 is also a prognostic gene, whose high expression is associated with favorable outcomes in renal cancer [101][100].

Dipeptidyl Peptidase Like 6 (DPP6) is a type II membrane glycoprotein known to regulate potassium channels and is mainly expressed in the central nervous system [102][101]. The methylation of CG sites in the DPP6 promoter was reported to be in greater numbers in tumor samples as when compared to normal samples from pancreatic ductal carcinoma, thus, the hypermethylation of DPP6 promoter is associated with poor overall survival [103][102]. The hypermethylation of DPP6 was associated with a high grade tumor in ccRCC [52][51]. Also, high expression of DPP6 was reported to be correlated with good prognoses in patients with breast cancer [104][103].

Guanine Nucleotide Binding Protein Beta Polypeptide 3 (GNB3) is involved in various transmembrane signaling systems such as in GTPase activity. Some studies associate the polymorphism GNB3-C825T with cholangiocarcinoma [106][105] and thyroid carcinoma [106][104], but another study discarded a relationship with the risk for breast cancer [107][106].

Serum Amyloid A 1 (SAA1) is an acute-phase protein mainly produced by hepatocytes in response to infection, tissue injury and malignancy. SAA1 modulates neutrophil function in the context of cancer [108][107]. SAA1 gene expression in patients with RCC is associated with poor prognosis [109][108]. According to survival analysis of Human Protein Atlas [100][99], SAA1 is also a prognostic gene with high expression for unfavorable outcomes in renal cancer [110][109]. Moreover, multiple mutation variants of SAA1 have been identified in patients with RCC [111][110].

LINC01732 is affiliated with the long non-coding RNAs (lncRNAs) class. To the best of our knowledge, there are no publications regarding LINC01732 at this time. Nevertheless, increasing evidence suggests that lncRNAs play critical roles in tumor development of RCC [112][111]. Further research could be executed to understand other lncRNAs, including LINC01732.

Since alterations in expression of different genes from the same pathway have higher impacts on gene function, we performed an enrichment analysis and

546 identified the pathways of Urothelial Carcinoma, Chronic Kidney disease, and  
547 Transitional cell carcinoma, Nephrolithiasis. Although the concurrence of RCC and  
548 Urothelial Carcinoma is clinically rare [113][112], previous studies reported the  
549 identification of clear cell tumors in general bladder carcinomas [114,115][113,114].  
550 On Nephrolithiasis, studies showed that Kidney stones were associated with  
551 increased papillary RCC risk but not clear-cell RCC risk [116][115].

552 We compared our signature in a benchmark with fourteen other signatures  
553 already published in the literature. All of the gene signatures [2,8–13,48,116–121]  
554 compared in this work use TCGA as their main training set to build their models.  
555 The studies reviewed have AUC-ROC between 0.568 to 0.884 with possible values  
556 ranging from 0 to 1, and the number of genes in each signature range from 3 to 66.  
557 Some studies use a different number of patients due to the distinct filtering  
558 approaches that the authors adopted, in addition to the updates of versions of  
559 TCGA-KIRC clinical data. The least absolute shrinkage and selection operator  
560 (LASSOLasso-Cox) was the most-used model approach to build the signatures  
561 [2,8,9,11–13,48,118–120], but network-based models with protein-protein interaction  
562 (PPI), aside from being an elegant approach to retrieve information from data, can  
563 also be used for this purpose [117,118][116,117].

564 This work consists of a pure *in silico* and data-driven study, and other analyses  
565 could be corroborated in the future with *in vitro* or *in vivo* experiments [123][122]. As  
566 next steps, we are applying this approach of machine learning and feature selection  
567 to find potential cancer biomarkers in multiples levels of biological data available in  
568 TCGA, such as long non-coding RNAs (lncRNAs), methylation, single-nucleotide  
569 variants (SNV), and copy number variants (CNV). In future works, we are  
570 expanding the machine learning approach presented in this work to find potential  
571 cancer biomarkers using multiples levels of biological data available in TCGA, by  
572 analyzing and integrating data of long non-coding RNAs (lncRNAs), methylation,  
573 single-nucleotide variants (SNV), and copy number variants (CNV).

## 574 5. Conclusions

575 We obtained satisfactory results combining RNASeq data (ICGC and TCGA),  
576 survival information, and machine learning strategies. We identified two clusters of  
577 genes with high expression (SAA1, OTX1, ZIC2, LINC01732, GNB3 and IL4) and  
578 low expression (AL353637.1, AR, HHLA2, LIMCH1, SEMA3G, DPP6, and FOXJ1)  
579 which are correlated with poor prognoses. We validated our 13-gene signature for  
580 ccRCC and confirmed our results with the literature, and by comparing each cancer  
581 stage of ccRCC with CPTAC and the survival effects of gene expression of  
582 individual genes in TCGA. Our signature for ccRCC indicates potential applications  
583 for personalized medicine in clinical practice.

584 Our main goal was to compare distinct gene signatures from literature and  
585 generate new gene signatures using feature selection methods. We contributed by  
586 providing a list of new genes, some of them not previously reported as biomarkers  
587 for ccRCC. The gene signature created by the mRMR method achieved a score of  
588 0.82 with AUC, being the best performer. We identified two clusters of genes with  
589 high expression (SAA1, OTX1, ZIC2, LINC01732, GNB3 and IL4) and low expression  
590 (AL353637.1, AR, HHLA2, LIMCH1, SEMA3G, DPP6, and FOXJ1) which are  
591 correlated with poor prognosis. We validated our 13-gene signature for ccRCC and  
592 confirmed our results with the literature, and by comparing each cancer stage of  
593 ccRCC with CPTAC and the survival effects of gene expression of individual genes  
594 in TCGA. We believe that further studies on the involvement of these genes in renal  
595 carcinogenic processes can improve the understanding of cancer biology. After  
596 experimental validations, new possible applications in clinical practices can benefit  
597 from the biomarker found with machine learning and feature selection.

598  
599  
600  
601  
602  
603  
604  
605  
606

**Author Contributions:** Conceptualization, Patrick Terrematte, Beatriz Stransky and Adrião Dória Neto; Data curation, Patrick Terrematte and Josivan Justino; Formal analysis, Patrick Terrematte and Dhiego Andrade; Funding acquisition, Adrião Dória Neto; Investigation, Patrick Terrematte; Methodology, Patrick Terrematte and Josivan Justino; Project administration, Adrião Dória Neto; Resources, Patrick Terrematte; Software, Patrick Terrematte; Supervision, Beatriz Stransky and Daniel Araújo; Validation, Patrick Terrematte, Beatriz Stransky and Daniel Araújo; Visualization, Patrick Terrematte; Writing – original draft, Patrick Terrematte and Dhiego Andrade; Writing – review & editing, Dhiego Andrade, Beatriz Stransky, Daniel Araújo and Adrião Dória Neto.

607  
608  
609  
610

**Funding:** Patrick Terrematte was funded by the Federal Rural University of Semi-arid. Dhiego Souto was funded by grants numbers 88887.161820/2017-0 and 88887.600071/2021-0 of Brazilian Funding agency CAPES - National Coordination of High Education Personnel Formation Programs. The APC was funded by the Federal University of Rio Grande do Norte.

611  
612

**Institutional Review Board Statement:** This study did not require ethical review and approval because it performed a secondary analysis of publicly available data.

613

**Informed Consent Statement:** Not applicable.

614  
615  
616  
617  
618  
619  
620  
621

**Data Availability Statement:** Publicly available datasets were analyzed in this study. The results shown here are based upon data generated by the TCGA Research Network [124][123]. The TCGA-KIRC (version 07-19-2019) [125][124] is available via UCSC Xena Browser [18,126] [17,125] (accessed on 05 March 2022), and the ICGC-RECA is available via ICGC Data Portal [19,20][18,19] (accessed on 05 March 2022). Code used for analyses and to produce the figures is available at: [http://hungria.imd.ufrn.br/~terrematte/gene\\_signature](http://hungria.imd.ufrn.br/~terrematte/gene_signature). Once this manuscript is accepted, the code will be publicly available at: [https://github.com/terrematte/gene\\_signature](https://github.com/terrematte/gene_signature).

622 |  
623  
624  
625  
626  
627

**Acknowledgments:** The authors would like to thank Isa Goldberg, Tayná da Silva Fiúza, Iara Dantas de Souza and Raul Maia Falcão for their suggestions and critical reading of the draft manuscript. The authors also thank the Metropole Digital Institute (IMD), the Center for High Performance Computing (NPAD) (<https://npad.ufrn.br>) and the Multidisciplinary Bioinformatics Environment (BioME) at UFRN for providing computing resources for data processing.

629  
630  
631  
632

**Conflicts of Interest:** The authors declare that they have no competing interests, and that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

633

## Abbreviations

AJCC	American Joint Committee on Cancer
TCGA	The Cancer Genome Atlas
ICGC	International Cancer Genome Consortium
KIRC	Kidney Renal Clear Cell Carcinoma
RECA	Renal Cell Cancer
ccRCC	clear cell Renal cell carcinoma
mRMR	Minimum Redundancy Maximum Relevance
AUC	Area Under the Curve
ROC	Receiving Operator Characteristics
PCA	Principal Component Analysis
RFE	Recursive Feature Elimination
GBM	Generalized Boosted Regression Model
Rpart	Recursive Partitioning and Regression Trees
XGBoost	eXtreme Gradient Boosting
CoxPH	Cox proportional hazards regression model

## 635 | Appendix A

636 | **Table A1.** Gene signatures of ccRCC after exclusion criteria. The PubMed query—was  
 637 | conducted in January 2021ef using the terms was: (renal OR kidney) AND (clear cell) AND  
 638 | (cancer) AND (prognosis OR survival OR outcomes) AND (gene signature AND regression),  
 639 | and filtering the years of 2015 to 2020.

Title & Code in Figure 3	Gene Signature
Prognostic gene signature identification using causal structure learning: applications in kidney cancer [117] <b>Code:</b> Ha.2014.CaInfo	ETV5, CREB3L1, GMPS, RBM15, SEPT6, TTL, ARID1A, ERCC5, TFG, FLT3, SLC34A2, FAM46C, PER1, DDB2, NACA, MLLT10, HMGA1, TCF12, RUNX1, CANT1, REL, ZNF331, JAZF1, ASPSCR1, PLAG1, NOTCH1, TAL2, ERCC2, SMARCA4, DNMT3A, HOXA11, GNAS, CHEK2, HLF, GNAQ, ETV6, SET, KIF5B, TRRAP, CDKN2C, VHL, RPL22, CHN1, STAT3, CDK4, CD274, KTN1, CYLD, BRD3, TRIM33
A Five-Gene Signature Predicts Prognosis in Patients with Kidney Renal Clear Cell Carcinoma [8] <b>Code:</b> Zhan.2015.CMMM	CKAP4, ISPD, MAN2A2, OTOF, SLC40A1
A four-gene signature predicts survival in clear-cell renal-cell carcinoma [123] <b>Code:</b> Dai.2016.Oncotarget	PTEN, PIK3C2A, ITPA, BCL3
Identification and validation of an eight-gene expression signature for predicting high Fuhrman grade renal cell carcinoma [17] <b>Code:</b> Wan.2017.IJ Cancer	ATOH8, ATP1A3, C10orf4, C17orf79, CHMP4C, CNGA1, EDA, FBXL3, GMDS, ISL2, KISS1, KLF2, MYADML2, NCRNA00116, OAZ1, ODZ3, PLA2G15, PPP1R1A, RAB40A, RRAS, SPOCK1, SQSTM1, TXNDC16, VAMP3
Comprehensive assessment gene signatures for clear cell renal cell carcinoma prognosis [9] <b>Code:</b> Chang.2018.Medicine	INTS8, GTPBP2, ANK3, SLC16A12, LIMCH1, Hsa-mir-374a
A five-gene signature may predict sunitinib sensitivity and serve as prognostic biomarkers for renal cell carcinoma [118] <b>Code:</b> YuanLeiChen.2018.JCP	BIRC5, CD44, MUC1, TF, CCL5
A Gene Signature of Survival Prediction for Kidney Renal Cell Carcinoma by Multi-Omic Data Analysis [18] <b>Code:</b> Hu.2019.IJMS	BID, CCNF, DLX4, FAM72D, PYCR1, RUNX1, TRIP13
Prognostic value of a gene signature in clear cell renal cell carcinoma [10] <b>Code:</b> LiangChen.2018.JCP	CENPW, FOXM1, NUF2
Identification of a 5-Gene Signature Predicting Progression and Prognosis of Clear Cell Renal Cell Carcinoma [12] <b>Code:</b> Pan.2019.MSM	OTX1, FOXE1, FAM83A, HMGA2, KRT6A, DPYSL5, ANXA8, MATN4, ROS1, CSMD3, MAGEC3, AMER2, CPLX2, PI3, KRT13, ERVV-2, ERVFRDE1, ANKFN1, VTN, NFE4, ZNF114
Construction and Validation of a 9-Gene Signature for Predicting Prognosis in Stage III Clear Cell Renal Cell Carcinoma [13] <b>Code:</b> Wu.2019.FrontiersOncology	ATP6V1C2, PCSK1N, PREX1, ANK3, HLA-DRA, SELENBP1, TYRP1, GABRA2, SERPINA5
Construction and validation of a seven-gene signature for predicting overall survival in patients with kidney renal clear cell carcinoma via an integrated bioinformatics analysis [11] <b>Code:</b> Jiang.2020.ACS	PODXL, SLC16A12, ZIC2, ATP2B3, KRT75, C20orf141, CHGA
A 14 immune-related gene signature predicts clinical outcomes of kidney renal clear cell carcinoma [19] <b>Code:</b> Zou.2020.PeerJ	TXLNA, SEMA3G, AR, BID, IL20RB, CCR10, BMP8A, SEMA3A, CCL7, GDF1, KLRC2, LHB, FGF17, IL4
A seven-gene signature model predicts overall survival	APOLD1, C9orf66, G6PC, PPP1R1A, CNN1G, TIMP1, TUBB2B

in kidney renal clear cell carcinoma [2]

**Code:** LingChen.2020.Hereditas

Identification of gene signature for treatment response to guide precision oncology in clear-cell renal cell carcinoma

[124]

**Code:** DCosta.2020.SciReports

ANGPT4, EDN1, VEGFA, ESM1, FLT1, KDR, CD34, PECAM1, NOTCH1, EDNRB, STIM2, FYN, VWF, GJA1, MCF2L, PPM1F, PTPRB, HEY1, ETS1, EXOC3L2, TBXA2R, TCF4, S1PR1, SLC9A3R2, NES, NFATC1, NOS3, PDE2A, CORO1A, CCR5, CXCR3, PTK2B, WAS, CD72, IL16, FYB1, FASLG, FERM1T3, FOXP3, XCL2, CD3E, CD7, LAX1, CD38, LCP1, LCP2, ITK, LAT, LCK, GRK2, CCL4, CCL5, CD2, PRF1, TIGIT, GZMA, GZMB, CD8A, CTLA4, EOMES, PDCD1, PYHIN1, SLA2, LTA, PSMB8, PSMB9

640  
641  
642

**Table A2.** New gene signatures of ccRCC obtain by the state-of-art of Machine Learning for Feature Selection methods: Recursive Feature Elimination, Boruta, Rpart, GBM and XGBoost for Survival.

Code	Method	N. Genes	Gene Signature
GBM	Filtering with Generalized Boosted Regression Models for Cox Proportional Hazard	30	AC084117.1, CRHBP, LINC00973, ITPKA, IGFN1, C14orf37, OTX1, LINC02446, HOTTIP, NEIL3, ZIC5, CCDC154, IL4, AC008663.1, FER1L4, DUSP5P1, AL078604.2, KRT6A, SPATC1L, RTL1, LINC01597, CRABP1, RASGRP3, C3orf85, AL034399.1, TRIM4, LINC00475, ADAMTS14, DPP6
Rpart	Filtering with Recursive partitioning for survival trees	30	TROAP, KIF18B, AURKB, LINC00973, AC003092.1, G6PC, ZNF181, MYBL2, FOXM1, NUF2, POU4F1, APOM, AR, NPHS1, AC018638.2, MERTK, AC098679.1, AL353637.1, IYD, C17orf80, SLC12A3, CDCA2, LINC02362, SRD5A3, EIF3F, AC138393.1, MCC, WFIKK1, ALDOB, APOL5
XGBoost	Filtering with XGBoost for Survival Analysis	30	LINC00973, LINC01271, CHAT, SPIC, AL355796.1, DLK1, ZIC5, LINC01700, ENTPD6, ATOH8, C14orf37, WNT7B, THEG, AC084117.1, ADA2, DCSTAMP, AL450311.2, A3GALT2, CNTNAP3B, TBC1D27, BIRC7, LINC00943, LINC01529, OR4C6, FAM47E, BCL3, AC105118.1, AL359736.1, SLC44A3, LINP1
Boruta	Wrapper Boruta with XGBoost for Survival Data	43	Age, ZIC2, CHAT, AMH, OTX1, BARX1, TROAP, CKAP4, ITPKA, NUF2, KRT75, KIF18B, SLC18A3, AL355796.1, RPL10P19, LINC02154, LINC00973, IL4, HOTAIRM1, Z84485.1, LINC02362, CASP9, CCNF, RTL1, BID, CHGA, RANBP3L, ZIC5, SLC16A12, SPATC1L, CD44, KRI1, RUFY4, AC073324.1, AC091812.1, AC156455.1, AGAP6, AC128685.1, SEMA3G, IGFN1, KLRC2, ANXA8, AURKB
RFE	Wrapper with Recursive Feature Elimination	89	A3GALT2, AC006450.2, AC073324.1, AC093520.1, AC103925.1, AC120498.6, AC128685.1, AC156455.1, ADAMTS14, AL355796.1, AL592494.1, AL606519.1, AMH, ANK3, ANXA8, AP000697.1, AP001029.1, AURKB, BARX1, BIRC5, C20orf141, CCNF, CDC42P2, CENPW, CHAT, CHGA, CKAP4, CRHBP, DLX4, DMRT3, DUSP5P1, G6PC, GOLGA6L2, GOLGA6L7P, HAMP, HAO1, HOTAIRM1, HP, IGFN1, IGHJ3P, IL20RB, IL4, ISL2, ITPKA, KIF18B, KLRC2, KRT75, KRT78, LINC00051, LINC00460, LINC00524, LINC00896, LINC00973, LINC01234, LINC01501, LINC01655, LINC01700, LINC01956, LINC02154, LINC02362, NEIL3, NFE4, NUF2, OTX1, PAEP, PGLYRP2, PI3, PITX1, PLG, PTPRB, RALYL, RPL10P19, RTL1, SAA1, SAA2, SAA4, SIM2, SLC16A12, SLC18A3, TGM3, TRIP13, TROAP, VSX1, WFDC10B, Z84485.1, ZIC2, ZIC5, ZPLD1
mRMR	Ensemble of Min-redundancy and Max-relevance with survival data	65	AR, AL353637.1, DPP6, FOXJ1, GNB3, HHLA2, IL4, LIMCH1, LINC01732, OTX1, SAA1, SEMA3G, ZIC2

643

644

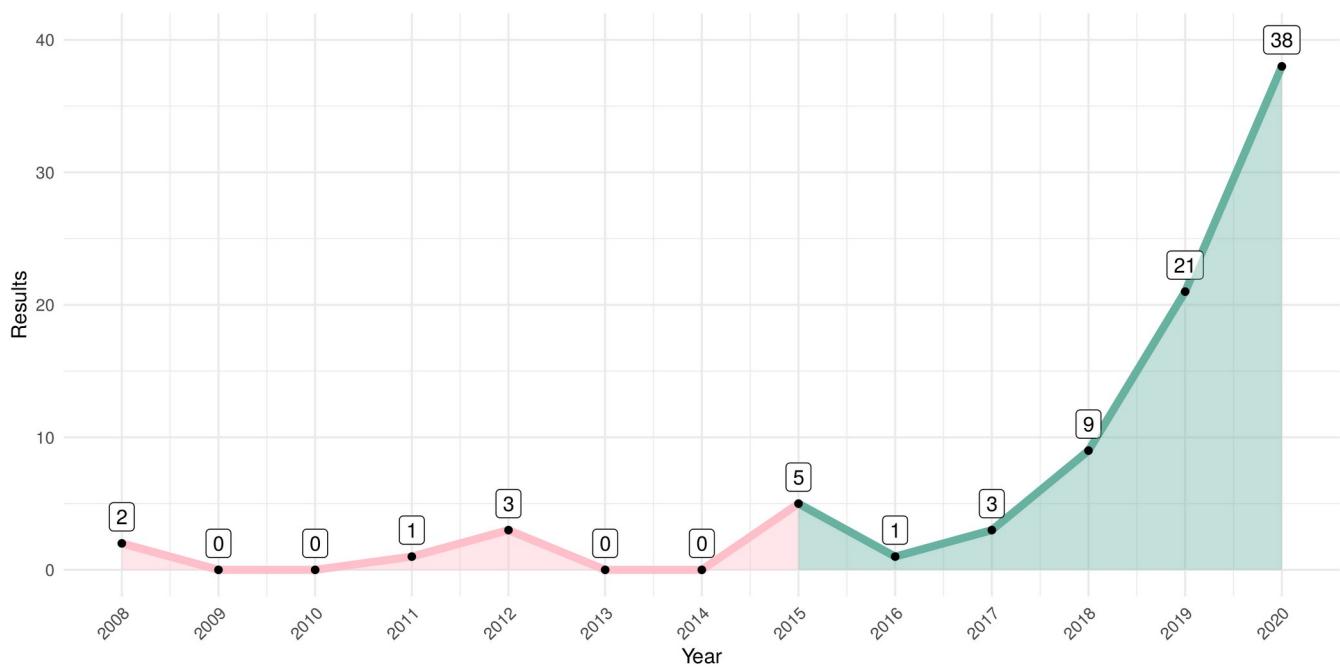
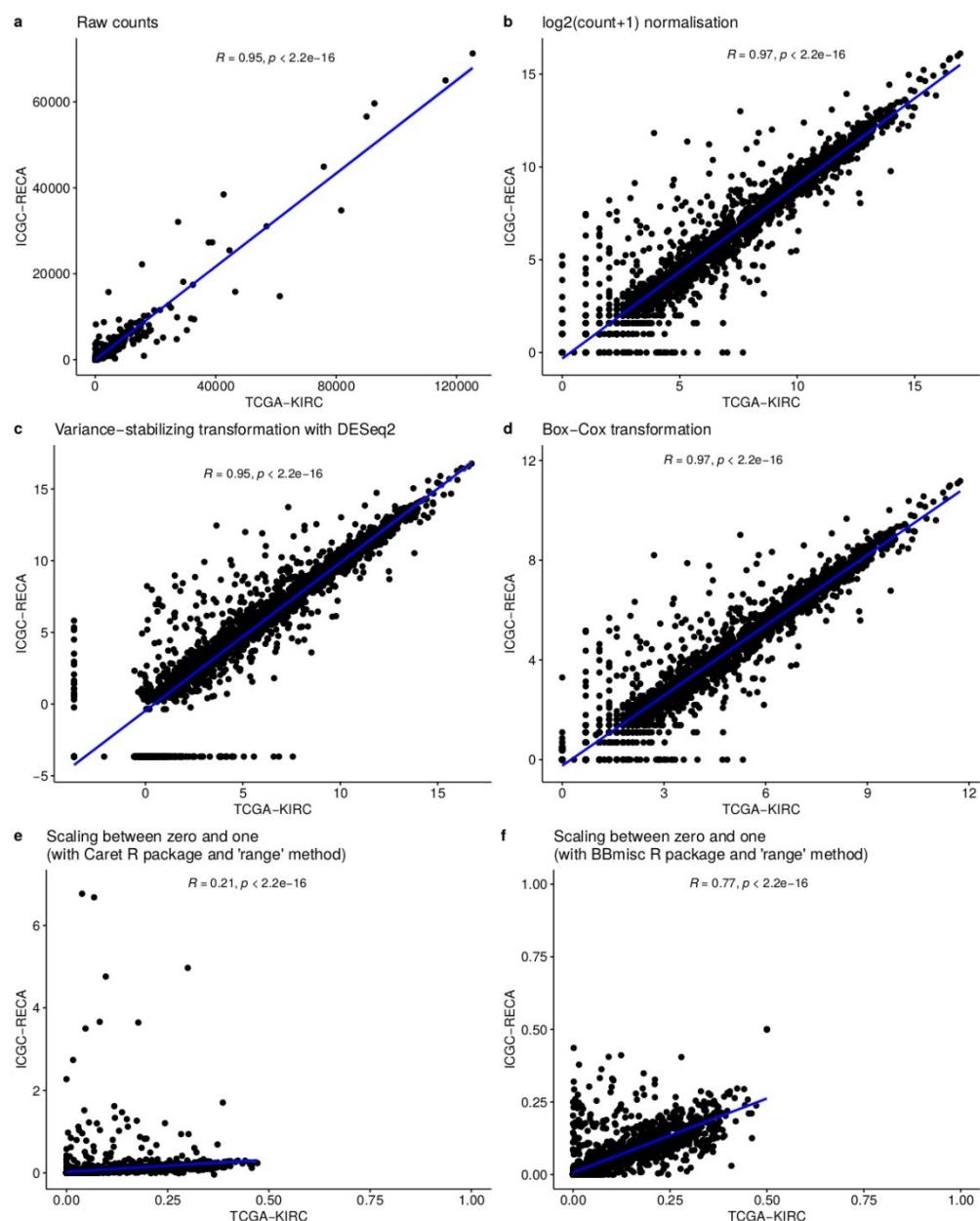


Figure A1. Number of papers published on PubMed by year on query performed in January 2021. Initially, in green, the gene signatures published in the period of 2015 to 2020 were selected to be compared. After the exclusion criteria, we obtained the 14 gene signatures.

645  
646  
647

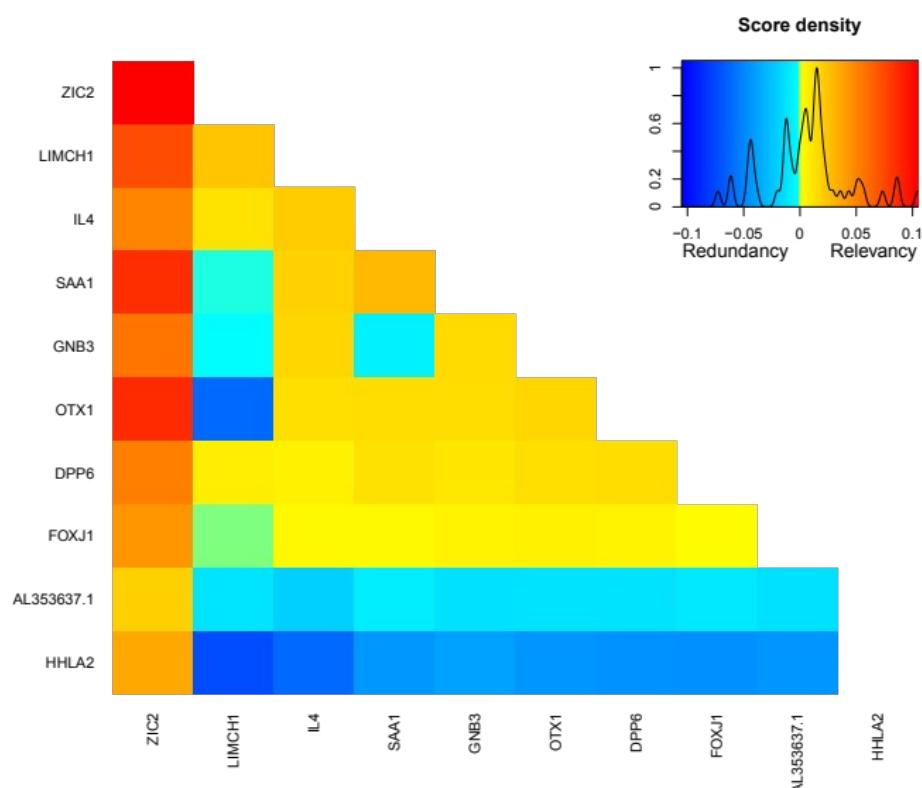
648 |



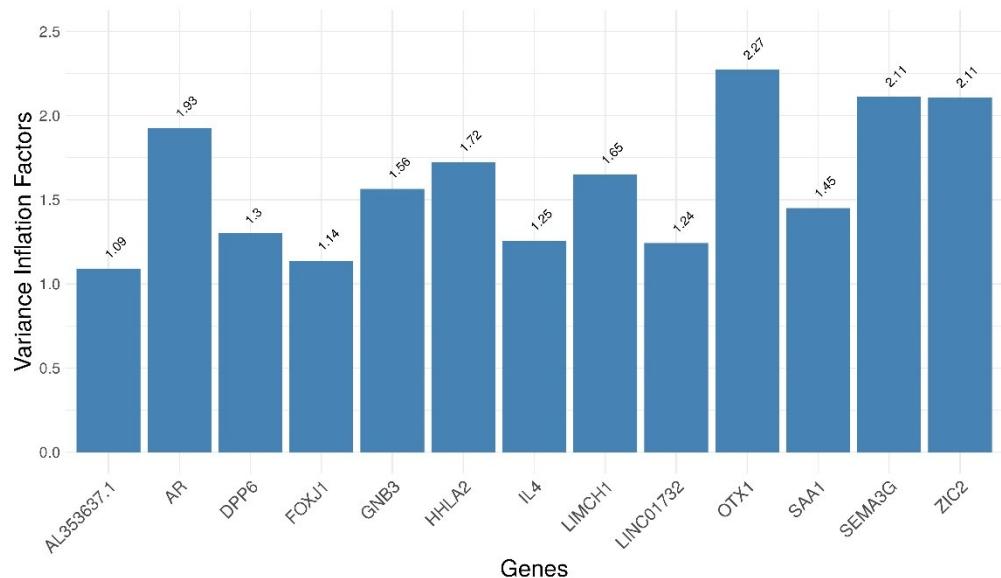
**Figure A21.** Scatter plot of median of gene expression comparing TCGA-KIRC and ICGC-RECA gene expression. (a) Raw counts. (b)  $\log_2(\text{count}+1)$  normalisation. (c) Variance-stabilizing transformation with DESeq2. (d) Box-Cox transformation. (e) Scaling between zero and one (with Caret R package and 'range' method). (f) Scaling between zero and one (with BBmisc R package and 'range' method).

649  
650  
651  
652  
653

654



**Figure A3.** Variable Ranking Based on Mutual Information of 10 most important genes of mRMR 13-gene signature of ccRCC. The most representative genes with respect to AJCC Staging of TCGA dataset.

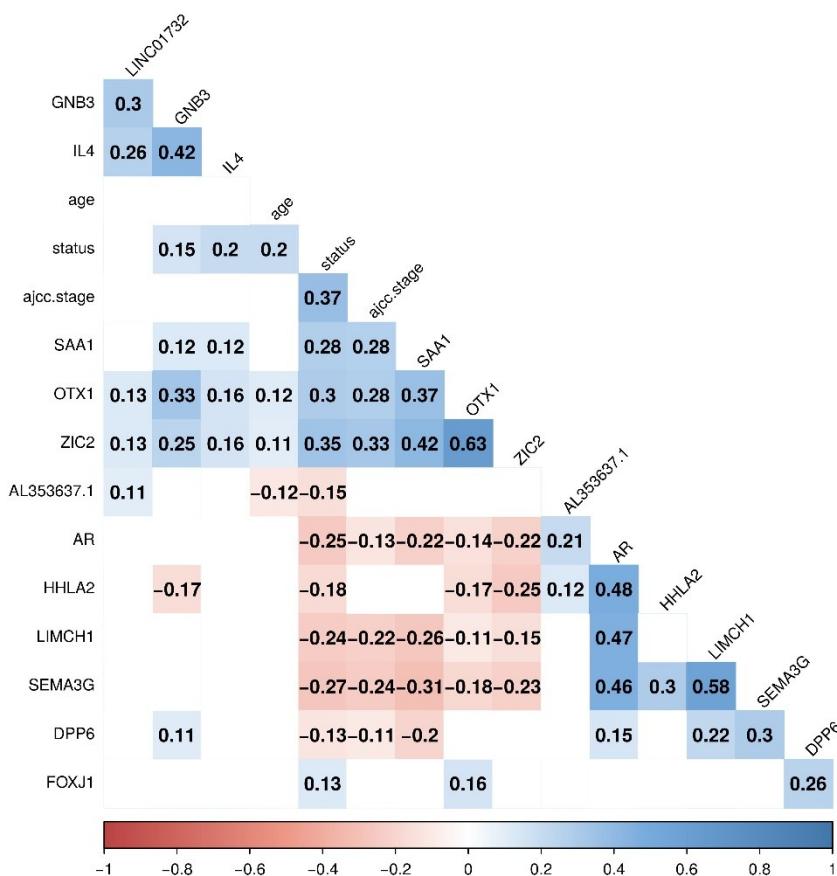


**Figure A24.** Collinearity Analysis with Variance Inflation Factors 13-gene signature of ccRCC. None of genes had Variance Inflation Factors  $\geq 5$ , indicating no collinearity or redundancy on the signature.

658 |

659

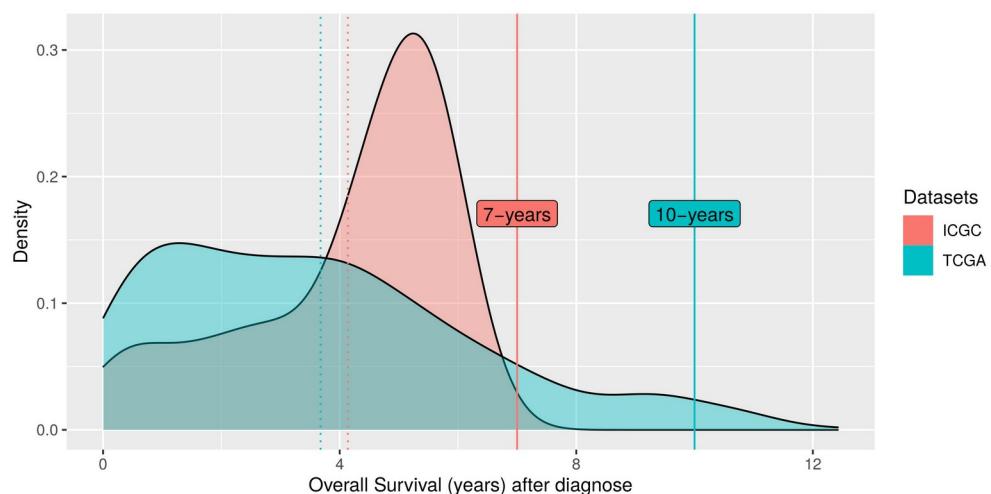
660



**Figure A5. Correlation Analysis between genes of mRMR 13-gene signature of ccRCC. No strong correlation between gene  $\geq 0.70$  were found, including the clinical data of Age, Overall survival status and AJCC Staging.**

**Figure A3. Correlation Analysis between genes of mRMR 13-gene signature of ccRCC. No strong correlation between gene  $\geq 0.70$  were found, including the clinical data of Age, Overall survival status and AJCC Staging.**

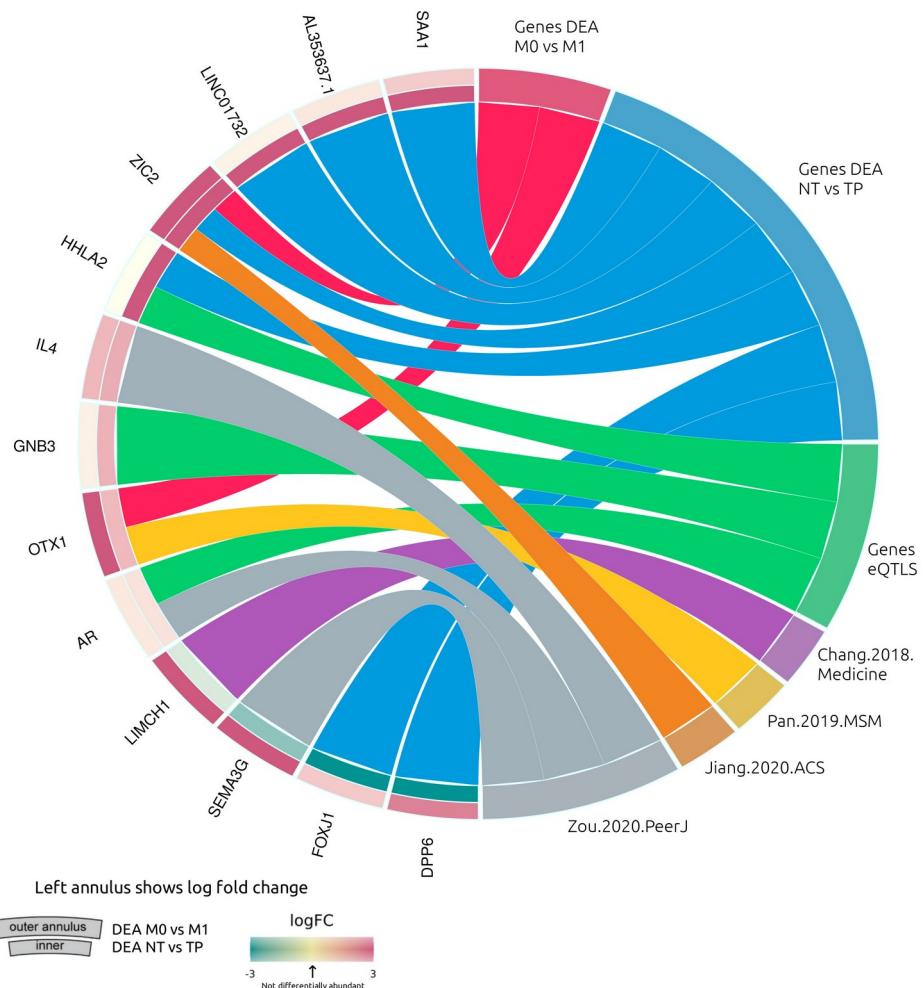
**Figure A4. Variable Ranking Based on Mutual Information of 10 most important genes of mRMR 13-gene signature of ccRCC. The most representative genes with respect to AJCC Staging of TCGA dataset.**



**Figure A6. Density plot of the distribution of the patient's overall survival in TCGA-KIRC and ICGC-RECA.** The dotted line indicates the mean of distributions, and the solid lines indicate the time prediction used for internal and external validations. We restrict the 10-years prediction for TCGA-KIRC to exclude outliers in the long tail of the density plot of the patient's overall survival. For the ICGC-RECA dataset, we decided to maintain a 7-years prediction in order to include all samples, and limit the time prediction to the range of distribution of this dataset for external validation.

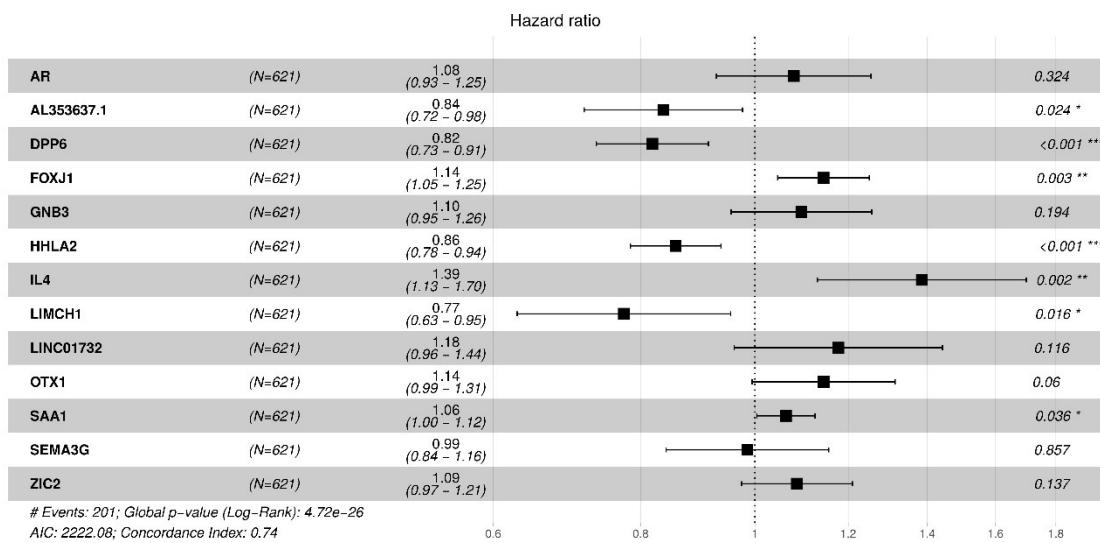
671  
672  
673  
674  
675  
676  
677

678 |  
679 |



**Figure A7. Circular diagram of mRMR gene signature and the source of genes DEA, genes from GTEx portal of expression quantitative trait loci (eQTLs) in Kidney Cortex, and gene signatures from the literature**

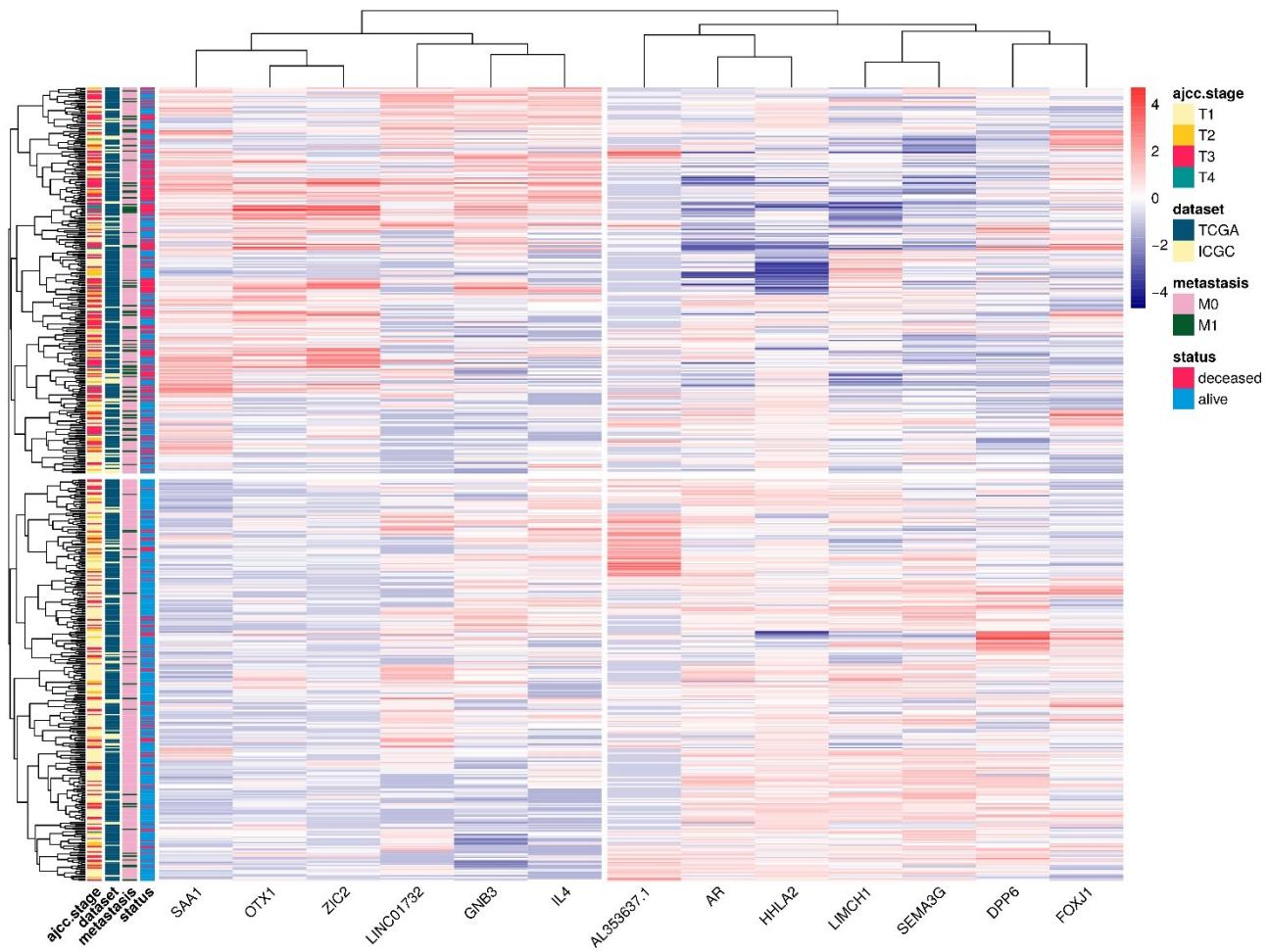
683



684

685  
686

**Figure A85.** Forest Plot for Cox Proportional Hazards Model [displaying the significative genes \(AL353637.1, DPP6, FOXJ1, HHLA2, and SAA1\).](#)



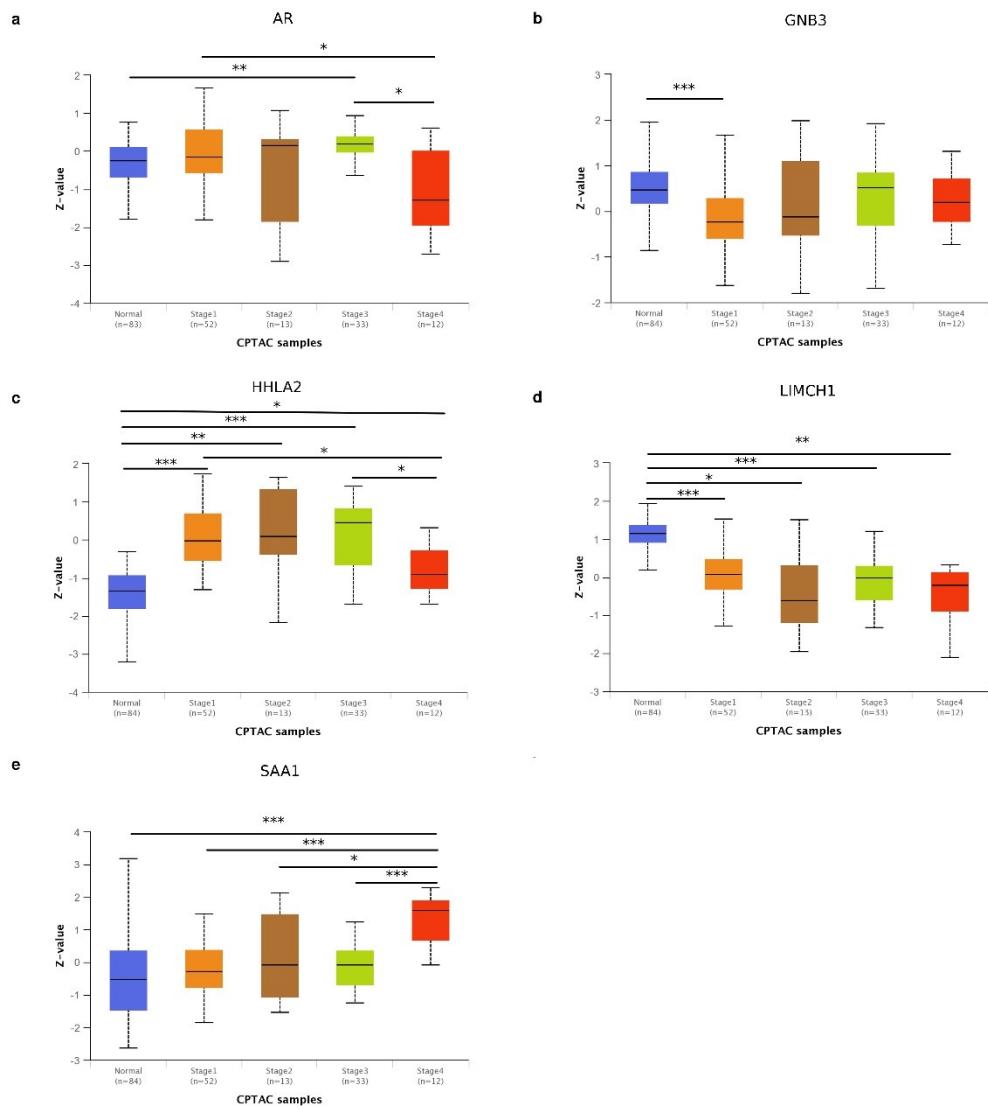
687

688  
689

**Figure A6.** Heatmap with Hierarchical clustering combining RNA-seq expression of patients on TCGA-KIRC and ICGC-RECA. Columns are genes of the mRMR signature. Rows indicate RNA-seq expression of 590

690  
691

**patients of TCGA-KIRC and ICGC-RECA. Data of patients with distant metastasis that cannot be assessed (MX) were removed in order to clarify the clustering.**

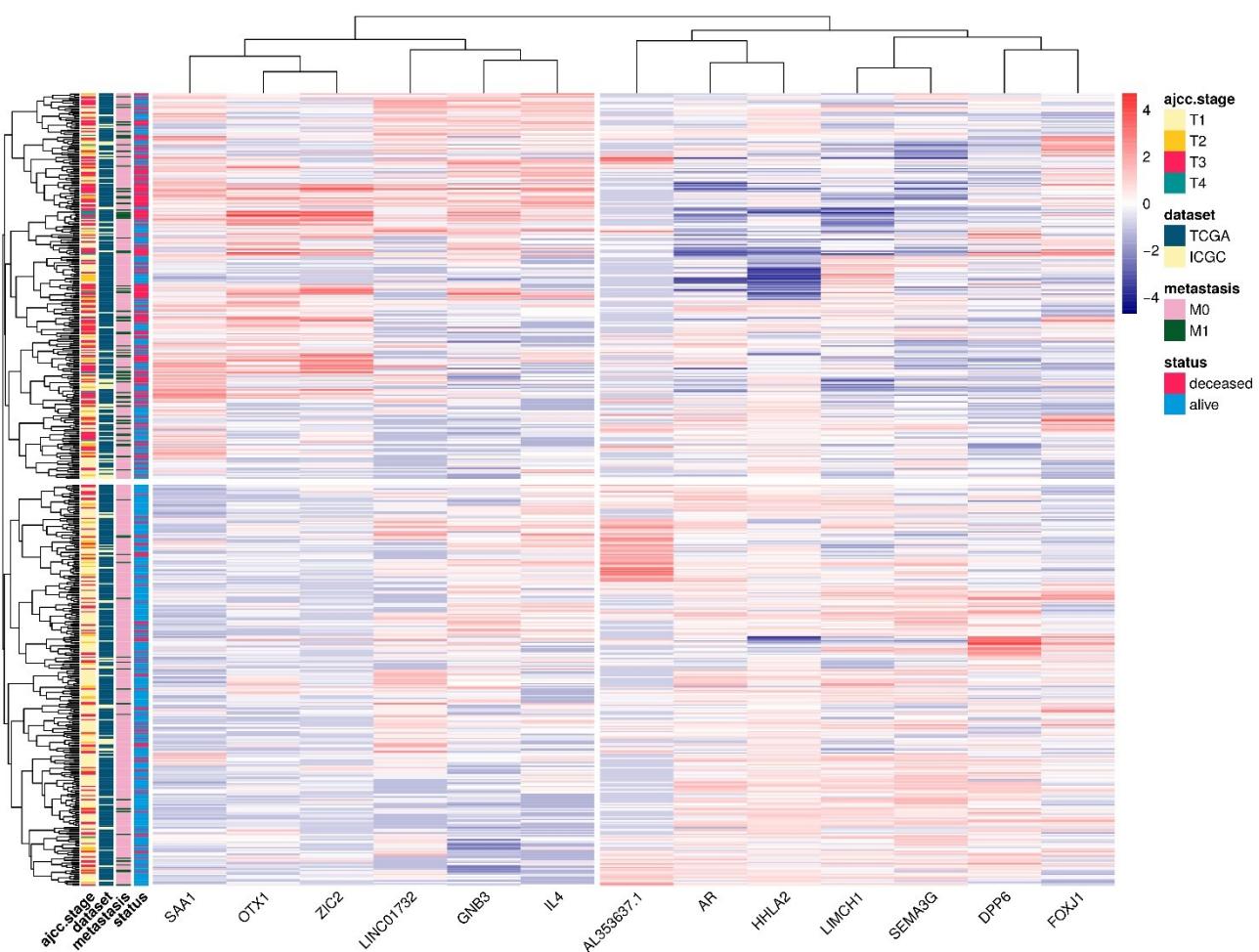


692

693  
694  
695  
696

**Figure A79.** Analysis by UALCAN with data of ccRCC from Clinical Proteomic Tumor Analysis Consortium (CPTAC) [47][46] (<http://ualcan.path.uab.edu/>). Protein expression of Genes identified in CPTAC: AR, GNB3, HHLA2, LIMCH1, and SAA1. \*p-value < 0.05, \*\* p-value < 0.01, \*\*\* p-value < 0.001.

697

698  
699  
700  
701  
702

**Figure A10. Heatmap with Hierarchical clustering combining RNA-seq expression of patients on TCGA-KIRC and ICGC-RECA.** Columns are genes of the mRMR signature. Rows indicate RNA-seq expression of 590 patients of TCGA-KIRC and ICGC-RECA. Data of patients with distant metastasis that cannot be assessed (MX) were removed in order to clarify the clustering.

## 703 References

1. Hsieh, J.J.; Purdue, M.P.; Signoretti, S.; Swanton, C.; Albiges, L.; Schmidinger, M.; Heng, D.Y.; Larkin, J.; Ficarra, V. Renal Cell Carcinoma. *Nature reviews. Disease primers* **2017**, *3*, 17009, doi:10.1038/nrdp.2017.9.
2. Chen, L.; Xiang, Z.; Chen, X.; Zhu, X.; Peng, X. A Seven-Gene Signature Model Predicts Overall Survival in Kidney Renal Clear Cell Carcinoma. *Hereditas* **2020**, *157*, 38, doi:10.1186/s41065-020-00152-y.
3. Cui, H.; Shan, H.; Miao, M.Z.; Jiang, Z.; Meng, Y.; Chen, R.; Zhang, L.; Liu, Y. Identification of the Key Genes and Pathways Involved in the Tumorigenesis and Prognosis of Kidney Renal Clear Cell Carcinoma. *Scientific reports* **2020**, *10*, 1–10.
4. Society, A.C. Facts & Figures: 2020 Edition 2020.
5. Padala, S.A.; Barsouk, A.; Thandra, K.C.; Saginala, K.; Mohammed, A.; Vakiti, A.; Rawla, P.; Barsouk, A. Epidemiology of Renal Cell Carcinoma. *World journal of oncology* **2020**, *11*, 79–87, doi:10.14740/wjon1279.
6. Kann, B.H.; Hosny, A.; Aerts, H.J.W.L. Artificial Intelligence for Clinical Oncology. *Cancer cell* **2021**, *39*, 916–927, doi:10.1016/j.ccr.2021.04.002.
7. Chibon, F. Cancer Gene Expression Signatures - the Rise and Fall? *European journal of cancer* **2013**, *49*, 2000–2009, doi:10.1016/j.ejca.2013.02.021.

8. Zhan, Y.; Guo, W.; Zhang, Y.; Wang, Q.; Xu, X.-J.; Zhu, L. A Five-Gene Signature Predicts Prognosis in Patients with Kidney Renal Clear Cell Carcinoma. *Computational and mathematical methods in medicine* **2015**, *2015*, 842784, doi:10.1155/2015/842784.
9. Chang, P.; Bing, Z.; Tian, J.; Zhang, J.; Li, X.; Ge, L.; Ling, J.; Yang, K.; Li, Y. Comprehensive Assessment Gene Signatures for Clear Cell Renal Cell Carcinoma Prognosis. *Medicine* **2018**, *97*, e12679, doi:10.1097/MD.00000000000012679.
10. Chen, L.; Luo, Y.; Wang, G.; Qian, K.; Qian, G.; Wu, C.-L.; Dan, H.C.; Wang, X.; Xiao, Y. Prognostic Value of a Gene Signature in Clear Cell Renal Cell Carcinoma. *Journal of cellular physiology* **2019**, *234*, 10324–10335, doi:10.1002/jcp.27700.
11. Jiang, H.; Chen, H.; Chen, N. Construction and Validation of a Seven-Gene Signature for Predicting Overall Survival in Patients with Kidney Renal Clear Cell Carcinoma via an Integrated Bioinformatics Analysis. *Animal cells and systems* **2020**, *24*, 160–170, doi:10.1080/19768354.2020.1760932.
12. Pan, Q.; Wang, L.; Zhang, H.; Liang, C.; Li, B. Identification of a 5-Gene Signature Predicting Progression and Prognosis of Clear Cell Renal Cell Carcinoma. *Medical science monitor: international medical journal of experimental and clinical research* **2019**, *25*, 4401–4413, doi:10.12659/MSM.917399.
13. Wu, J.; Jin, S.; Gu, W.; Wan, F.; Zhang, H.; Shi, G.; Qu, Y.; Ye, D. Construction and Validation of a 9-Gene Signature for Predicting Prognosis in Stage III Clear Cell Renal Cell Carcinoma. *Frontiers in oncology* **2019**, *9*, 152, doi:10.3389/fonc.2019.00152.
14. Kalantzakos, T.J.; Sullivan, T.B.; Gloria, T.; Canes, D.; Moinzadeh, A.; Rieger-Christ, K.M. MiRNA-424-5p Suppresses Proliferation, Migration, and Invasion of Clear Cell Renal Cell Carcinoma and Attenuates Expression of O-GlcNAc-Transferase. *Cancers* **2021**, *13*, doi:10.3390/cancers13205160.
15. Wang, P.; Li, Y.; Reddy, C.K. Machine Learning for Survival Analysis: A Survey. *ACM Comput. Surv.* **2019**, *51*, doi:10.1145/3214306.
16. Cox, D.R. Regression Models and Life-Tables. *Journal of the Royal Statistical Society. Series B (Methodological)* **1972**, *34*, 187–220.
17. Wan, F.; Zhu, Y.; Han, C.; Xu, Q.; Wu, J.; Dai, B.; Zhang, H.; Shi, G.; Gu, W.; Ye, D. Identification and Validation of an Eight-Gene Expression Signature for Predicting High Fuhrman Grade Renal Cell Carcinoma. *International journal of cancer. Journal international du cancer* **2017**, *140*, 1199–1208, doi:10.1002/ijc.30535.
18. Hu, F.; Zeng, W.; Liu, X. A Gene Signature of Survival Prediction for Kidney Renal Cell Carcinoma by Multi-Omic Data Analysis. *International journal of molecular sciences* **2019**, *20*, doi:10.3390/ijms20225720.
19. Zou, Y.; Hu, C. A 14 Immune-Related Gene Signature Predicts Clinical Outcomes of Kidney Renal Clear Cell Carcinoma. *PeerJ* **2020**, *8*, e10183, doi:10.7717/peerj.10183.
20. Peng, H.; Long, F.; Ding, C. Feature Selection Based on Mutual Information Criteria of Max-Dependency, Max-Relevance, and Min-Redundancy. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **2005**, *27*, 1226–1238, doi:10.1109/TPAMI.2005.159.
21. Network, C.G.A.R. Comprehensive Molecular Characterization of Clear Cell Renal Cell Carcinoma. *Nature* **2013**, *499*, 43–49, doi:10.1038/nature12222.
22. GDC TCGA Kidney Clear Cell Carcinoma (KIRC) 2022.
23. Zhang, J.; Bajari, R.; Andric, D.; Gerthoffert, F.; Lepsa, A.; Nahal-Bose, H.; Stein, L.D.; Ferretti, V. The International Cancer Genome Consortium Data Portal. *Nature biotechnology* **2019**, *37*, 367–369.
24. Renal Cell Cancer - EU/FR (RECA) - Data Release 28 2022.

25. Gao, G.F.; Parker, J.S.; Reynolds, S.M.; Silva, T.C.; Wang, L.-B.; Zhou, W.; Akbani, R.; Bailey, M.; Balu, S.; Berman, B.P.; et al. Before and After: Comparison of Legacy and Harmonized TCGA Genomic Data Commons' Data. *Cell systems* **2019**, *9*, 24–34.e10, doi:10.1016/j.cels.2019.06.006.
26. Love, M.I.; Huber, W.; Anders, S. Moderated Estimation of Fold Change and Dispersion for RNA-Seq Data with DESeq2. *Genome biology* **2014**, *15*, 550, doi:10.1186/s13059-014-0550-8.
27. Kuhn, M. Building Predictive Models in R Using the Caret Package. *Journal of statistical software* **2008**, *28*, 1–26, doi:10.18637/jss.v028.i05.
28. Consortium, Gte. The GTEx Consortium Atlas of Genetic Regulatory Effects across Human Tissues. *Science* **2020**, *369*, 1318–1330, doi:10.1126/science.aaz1776.
29. Consortium, Gte. Human Genomics. The Genotype-Tissue Expression (GTEx) Pilot Analysis: Multitissue Gene Regulation in Humans. *Science* **2015**, *348*, 648–660, doi:10.1126/science.1262110.
30. Spooner, A.; Chen, E.; Sowmya, A.; Sachdev, P.; Kochan, N.A.; Trollor, J.; Brodaty, H. A Comparison of Machine Learning Methods for Survival Analysis of High-Dimensional Clinical Data for Dementia Prediction. *Scientific Reports* **2020**, *10*, 20410, doi:10.1038/s41598-020-77220-w.
31. Simon, N.; Friedman, J.; Hastie, T.; Tibshirani, R. Regularization Paths for Cox's Proportional Hazards Model via Coordinate Descent. *Journal of Statistical Software* **2011**, *39*, 1–13.
32. Ding, C.; Peng, H. Minimum Redundancy Feature Selection from Microarray Gene Expression Data. *Journal of bioinformatics and computational biology* **2005**, *3*, 185–205, doi:10.1142/s0219720005001004.
33. Jay, N.D.; De Jay, N.; Papillon-Cavanagh, S.; Olsen, C.; El-Hachem, N.; Bontempi, G.; Haibe-Kains, B. MRMRe: An R Package for Parallelized MRMR Ensemble Feature Selection. *Bioinformatics* **2013**, *29*, 2365–2368, doi:10.1093/bioinformatics/btt383.
34. Lang, M.; Binder, M.; Richter, J.; Schratz, P.; Pfisterer, F.; Coors, S.; Au, Q.; Casalicchio, G.; Kotthoff, L.; Bischl, B. Mlr3: A Modern Object-Oriented Machine Learning Framework in R. *Journal of Open Source Software* **2019**, *4*, 1903, doi:10.21105/joss.01903.
35. Wei, T.; Simko, V. R Package “Corrplot”: Visualization of a Correlation Matrix; 2021;
36. Kratzer, G.; Furrer, R. Varrank: An R Package for Variable Ranking Based on Mutual Information with Applications to Observed Systemic Datasets. *arXiv preprint arXiv:1804.07134* **2018**.
37. Blanche, P.; Kattan, M.W.; Gerds, T.A. The C-Index Is Not Proper for the Evaluation of t-Year Predicted Risks. *Biostatistics* **2019**, *20*, 347–357, doi:10.1093/biostatistics/kxy006.
38. Uno, H.; Cai, T.; Tian, L.; Wei, L.J. Evaluating Prediction Rules For-Year Survivors With Censored Regression Models. *Journal of the American Statistical Association* **2007**, *102*, 527–537, doi:10.1198/016214507000000149.
39. Potapov, S.; Adler, W.; Schmid, M. SurvAUC: Estimators of Prediction Accuracy for Time-to-Event Data.; 2012;
40. Kassambara, A.; Kosinski, M.; Biecek, P. Survminer: Drawing Survival Curves Using “Ggplot2”; 2021;
41. Piñero, J.; Ramírez-Anguita, J.M.; Saúch-Pitarch, J.; Ronzano, F.; Centeno, E.; Sanz, F.; Furlong, L.I. The DisGeNET Knowledge Platform for Disease Genomics: 2019 Update. *Nucleic acids research* **2020**, *48*, D845–D855, doi:10.1093/nar/gkz1021.
42. Chen, H.; Boutros, P.C. VennDiagram: A Package for the Generation of Highly-Customizable Venn and Euler Diagrams in R. *BMC bioinformatics* **2011**, *12*, 35, doi:10.1186/1471-2105-12-35.
43. Walter, W.; Sánchez-Cabo, F.; Ricote, M. GOplot: An R Package for Visually Combining Expression Data with Functional Analysis. *Bioinformatics* **2015**, *31*, 2912–2914, doi:10.1093/bioinformatics/btv300.
44. Lê, S.; Josse, J.; Husson, F. FactoMineR: An R Package for Multivariate Analysis. *Journal of Statistical Software* **2008**, *25*, doi:10.18637/jss.v025.i01.

45. Kassambara, A.; Mundt, F. *Factoextra: Extract and Visualize the Results of Multivariate Data Analyses*; 2020;
46. Therneau, T.M. *A Package for Survival Analysis in R*; 2022;
47. Patil, I. Visualizations with Statistical Details: The “*ggstatsplot*” Approach. *Journal of Open Source Software* **2021**, *6*, 3167, doi:10.21105/joss.03167.
48. Wu, T.; Hu, E.; Xu, S.; Chen, M.; Guo, P.; Dai, Z.; Feng, T.; Zhou, L.; Tang, W.; Zhan, L.; et al. ClusterProfiler 4.0: A Universal Enrichment Tool for Interpreting Omics Data. *The Innovation* **2021**, *2*, 100141, doi:10.1016/j.x-inn.2021.100141.
49. Harrison, E.; Drake, T.; Ots, R. *Finalfit: Quickly Create Elegant Regression Results Tables and Plots When Modelling*; 2022;
50. Kolde, R. *Pheatmap: Pretty Heatmaps*; 2019;
51. Chandrashekhar, D.S.; Bashel, B.; Balasubramanya, S.A.H.; Creighton, C.J.; Ponce-Rodriguez, I.; Chakravarthi, B.V.S.K.; Varambally, S. UALCAN: A Portal for Facilitating Tumor Subgroup Gene Expression and Survival Analyses. *Neoplasia (New York, N.Y.)* **2017**, *19*, 649–658.
52. Wan, B.; Liu, B.; Huang, Y.; Yu, G.; Lv, C. Prognostic Value of Immune-Related Genes in Clear Cell Renal Cell Carcinoma. *Aging* **2019**, *11*, 11474–11489, doi:10.18632/aging.102548.
53. Gao, X.; Yang, J.; Chen, Y. Identification of a Four Immune-Related Genes Signature Based on an Immunogenomic Landscape Analysis of Clear Cell Renal Cell Carcinoma. *Journal of cellular physiology* **2020**, *235*, 9834–9850, doi:10.1002/jcp.29796.
54. Zhang, Z.; Lin, E.; Zhuang, H.; Xie, L.; Feng, X.; Liu, J.; Yu, Y. Construction of a Novel Gene-Based Model for Prognosis Prediction of Clear Cell Renal Cell Carcinoma. *Cancer cell international* **2020**, *20*, 27, doi:10.1186/s12935-020-1113-6.
55. Kang, H.W.; Park, H.; Seo, S.P.; Byun, Y.J.; Piao, X.M.; Kim, S.M.; Kim, W.T.; Yun, S.J.; Jang, W.; Shon, H.S.; et al. Methylation Signature for Prediction of Progression Free Survival in Surgically Treated Clear Cell Renal Cell Carcinoma. *Journal of Korean medical science* **2019**, *34*, e144, doi:10.3346/jkms.2019.34.e144.
56. Jia, Z.; Wan, F.; Zhu, Y.; Shi, G.; Zhang, H.; Dai, B.; Ye, D. Forkhead-Box Series Expression Network Is Associated with Outcome of Clear-Cell Renal Cell Carcinoma. *Oncology letters* **2018**, *15*, 8669–8680, doi:10.3892/ol.2018.8405.
57. Zhu, P.; Piao, Y.; Dong, X.; Jin, Z. Forkhead Box J1 Expression Is Upregulated and Correlated with Prognosis in Patients with Clear Cell Renal Cell Carcinoma. *Oncology letters* **2015**, *10*, 1487–1494, doi:10.3892/ol.2015.3376.
58. Li, C.-S.; Chae, S.-C.; Lee, J.-H.; Zhang, Q.; Chung, H.-T. Identification of Single Nucleotide Polymorphisms in FOXJ1 and Their Association with Allergic Rhinitis. *Journal of human genetics* **2006**, *51*, 292–297, doi:10.1007/s10038-006-0359-8.
59. Li, C.-S.; Zhang, Q.; Lim, M.-K.; Sheen, D.-H.; Shim, S.-C.; Kim, J.-Y.; Lee, S.-S.; Yun, K.-J.; Moon, H.-B.; Chung, H.-T.; et al. Association of FOXJ1 Polymorphisms with Systemic Lupus Erythematosus and Rheumatoid Arthritis in Korean Population. *Experimental & Molecular Medicine* **2007**, *39*, 805–811, doi:10.1038/emm.2007.87.
60. Srivatsan, S.; Peng, S.L. Cutting Edge: Foxj1 Protects against Autoimmunity and Inhibits Thymocyte Egress. *Journal of immunology* **2005**, *175*, 7805–7809, doi:10.4049/jimmunol.175.12.7805.
61. Xian, S.; Shang, D.; Kong, G.; Tian, Y. FOXJ1 Promotes Bladder Cancer Cell Growth and Regulates Warburg Effect. *Biochemical and biophysical research communications* **2018**, *495*, 988–994, doi:10.1016/j.bbrc.2017.11.063.

62. Chen, H.-W.; Huang, X.-D.; Li, H.-C.; He, S.; Ni, R.-Z.; Chen, C.-H.; Peng, C.; Wu, G.; Wang, G.-H.; Wang, Y.-Y.; et al. Expression of FOXJ1 in Hepatocellular Carcinoma: Correlation with Patients' Prognosis and Tumor Cell Proliferation. *Molecular carcinogenesis* **2013**, *52*, 647–659, doi:10.1002/mc.21904.
63. Liu, K.; Fan, J.; Wu, J. Forkhead Box Protein J1 (FOXJ1) Is Overexpressed in Colorectal Cancer and Promotes Nuclear Translocation of B-Catenin in SW620 Cells. *Medical science monitor: international medical journal of experimental and clinical research* **2017**, *23*, 856–866, doi:10.12659/msm.902906.
64. Wang, J.; Cai, X.; Xia, L.; Zhou, J.; Xin, J.; Liu, M.; Shang, X.; Liu, J.; Li, X.; Chen, Z.; et al. Decreased Expression of FOXJ1 Is a Potential Prognostic Predictor for Progression and Poor Survival of Gastric Cancer. *Annals of surgical oncology* **2015**, *22*, 685–692, doi:10.1245/s10434-014-3742-2.
65. Abedalthagafi, M.S.; Wu, M.P.; Merrill, P.H.; Du, Z.; Woo, T.; Sheu, S.-H.; Hurwitz, S.; Ligon, K.L.; Santagata, S. Decreased FOXJ1 Expression and Its Ciliogenesis Programme in Aggressive Ependymoma and Choroid Plexus Tumours. *The Journal of pathology* **2016**, *238*, 584–597, doi:10.1002/path.4682.
66. Lin, B.M.; Nadkarni, G.N.; Tao, R.; Graff, M.; Fornage, M.; Buyske, S.; Matise, T.C.; Highland, H.M.; Wilkens, L.R.; Carlson, C.S.; et al. Genetics of Chronic Kidney Disease Stages Across Ancestries: The PAGE Study. *Frontiers in genetics* **2019**, *10*, 494, doi:10.3389/fgene.2019.00494.
67. Shirota, H.; Klinman, D.M.; Ito, S.-E.; Ito, H.; Kubo, M.; Ishioka, C. IL4 from T Follicular Helper Cells Down-regulates Antitumor Immunity. *Cancer Immunology Research* **2017**, *5*, 61–71, doi:10.1158/2326-6066.cir-16-0113.
68. Ito, S.-E.; Shirota, H.; Kasahara, Y.; Saijo, K.; Ishioka, C. IL-4 Blockade Alters the Tumor Microenvironment and Augments the Response to Cancer Immunotherapy in a Mouse Model. *Cancer Immunology, Immunotherapy* **2017**, *66*, 1485–1496, doi:10.1007/s00262-017-2043-6.
69. Jia, Y.; Xie, X.; Shi, X.; Li, S. Associations of Common IL-4 Gene Polymorphisms with Cancer Risk: A Meta-Analysis. *Molecular Medicine Reports* **2017**, *16*, 1927–1945, doi:10.3892/mmr.2017.6822.
70. Cheng, H.; Borczuk, A.; Janakiram, M.; Ren, X.; Lin, J.; Assal, A.; Halmos, B.; Perez-Soler, R.; Zang, X. Wide Expression and Significance of Alternative Immune Checkpoint Molecules, B7x and HHLA2, in PD-L1-Negative Human Lung Cancers. *Clinical Cancer Research* **2018**, *24*, 1954–1964, doi:10.1158/1078-0432.ccr-17-2924.
71. Zhao, R.; Chinai, J.M.; Buhl, S.; Scandiuzzi, L.; Ray, A.; Jeon, H.; Ohaegbulam, K.C.; Ghosh, K.; Zhao, A.; Scharff, M.D.; et al. HHLA2 Is a Member of the B7 Family and Inhibits Human CD4 and CD8 T-Cell Function. *Proceedings of the National Academy of Sciences of the United States of America* **2013**, *110*, 9879–9884, doi:10.1073/pnas.1303524110.
72. Byun, J.M.; Cho, H.J.; Park, H.Y.; Lee, D.S.; Choi, I.H.; Kim, Y.N.; Jeong, C.H.; Kim, D.H.; Hwa Im, D.; Min, B.J.; et al. The Clinical Significance of HERV-H LTR -Associating 2 Expression in Cervical Adenocarcinoma. *Medicine* **2021**, *100*, e23691, doi:10.1097/MD.00000000000023691.
73. Boor, P.P.C.; Sideras, K.; Biermann, K.; Hosein Aziz, M.; Levink, I.J.M.; Mancham, S.; Erler, N.S.; Tang, X.; van Eijck, C.H.; Bruno, M.J.; et al. HHLA2 Is Expressed in Pancreatic and Ampullary Cancers and Increased Expression Is Associated with Better Post-Surgical Prognosis. *British journal of cancer* **2020**, *122*, 1211–1218, doi:10.1038/s41416-020-0755-4.
74. Cheng, H.; Janakiram, M.; Borczuk, A.; Lin, J.; Qiu, W.; Liu, H.; Chinai, J.M.; Halmos, B.; Perez-Soler, R.; Zang, X. HHLA2, a New Immune Checkpoint Member of the B7 Family, Is Widely Expressed in Human Lung Cancer and Associated with EGFR Mutational Status. *Clinical cancer research: an official journal of the American Association for Cancer Research* **2017**, *23*, 825–832, doi:10.1158/1078-0432.CCR-15-3071.
75. Shimonosono, M.; Arigami, T.; Yanagita, S.; Matsushita, D.; Uchikado, Y.; Kijima, Y.; Kurahara, H.; Kita, Y.; Mori, S.; Sasaki, K.; et al. The Association of Human Endogenous Retrovirus-H Long Terminal Repeat-Associ-

- ating Protein 2 (HHLA2) Expression with Gastric Cancer Prognosis. *Oncotarget* **2018**, *9*, 22069–22078, doi:10.18632/oncotarget.25179.
76. Chen, L.; Zhu, D.; Feng, J.; Zhou, Y.; Wang, Q.; Feng, H.; Zhang, J.; Jiang, J. Overexpression of HHLA2 in Human Clear Cell Renal Cell Carcinoma Is Significantly Associated with Poor Survival of the Patients. *Cancer Cell International* **2019**, *19*, doi:10.1186/s12935-019-0813-2.
77. Reidy, K.; Tufro, A. Semaphorins in Kidney Development and Disease: Modulators of Ureteric Bud Branching, Vascular Morphogenesis, and Podocyte-Endothelial Crosstalk. *Pediatric nephrology* **2011**, *26*, 1407–1412, doi:10.1007/s00467-011-1769-1.
78. Xia, J.; Worzfeld, T. Semaphorins and Plexins in Kidney Disease. *Nephron* **2016**, *132*, 93–100, doi:10.1159/000443645.
79. Neufeld, G.; Mumblat, Y.; Smolkin, T.; Toledano, S.; Nir-Zvi, I.; Ziv, K.; Kessler, O. The Role of the Semaphorins in Cancer. *Cell adhesion & migration* **2016**, *10*, 652–674, doi:10.1080/19336918.2016.1197478.
80. Karayan-Tapon, L.; Wager, M.; Guilhot, J.; Levillain, P.; Marquant, C.; Clarhaut, J.; Potiron, V.; Roche, J. Semaphorin, Neuropilin and VEGF Expression in Glial Tumours: SEMA3G, a Prognostic Marker? *British journal of cancer* **2008**, *99*, 1153–1160, doi:10.1038/sj.bjc.6604641.
81. Wu, H.; Malone, A.F.; Donnelly, E.L.; Kirita, Y.; Uchimura, K.; Ramakrishnan, S.M.; Gaut, J.P.; Humphreys, B.D. Single-Cell Transcriptomics of a Human Kidney Allograft Biopsy Specimen Defines a Diverse Inflammatory Response. *Journal of the American Society of Nephrology: JASN* **2018**, *29*, 2069–2080, doi:10.1681/ASN.2018020125.
82. Liang, J.; Liu, Z.; Zou, Z.; Tang, Y.; Zhou, C.; Yang, J.; Wei, X.; Lu, Y. The Correlation Between the Immune and Epithelial-Mesenchymal Transition Signatures Suggests Potential Therapeutic Targets and Prognosis Prediction Approaches in Kidney Cancer. *Scientific Reports* **2018**, *8*, doi:10.1038/s41598-018-25002-w.
83. Balk, S.P.; Knudsen, K.E. AR, the Cell Cycle, and Prostate Cancer. *Nuclear Receptor Signaling* **2008**, *6*, nrs.06001, doi:10.1621/nrs.06001.
84. Sun, M.; Abdollah, F. Re: AR-V7 and Resistance to Enzalutamide and Abiraterone in Prostate Cancer. *European Urology* **2015**, *68*, 162–163, doi:10.1016/j.eururo.2015.03.054.
85. Huang, Q.; Sun, Y.; Zhai, W.; Ma, X.; Shen, D.; Du, S.; You, B.; Niu, Y.; Huang, C.-P.; Zhang, X.; et al. Androgen Receptor Modulates Metastatic Routes of VHL Wild-Type Clear Cell Renal Cell Carcinoma in an Oxygen-Dependent Manner. *Oncogene* **2020**, *39*, 6677–6691, doi:10.1038/s41388-020-01455-0.
86. Chen, Y.; Sun, Y.; Rao, Q.; Xu, H.; Li, L.; Chang, C. Androgen Receptor (AR) Suppresses MiRNA-145 to Promote Renal Cell Carcinoma (RCC) Progression Independent of VHL Status. *Oncotarget* **2015**, *6*, 31203–31215, doi:10.18632/oncotarget.4522.
87. Lee, K.-H.; Kim, B.-C.; Jeong, S.-H.; Jeong, C.W.; Ku, J.H.; Kwak, C.; Kim, H.H. Histone Demethylase LSD1 Regulates Kidney Cancer Progression by Modulating Androgen Receptor Activity. *International journal of molecular sciences* **2020**, *21*, doi:10.3390/ijms21176089.
88. Wang, K.; Sun, Y.; Tao, W.; Fei, X.; Chang, C. Androgen Receptor (AR) Promotes Clear Cell Renal Cell Carcinoma (CcRCC) Migration and Invasion via Altering the CircHIAT1/MiR-195-5p/29a-3p/29c-3p/CDC42 Signals. *Cancer letters* **2017**, *394*, 1–12, doi:10.1016/j.canlet.2016.12.036.
89. You, B.; Sun, Y.; Luo, J.; Wang, K.; Liu, Q.; Fang, R.; Liu, B.; Chou, F.; Wang, R.; Meng, J.; et al. Androgen Receptor Promotes Renal Cell Carcinoma (RCC) Vasculogenic Mimicry (VM) via Altering TWIST1 Nonsense-Mediated Decay through LncRNA-TANAR. *Oncogene* **2021**, *40*, 1674–1689, doi:10.1038/s41388-020-01616-1.

90. Larsen, K.B.; Lutterodt, M.C.; Møllgård, K.; Møller, M. Expression of the Homeobox Genes OTX2 and OTX1 in the Early Developing Human Brain. *The journal of histochemistry and cytochemistry: official journal of the Histochemistry Society* **2010**, *58*, 669–678, doi:10.1369/jhc.2010.955757.
91. García-Frigola, C.; Carreres, M.I.; Vigar, C.; Mason, C.; Herrera, E. Zic2 Promotes Axonal Divergence at the Optic Chiasm Midline by EphB1-Dependent and -Independent Mechanisms. *Development* **2008**, *135*, 1833–1841, doi:10.1242/dev.020693.
92. Grinberg, I.; Millen, K.J. The ZIC Gene Family in Development and Disease. *Clinical genetics* **2005**, *67*, 290–296, doi:10.1111/j.1399-0004.2005.00418.x.
93. Marchini, S.; Poynor, E.; Barakat, R.R.; Clivio, L.; Cinquini, M.; Fruscio, R.; Porcu, L.; Bussani, C.; D’Incalci, M.; Erba, E.; et al. The Zinc Finger Gene ZIC2 Has Features of an Oncogene and Its Overexpression Correlates Strongly with the Clinical Course of Epithelial Ovarian Cancer. *Clinical cancer research: an official journal of the American Association for Cancer Research* **2012**, *18*, 4313–4324, doi:10.1158/1078-0432.CCR-12-0037.
94. Liu, Z.-H.; Chen, M.-L.; Zhang, Q.; Zhang, Y.; An, X.; Luo, Y.-L.; Liu, X.-M.; Liu, S.-X.; Liu, Q.; Yang, T.; et al. ZIC2 Is Downregulated and Represses Tumor Growth via the Regulation of STAT3 in Breast Cancer. *International journal of cancer. Journal international du cancer* **2020**, *147*, 505–518, doi:10.1002/ijc.32922.
95. Wu, C.-Y.; Li, L.; Chen, S.-L.; Yang, X.; Zhang, C.Z.; Cao, Y. A Zic2/Runx2/NOLC1 Signaling Axis Mediates Tumor Growth and Metastasis in Clear Cell Renal Cell Carcinoma. *Cell death & disease* **2021**, *12*, 319, doi:10.1038/s41419-021-03617-8.
96. Lin, Y.-H.; Zhen, Y.-Y.; Chien, K.-Y.; Lee, I.-C.; Lin, W.-C.; Chen, M.-Y.; Pai, L.-M. LIMCH1 Regulates Non-muscle Myosin-II Activity and Suppresses Cell Migration. *Molecular biology of the cell* **2017**, *28*, 1054–1065, doi:10.1091/mbc.E15-04-0218.
97. Karlsson, T.; Kvärnström, S.; Holmlund, C.; Botling, J.; Micke, P.; Henriksson, R.; Johansson, M.; Hedman, H. LMO7 and LIMCH1 Interact with LRIG Proteins in Lung Cancer, with Prognostic Implications for Early-Stage Disease. *Lung cancer* **2018**, *125*, 174–184, doi:10.1016/j.lungcan.2018.09.017.
98. Cizkova, M.; Cizeron-Clairac, G.; Vacher, S.; Susini, A.; Andrieu, C.; Lidereau, R.; Bièche, I. Gene Expression Profiling Reveals New Aspects of PIK3CA Mutation in ERalpha-Positive Breast Cancer: Major Implication of the Wnt Signaling Pathway. *PloS one* **2010**, *5*, e15647, doi:10.1371/journal.pone.0015647.
99. Halle, M.K.; Sødal, M.; Forsse, D.; Engerud, H.; Woie, K.; Lura, N.G.; Wagner-Larsen, K.S.; Trovik, J.; Bertelsen, B.I.; Haldorsen, I.S.; et al. A 10-Gene Prognostic Signature Points to LIMCH1 and HLA-DQB1 as Important Players in Aggressive Cervical Cancer Disease. *British journal of cancer* **2021**, *124*, 1690–1698, doi:10.1038/s41416-021-01305-0.
100. Uhlen, M.; Zhang, C.; Lee, S.; Sjöstedt, E.; Fagerberg, L.; Bidkhor, G.; Benfeitas, R.; Arif, M.; Liu, Z.; Edfors, F.; et al. A Pathology Atlas of the Human Cancer Transcriptome. *Science* **2017**, *357*, doi:10.1126/science.aan2507.
101. Expression of LIMCH1 in Renal Cancer - Interactive Survival Scatter Plot - The Human Protein Atlas 2022.
102. Clark, B.D.; Kwon, E.; Maffie, J.; Jeong, H.-Y.; Nadal, M.; Strop, P.; Rudy, B. DPP6 Localization in Brain Supports Function as a Kv4 Channel Associated Protein. *Frontiers in molecular neuroscience* **2008**, *1*, 8, doi:10.3389/neuro.02.008.2008.
103. Zhao, X.; Cao, D.; Ren, Z.; Liu, Z.; Lv, S.; Zhu, J.; Li, L.; Lang, R.; He, Q. Dipeptidyl Peptidase like 6 Promoter Methylation Is a Potential Prognostic Biomarker for Pancreatic Ductal Adenocarcinoma. *Bioscience reports* **2020**, *40*, doi:10.1042/BSR20200214.

104. Choy, T.-K.; Wang, C.-Y.; Phan, N.N.; Khoa Ta, H.D.; Anuraga, G.; Liu, Y.-H.; Wu, Y.-F.; Lee, K.-H.; Chuang, J.-Y.; Kao, T.-J. Identification of Dipeptidyl Peptidase (DPP) Family Genes in Clinical Breast Cancer Patients via an Integrated Bioinformatics Approach. *Diagnostics (Basel, Switzerland)* **2021**, *11*, doi:10.3390/diagnostics11071204.
105. Wang, Y.; Zhang, S. Quantitative Assessment of the Association between GNB3 C825T Polymorphism and Cancer Risk. *Journal of B.U.ON.: official journal of the Balkan Union of Oncology* **2014**, *19*, 1092–1095.
106. Fingas, C.D.; Katsounas, A.; Kahraman, A.; Siffert, W.; Jochum, C.; Gerken, G.; Nückel, H.; Canbay, A. Prognostic Assessment of Three Single-Nucleotide Polymorphisms (GNB3 825C>T, BCL2-938C>A, MCL1-386C>G) in Extrahepatic Cholangiocarcinoma. *Cancer investigation* **2010**, *28*, 472–478, doi:10.3109/07357900903095714.
107. Paleari, R.G.; Peres, R.M.R.; Florentino, J.O.; Heinrich, J.K.; Bragance, W.O.; Del Valle, J.C.T.; Zeferino, L.C.; Derchain, S.F.M.; Sarian, L.O. Reduced Prevalence of the C825T Polymorphism of the G-Protein Beta Subunit Gene in Women with Breast Cancer. *The International journal of biological markers* **2011**, *26*, 234–240, doi:10.5301/JBM.2011.8751.
108. Santo, C.D.; De Santo, C.; Arscott, R.; Booth, S.; Karydis, I.; Jones, M.; Asher, R.; Salio, M.; Middleton, M.; Cerundolo, V. Invariant NKT Cells Modulate the Suppressive Activity of IL-10-Secreting Neutrophils Differentiated with Serum Amyloid A. *Nature Immunology* **2010**, *11*, 1039–1046, doi:10.1038/ni.1942.
109. Paret, C.; Schön, Z.; Szponar, A.; Kovacs, G. Inflammatory Protein Serum Amyloid A1 Marks a Subset of Conventional Renal Cell Carcinomas with Fatal Outcome. *European Urology* **2010**, *57*, 859–866, doi:10.1016/j.eururo.2009.08.014.
110. Expression of SAA1 in Renal Cancer - Interactive Survival Scatter Plot - The Human Protein Atlas 2022.
111. Marshall, F.F. Serum Protein Profiling by SELDI Mass Spectrometry: Detection of Multiple Variants of Serum Amyloid Alpha in Renal Cancer Patients. *The Journal of urology* **2005**, *173*, 1919–1920.
112. Guo, R.; Zou, B.; Liang, Y.; Bian, J.; Xu, J.; Zhou, Q.; Zhang, C.; Chen, T.; Yang, M.; Wang, H.; et al. LncRNA RCAT1 Promotes Tumor Progression and Metastasis via MiR-214-5p/E2F2 Axis in Renal Cell Carcinoma. *Cell death & disease* **2021**, *12*, 689, doi:10.1038/s41419-021-03955-7.
113. Qi, N.; Chen, Y.; Gong, K.; Li, H. Concurrent Renal Cell Carcinoma and Urothelial Carcinoma: Long-Term Follow-up Study of 27 Cases. *World journal of surgical oncology* **2018**, *16*, 16, doi:10.1186/s12957-018-1321-x.
114. Knez, V.M.; Barrow, W.; Lucia, M.S.; Wilson, S.; La Rosa, F.G. Clear Cell Urothelial Carcinoma of the Urinary Bladder: A Case Report and Review of the Literature. *Journal of medical case reports* **2014**, *8*, 275, doi:10.1186/1752-1947-8-275.
115. Rotellini, M.; Fondi, C.; Paglierani, M.; Stomaci, N.; Raspollini, M.R. Clear Cell Carcinoma of the Bladder in a Patient with a Earlier Clear Cell Renal Cell Carcinoma: A Case Report with Morphologic, Immunohistochemical, and Cytogenetical Analysis. *Applied immunohistochemistry & molecular morphology: AIMM / official publication of the Society for Applied Immunohistochemistry* **2010**, *18*, 396–399, doi:10.1097/PAI.0b013e3181d57dce.
116. van de Pol, J.A.A.; van den Brandt, P.A.; Schouten, L.J. Kidney Stones and the Risk of Renal Cell Carcinoma and Upper Tract Urothelial Carcinoma: The Netherlands Cohort Study. *British journal of cancer* **2018**, *120*, 368–374, doi:10.1038/s41416-018-0356-7.
117. Ha, M.J.; Baladandayuthapani, V.; Do, K.-A. Prognostic Gene Signature Identification Using Causal Structure Learning: Applications in Kidney Cancer. *Cancer informatics* **2015**, *14*, 23–35, doi:10.4137/CIN.S14873.
118. Chen, Y.-L.; Ge, G.-J.; Qi, C.; Wang, H.; Wang, H.-L.; Li, L.-Y.; Li, G.-H.; Xia, L.-Q. A Five-Gene Signature May Predict Sunitinib Sensitivity and Serve as Prognostic Biomarkers for Renal Cell Carcinoma. *Journal of cellular physiology* **2018**, *233*, 6649–6660, doi:10.1002/jcp.26441.

119. Jafari, M.; Guan, Y.; Wedge, D.C.; Ansari-Pour, N. Re-Evaluating Experimental Validation in the Big Data Era: A Conceptual Argument. *Genome biology* **2021**, *22*, 71, doi:10.1186/s13059-021-02292-4.
120. The Cancer Genome Atlas Program **2022**.
121. TCGA/GDC Data Portal - Data Release 18.0 **2019**.
122. Goldman, M.J.; Craft, B.; Hastie, M.; Repečka, K.; McDade, F.; Kamath, A.; Banerjee, A.; Luo, Y.; Rogers, D.; Brooks, A.N.; et al. Visualizing and Interpreting Cancer Genomics Data via the Xena Platform. *Nature Biotechnology* **2020**, *38*, 675–678, doi:10.1038/s41587-020-0546-8.
123. Dai, J.; Lu, Y.; Wang, J.; Yang, L.; Han, Y.; Wang, Y.; Yan, D.; Ruan, Q.; Wang, S. A Four-Gene Signature Predicts Survival in Clear-Cell Renal-Cell Carcinoma. *Oncotarget* **2016**, *7*, 82712–82726, doi:10.18632/oncotarget.12631.
124. D’Costa, N.M.; Cina, D.; Shrestha, R.; Bell, R.H.; Lin, Y.-Y.; Asghari, H.; Monjaras-Avila, C.U.; Kollmannsberger, C.; Hach, F.; Chavez-Munoz, C.I.; et al. Identification of Gene Signature for Treatment Response to Guide Precision Oncology in Clear-Cell Renal Cell Carcinoma. *Scientific reports* **2020**, *10*, 2026, doi:10.1038/s41598-020-58804-y.