

# **Prognostic and biological relevance of collagen biosynthesis pathway in prostate cancer**

**Supplementary material, transcriptome part**

2024-01-01

## Supplementary Methods

### Software

The analysis was done with R version 4.2.3 (R Foundation).

Tabular data were handled with the packages *tidyverse* (1), *rlang* (2) and *trafo*. Text data were handled with *stringi* (3).

Import of the TCGA and DKFZ data sets from the cBioportal repository was accomplished with in-house-developed R scripts. Transcriptome datasets from the Gene Expression Omnibus were fetched with the *GEOquery* package (4). Gene and probe annotation was accomplished with the *AnnotationDbi* (5) and *org.Hs.eg.db* packages (6). For prediction of non-malignant cell counts and fractions in cancer samples, the R implementation of the *MCP Counter* and *xCell* algorithms provided by the *immunedecconv* package was used (7–9). Single sample gene set enrichment analysis scores (ssGSEA) were computed with the *GSVA* algorithm (10) implemented by the package *gseaTools*.

Semi-supervised clustering was done with the package *clustTools* employing algorithms from the *philentropy*, *cluster*, *factoextra* and *umap* packages (11–14).

For statistical hypothesis testing, effect size calculation, gene set variation analysis (GSVA) and differential gene expression analysis, the packages *rstatix* (15), *ExDA* and *microViz* were employed. Biological process gene ontology (GO) enrichment analysis was performed with the *goana* algorithm implemented by the *limma* (16) and *microViz* packages. Clustering of GO terms by their semantic similarity was accomplished with *GOSemSim* (17), *microViz* and *clustTools*. Differential modulation of transcriptional regulons and signaling pathways was analyzed with the *collectRI* (18) and *PROGENy* (19) databases by linear modeling tools of the *decoupleR* package (20). Metabolic reaction activity modeling and metabolic subsystem enrichment analysis were performed with the packages *BiGGR* (21,22) and *biggrExtra*.

For multi-parameter Elastic Net Cox modeling, the packages *glmnet* (23), *survival* (24), *survminer* (25) and *coxExtensions* were used.

For visualization of the results, the packages *ggplot2* (1) (scatter, stack, bubble and bar plots, heat maps), *ExDA* (violin, stack and ribbon plots), *survminer* and *coxExtensions* (Kaplan-Meier plots), and *microViz* (heat maps) were used. Figures and tables were created with the packages *cowplot* (26) and *flextable* (27).

## Data sources and data import

Transcriptome cohorts for analysis in the current report were selected from studies deposited at cBioportal and Gene Expression Omnibus (GEO) with the following criteria: availability of information on Gleason scoring, biochemical relapse and biochemical relapse-free survival, and expression data for 55 collagen-related genes investigated also in the proteomic part of the project (**Supplementary Table S1** and **S2**).

The TCGA prostate cancer data set (28,29) and DKFZ data set (30) consisted of normalized RNA sequencing data for 493 and 118 cancer samples, respectively, with accompanying clinical information and were obtained from the cBioportal repository with in-house-developed R scripts. The TCGA cohort included also 52 donor-matched samples of the benign prostate and prostate cancer tissue.

The GSE54460 (31) (RNA sequencing, n = 106 cancer samples), GSE70768 (32) (Illumina HumanHT-12 V4.0 expression beadchip microarray, n = 125), GSE70769 (32) (Illumina HumanHT-12 V4.0 expression beadchip microarray, n = 94), and GSE220095 (33) (RNA sequencing, n = 176) data sets included normalized whole-transcriptome and basic clinical information and were fetched from GEO with the `getGEO()` function from the *GEOquery* package. The GSE70768 data set (32) included also 73 donor-matched normal prostate and cancer tissue specimens. Demographic and clinical characteristic of the investigated transcriptome data sets is provided in **Supplementary Table S1**.

Gene expression levels were transformed with the  $\log_2(\text{expression})$  and  $\log_2(\text{count}+1)$  function for the microarray and RNA sequencing data sets prior to further analyses. The *MCP Counter* and *xCell* estimates of non-malignant cell content in cancer samples were computed with the `deconvolute()` function from the *immunedconv* package. Gene signatures of the [Reactome pathways](#) were obtained from the [MSig database](#), version 7.5.1. ssGSEA scores were computed with the `calculate()` function from the *gseaTools* package.

The investigated collagen-related genes included the published Collagen Signature (34) and the Reactome Collagen formation pathway genes (R-HSA-1474290) and were constrained to features covered by proteomic analyses of the current report (**Supplementary Table S2**).

## Statistical hypothesis testing, effect size and statistical significance

Differences in numeric variables between two groups were investigated by paired and unpaired two-tailed T tests with Cohen's d effect size statistic or Mann-Whitney tests with r effect size statistic. Differences in numeric variables between three or more groups were assessed by one-way ANOVA with  $\eta^2$  effect size statistic. Differences in distribution of categorical variables were assessed by  $\chi^2$  test with Cramer's V effect size statistic. Odds

ratio served as an effect size statistic in enrichment analyses. Effect sizes were interpreted as follows (35):

- Cohen's d, weak: 0.2 - 0.5, moderate: 0.5 - 0.8, large:  $\geq 0.8$
- r, weak: 0.1 - 0.3, moderate: 0.3 - 0.5, large:  $\geq 0.5$
- $\eta^2$ : weak: 0.02 - 0.13, moderate: 0.13 - 0.26, large:  $\geq 0.26$
- Cramer's V, weak: 0.1 - 0.3, moderate: 0.3 - 0.5, large:  $\geq 0.5$
- OR, weak: 1.44 - 2.48, moderate: 2.48 - 4.27, large:  $\geq 4.27$

P values were corrected for multiple testing with the false discovery rate method (FDR) (36) within each analysis step and cohort. Effects were considered statistically significant for FDR-corrected p values  $< 0.05$ .

## Comparison of gene expression between the benign and cancer tissue, and cancer samples stratified by Gleason scores

Differences in  $\log_2$ -transformed expression levels of the collagen pathway genes between donor-matched cancer and benign prostate specimens were investigated with paired two-tailed T test with Cohen's effect size statistic in the GSE70768 and TCGA cohorts (**Figure 1B, Supplementary Figure S1 and Supplementary Table S3**). Differences in gene expression between cancer samples stratified by Gleason score (5 - 6, 7,  $\geq 8$ ) were explored by one-way ANOVA with  $\eta^2$  effect size statistic (**Figure 1C, Supplementary Table S4**). Differentially regulated genes were defined by pFDR  $< 0.05$  and at least weak effect size of the difference ( $d \geq 0.2$  or  $\eta^2 \geq 0.02$ ). The analyses were done with the function `compare_variables()` from the *ExDA* package.

## Collagen clusters of prostate cancer samples

Prostate cancer samples in the TCGA cohort were clustered by normalized  $\log_2$ -transformed expression levels of the collagen-related genes of interest with the PAM (partition around medoids) algorithm with cosine distance between the samples (function `kcluster()`, package *clustTools*). The choice of the clustering algorithm was motivated by the following criteria (**Supplementary Figure S2A**):

- cluster separation measured by mean silhouette statistic (37)
- fraction of potentially misclassified observations (observations with negative silhouette widths) (37)

- explanatory performance gauged by fraction of explained clustering variance (ratio of the total between-cluster sum of squares to the total sum of squares)
- preservation of the nearest neighborhood (mean fraction of the five nearest neighbors placed in the same cluster) (38)
- cluster re-assignment accuracy in 5-fold cross-validation (39,40). Computation of those numeric statistics of clustering quality in a cross-validation setting was done with the methods `summary()` and `cv()` provided by the *clustTools* package. The cluster number choice was based on the bend of the peak of the average silhouette statistic (37) (method `plot()`, package *clustTools*, **Supplementary Figure S2B**). By this means, two clusters of cancer samples were identified: ‘collagen high’ and ‘collagen low’ cancers.

Assignment of cancers samples of the GSE54460, GSE70768, GSE70769, GSE220095, and DKFZ cohorts to the collagen clusters based on normalized  $\log_2$ -transformed expression levels of the collagen-related genes was done with an inverse distance weighted k-nearest neighbor classifier (method `predict()`, package *clustTools*) (40). The quality of the cluster assignment was assessed by comparing the collagen cluster distribution and numeric statistics of clustering quality described above the training TCGA data set and the test cohorts (**Supplementary Figure S2CD**). In addition, separation of the clusters was assessed by a visual inspection of the UMAP layout plots (method `plot()`, package *clustTools*) and heat maps of the mean levels of the clustering factors in the collagen clusters (**Supplementary Figure S3**).

Differences in the  $\log_2$ -transformed collagen pathway gene expression between the collagen clusters were investigated by two-tailed T test with Cohen’s d effect size statistic (**Supplementary Table S5**).

## Clinical characteristic of the collagen clusters

Differences in numeric clinical variables (age and PSA) between the collagen clusters were assessed by Mann-Whitney test with r effect size statistic (`test_two_groups()`, *microViz*). Differences in qualitative clinical variables between the collagen clusters were investigated by  $\chi^2$  test with Cramer’s V effect size statistic (function `compare_variables()`, package *ExDA*). The demographic and clinical characteristic of the collagen clusters is presented in **Supplementary Figure S4** and **Supplementary Table S6**.

## Non-malignant cell infiltration in the collagen clusters

Non-malignant cell counts and non-malignant cell fractions in the cancer samples were predicted by the *MCP Counter* (9) and *xCell* (8) algorithms, respectively. Differences in the predicted infiltration levels between the collagen clusters were investigated by Mann-

Whitney test with r effect size statistic (function `compare_variables()`, package *ExDA*). Populations differing between the clusters with pFDR < 0.05 and at least weak effect size ( $r \geq 0.1$ ) were considered biologically relevant. The analysis results are shown in **Supplementary Figure S5** and **Supplementary Table S7**.

## Gene set variance analysis of Reactome pathway gene signatures

ssGSEA scores (10) of the Reactome pathway gene signatures were compared between the collagen clusters by two-tailed T test with Cohen's d effect size statistic. Signatures found to be significantly regulated (pFDR < 0.05) with at least weak effect size ( $d \geq 0.2$ ) in at least five cohorts were further investigated. To get additional functional insight and for purposes of visualization, these common regulated signatures were classified by their co-expression patterns in the TCGA data set by unsupervised KMEANS clustering (function `kcluster()`, package *clustTools*). The analysis results are presented in **Supplementary Figure S6** and **Supplementary Table S8**.

## Differential gene expression and GO enrichment analysis

Differences in whole-transcriptome gene expression between the collagen clusters were investigated by FDR-corrected two-tailed T test with Cohen's d effect size statistic (function `test_two_groups()`, package *microViz*). Genes significantly differentially expressed between the clusters were identified by the significance cutoff (pFDR < 0.05) and at least weak effect size of the expression differences ( $d \geq 0.2$ ). Genes found to be differentially regulated between the clusters are listed in **Supplementary Table S9**.

Biological process GO enrichment within genes upregulated in each of the collagen high and collagen low clusters was analyzed with the *limma* package *goana* algorithm (16), re-implemented by the development package *microViz* (function `GOana()`). Enrichment p values were adjusted for multiple testing with the FDR method. Odds ratio (OR) calculated for a GO term with the following formula served as a measure of effect size:

$$OR_{GO} = \frac{N_{GO \cap DGE} \times N_{total}}{N_{GO \cap total} \times N_{DGE}}$$

where  $N_{GO \cap DGE}$  stands for the number of differentially regulated genes assigned to the GO term,  $N_{total}$  is the total number of investigated genes,  $N_{GO \cap total}$  is the number of investigated genes assigned to the GO term, and  $N_{DGE}$  stands for the number of differentially regulated genes. GO terms found to be significantly enriched with at least weak effect size (pFDR < 0.05,  $OR \geq 1.44$ ) in at least five cohorts were further investigated. For additional functional insight, such common enriched GO terms were subjected to multi-dimensional scaling and unsupervised KMEANS clustering with pairwise Wang distance as a measure of semantic similarity (functions `go_sem()` and `kcluster()`, packages

*GOSemSim*, *microViz* and *clustTools*) (17). Results of the GO enrichment analysis are shown in **Supplementary Figure S6** and **Supplementary Table S10**.

## Transcriptional regulons and signaling pathway activity in the collagen clusters

Modulation of transcriptional regulons, i.e. sets of genes whose expression is controlled by a common transcriptional factor, and activity of selected signaling pathways in collagen high cancers as compared with collagen low tumors were explored by linear modeling with the *collectRI* (18) and *PROGENy* (19) databases, respectively. Linear models were constructed for whole-transcriptome effect sizes of gene regulation (Cohen's d, see: **Differential gene expression and GO enrichment analysis**) with the `run_ulm()` and `run_mlm()` tools from the *decoupleR* package (20). In each case, linear model score (LM score) served as a metric of activity magnitude. P values (LM score  $\neq 0$ ) were corrected for multiple testing with the FDR method. Regulons and signaling pathways found to be significantly activated or inhibited in at least five cohorts were further discussed. The analysis results are presented in **Supplementary Figure S7**, and **Supplementary Tables S11** and **S12**.

## Activity of metabolic reactions in the collagen clusters

Rules of assignment of genes to biochemical reactions were retrieved from the Recon2 human metabolism model available via the BiGG database (22) and the R package *BiGGR* (21). Estimates and standard errors of  $\log_2$  fold-regulation of expression for all available genes between the collagen clusters were calculated by two-tailed T test as described in **Differential gene expression and GO enrichment analysis**. Estimates of biochemical pathway fold-regulation were computed by evaluation of the gene assignment rules in the Recon2 model. The 'gene<sub>A</sub> OR gene<sub>B</sub>' operator was interpreted as arithmetic mean of expression regulation estimates for the genes A and B. The 'gene<sub>A</sub> AND gene<sub>B</sub>' operator was interpreted as minimum of expression regulation estimates for the gene A and gene B (21). Standard deviation, 95% confidence intervals and p values for the predicted reaction regulation estimates were obtained by a Monte Carlo simulation with n = 3000 draws from normal distribution of gene expression regulation estimates (mean: gene expression regulation estimate, standard deviation: standard error of the gene expression regulation estimate) (21,41). P values were corrected for multiple testing with the FDR method. The analysis was done with the packages *BiGGR* (21) and *biggrExtra* (41) (function `build_geneSBML()`). Biochemical reactions predicted to be significantly activated or inhibited in the collagen high cluster as compared with the collagen low cluster are listed in **Supplementary Table S13**.



Enrichment of significantly activated or inhibited biochemical reactions within the Recon model metabolic subsystems (22) was investigated by comparing the frequency of significantly activated or inhibited reactions in the given subsystem with 10000 random samples from the total reaction pool (function `suba()`, package `biggrExtra`). The magnitude of enrichment was assessed with OR defined for the  $i$ -th metabolic subsystem as follows:

$$OR_i = \frac{n_{regulated,i} \times n_{total}}{n_i \times n_{regulated}}$$

where  $n_{regulated,i}$  denotes the number of significantly activated or inhibited biochemical reactions within the  $i$ -th metabolic subsystem,  $n_{total}$  denotes the total number of investigated reactions,  $n_i$  denotes the number of metabolic reactions within the  $i$ -th metabolic subsystem, and  $n_{regulated}$  is the total number of significantly activated or inhibited biochemical reactions. Common significantly activated or inhibited metabolic subsystems were identified as metabolic subsystems significantly activated or inhibited in at least four cohorts, with at least weak effect size defined by  $OR \geq 1.44$  (**Supplementary Figure S8, Supplementary Table S14**).

## Survival analysis and transcriptional collagen score

Biochemical relapse-free survival was modeled in the TCGA training cohort as a function of normalized  $\log_2$ -transformed expression of the collagen-related genes of interest by multi-parameter Elastic Net Cox modeling (42,43). Of note, both first and second order terms were included in the Elastic Net model to account for possible non-linear association of gene expression with survival probability. The optimal values of the shrinkage parameter  $\lambda$  was found by tuning in 200 repeats 10-fold cross-validation and corresponded to the model with the lowest deviance (function `cv.glmnet()`, package `glmnet`) (23). The Elastic Net model was trained for the TCGA data set for the tuned  $\lambda$  parameter with the `glmnet()` function from the `glmnet` package. The Transcriptional Collagen Score was defined as the linear predictor score of the Elastic Net model and was computed for the training TCGA cohort and the validation collectives GSE54460, GSE70768, GSE70769, GSE220095 and DKFZ with the `predict()` method. Gene expression values used for calculation of the Transcriptional Collagen Score, i.e. non-zero coefficients of the Elastic Net model, are presented as hazard ratios (HR) in **Figure 3A** and **Supplementary Table S15**.

Performance of the Transcriptional Collagen Score at prediction of biochemical relapse-free survival was assessed by uni-variable Cox regression (functions `coxph()` and `as_coxex()`, packages `survival` and `coxExtensions`) (24). The proportional hazard assumption of the Cox models was checked with the `summary(type = 'assumptions')` method, which implements the genuine `cox.zph()` algorithm (44). Harrell's concordance index (45) and integrated Brier score (46) were employed to investigate the overall concordance between the predicted and observed survival and to assess overall calibration (method `summary(type = 'fit')`, package `coxExtensions`, **Figure 3B** and **Supplementary**



**Table S16).** Statistical significance of differences in biochemical relapse-free survival between patients stratified by tertiles of the Transcriptional Collagen Score was determined by Peto-Peto test (functions `surv_fit()` and `surv_pvalue()`, package *survminer*, **Figure 3C**).

## Data and code availability

Publicly available data sets were analyzed. Formatted data sets used for analyses will be made available upon request to the corresponding author. The transcriptome R analysis pipeline is available from GitHub ([https://github.com/PiotrTymoszuk/collagen\\_pca](https://github.com/PiotrTymoszuk/collagen_pca)).

## Supplementary Tables

*Supplementary Table S1: Characteristic of the analyzed cohorts. Numeric variables are presented as medians with interquartile ranges (IQR) and ranges. Qualitative variables are presented as percentages of categories within the complete observation set.*

Variable <sup>a</sup>	GSE54460	GSE70768	GSE70769	GSE22009 5	TCGA	DKFZ
Age at diagnosis, years		62 [IQR: 57 - 67] range: 41 - 93 n = 124			61 [IQR: 56 - 66] range: 41 - 78 n = 493	48 [IQR: 46 - 49] range: 32 - 52 n = 118
PSA at diagnosis	7.2 [IQR: 5.5 - 13] range: 1.8 - 73 n = 103	8 [IQR: 6 - 12] range: 3.2 - 280 n = 123	8 [IQR: 5.9 - 11] range: 1.5 - 120 n = 90	8.2 [IQR: 5.6 - 14] range: 1 - 120 n = 176	0.1 [IQR: 0.03 - 0.11] range: 0 - 320 n = 436	8.1 [IQR: 5.9 - 23] range: 1.9 - 740 n = 116
Clinical tumor stage		T1: 56% (n = 62) T2: 30% (n = 33) T3: 14% (n = 16) n = 111	T1: 46% (n = 41) T2: 44% (n = 39) T3: 10% (n = 9) n = 89			
Pathological tumor stage	T1: 13% (n = 14) T2: 70% (n = 73) T3: 16% (n = 17) T4: 0.95% (n = 1) n = 105	T2: 31% (n = 34) T3: 68% (n = 76) T4: 0.9% (n = 1) n = 111	T2: 53% (n = 48) T3: 47% (n = 42) n = 90	T2: 66% (n = 117) T3: 26% (n = 46) T4: 7.4% (n = 13) n = 176	T1: 0% (n = 0) T2: 38% (n = 186) T3: 60% (n = 290) T4: 2.1% (n = 10) n = 486	T1: 0% (n = 0) T2: 64% (n = 74) T3: 30% (n = 35) T4: 6% (n = 7) n = 116
Pathological node stage		N0: 91% (n = 82) N1: 8.9% (n = 8) n = 90	N0: 100% (n = 18) n = 18	N0: 87% (n = 146) N1: 13% (n = 22) n = 168	N0: 81% (n = 342) N1: 19% (n = 78) n = 420	
Pathological metastasis stage		M0: 86% (n = 6)	M0: 87% (n = 26)			

Variable <sup>a</sup>	GSE54460	GSE70768	GSE70769	GSE22009 5	TCGA	DKFZ
		M1: 14% (n = 1) n = 7	M1: 13% (n = 4) n = 30			
Gleason score	5: 0.94% (n = 1) 6: 9.4% (n = 10) 7: 75% (n = 80) 8: 9.4% (n = 10) 9: 4.7% (n = 5) n = 106	1: 1.6% (n = 2) 6: 14% (n = 17) 7: 71% (n = 87) 8: 7.4% (n = 9) 9: 5.7% (n = 7) n = 122	1: 1.1% (n = 1) 5: 2.2% (n = 2) 6: 20% (n = 18) 7: 62% (n = 56) 8: 5.5% (n = 5) 9: 9.9% (n = 9) n = 91	1: 0.57% (n = 1) 6: 20% (n = 35) 7: 68% (n = 120) 8: 2.8% (n = 5) 9: 8.5% (n = 15) n = 176	7 [IQR: 7 - 9] range: 6 - 10 n = 493	6: 11% (n = 13) 7: 74% (n = 87) 8: 0.85% (n = 1) 9: 14% (n = 16) 10: 0.85% (n = 1) n = 118
	5 - 6: 10% (n = 11) 7: 75% (n = 80) 8+: 14% (n = 15) n = 106	5 - 6: 16% (n = 19) 7: 71% (n = 87) 8+: 13% (n = 16) n = 122	5 - 6: 23% (n = 21) 7: 62% (n = 56) 8+: 15% (n = 14) n = 91	5 - 6: 20% (n = 36) 7: 68% (n = 120) 8+: 11% (n = 20) n = 176	5 - 6: 9.1% (n = 45) 7: 50% (n = 245) 8+: 41% (n = 203) n = 493	5 - 6: 11% (n = 13) 7: 74% (n = 87) 8+: 15% (n = 18) n = 118
Surgical margins	negative: 60% (n = 61) positive: 40% (n = 40) n = 101	negative: 78% (n = 93) positive: 22% (n = 26) n = 119	negative: 55% (n = 51) positive: 45% (n = 42) n = 93			
Extracapsular extension			46% (n = 42) n = 91			
Death					2% (n = 10) n = 493	
Overall survival, months					30 [IQR: 17 - 48] range: 0.76 - 170 n = 493	
Biochemical	52% (n =	17% (n =	48% (n =	43% (n =	19% (n =	23% (n =

<b>Variable<sup>a</sup></b>	<b>GSE54460</b>	<b>GSE70768</b>	<b>GSE70769</b>	<b>GSE22009 5</b>	<b>TCGA</b>	<b>DKFZ</b>
relapse	55) n = 106	19) n = 112	45) n = 93	75) n = 176	93) n = 493	24) n = 105
Biochemical relapse-free survival, months	49 [IQR: 18 - 77] range: 0 - 170 n = 106	30 [IQR: 17 - 49] range: 1 - 65 n = 111	58 [IQR: 19 - 80] range: 0.36 - 100 n = 92	75 [IQR: 46 - 110] range: 0.66 - 130 n = 176	26 [IQR: 14 - 45] range: 0.76 - 170 n = 493	36 [IQR: 13 - 49] range: 0.5 - 76 n = 105

<sup>a</sup>PSA: prostate-specific antigen.

*Supplementary Table S2: Collagen-related genes and their classification.*

<b>Gene group</b>	<b>Gene symbol</b>	<b>Entrez ID</b>
proline turnover	<i>ALDH18A1</i>	5832
	<i>PEPD</i>	5184
	<i>PYCR1</i>	5831
collagen modification	<i>LOX</i>	4015
	<i>LOXL1</i>	4016
	<i>LOXL2</i>	4017
	<i>P4HA1</i>	5033
	<i>P4HA2</i>	8974
	<i>P4HB</i>	5034
	<i>PLOD1</i>	5351
	<i>PLOD2</i>	5352
	<i>PLOD3</i>	8985
ECM component	<i>PPIB</i>	5479
	<i>COL11A1</i>	1301
	<i>COL11A2</i>	1302
	<i>COL14A1</i>	7373
	<i>COL15A1</i>	1306
	<i>COL16A1</i>	1307
	<i>COL17A1</i>	1308
	<i>COL18A1</i>	80781
	<i>COL19A1</i>	1310
	<i>COL1A1</i>	1277
	<i>COL1A2</i>	1278
	<i>COL2A1</i>	1280
	<i>COL3A1</i>	1281

Gene group	Gene symbol	Entrez ID
	<i>COL4A1</i>	1282
	<i>COL4A2</i>	1284
	<i>COL4A3</i>	1285
	<i>COL4A5</i>	1287
	<i>COL4A6</i>	1288
	<i>COL5A1</i>	1289
	<i>COL5A2</i>	1290
	<i>COL6A1</i>	1291
	<i>COL6A2</i>	1292
	<i>COL6A3</i>	1293
	<i>COL7A1</i>	1294
	<i>COL9A1</i>	1297
	<i>COL9A2</i>	1298
	<i>COL9A3</i>	1299
	<i>LAMA3</i>	3909
	<i>LAMB3</i>	3914
	<i>LAMC2</i>	3918
ECM processing	<i>ADAMTS2</i>	9509
	<i>BMP1</i>	649
	<i>CTSS</i>	1520
	<i>MMP13</i>	4322
	<i>MMP7</i>	4316
	<i>MMP9</i>	4318
	<i>PCOLCE</i>	5118

Gene group	Gene symbol	Entrez ID
adhesion	<i>PCOLCE2</i>	26577
	<i>SERPINH1</i>	871
	<i>CD151</i>	977
	<i>DST</i>	667
	<i>ITGA6</i>	3655
	<i>ITGB4</i>	3691



*Supplementary Table S3: Expression of the collagen pathway genes in the malignant and benign tissue compared by paired T test with Cohen's d effect size statistic. P values were corrected for multiple testing with the false discovery rate method. log2-transformed expression values are presented as medians with interquartile ranges (IQR) and ranges. The table is available as a supplementary Excel file.*

*Supplementary Table S4: Expression of the collagen pathway genes in cancer samples stratified by Gleason scores compared by one-way ANOVA with eta-square effect size statistic. P values were corrected for multiple testing with the false discovery rate method. log2-transformed expression values are presented as medians with interquartile ranges (IQR) and ranges. The table is available as a supplementary Excel file.*

*Supplementary Table S5: Expression of the cluster-defining collagen pathway genes in the collagen clusters of prostate cancer. Statistical significance was assessed by two-tailed T test with Cohen's d effect size statistic. P values were corrected for multiple testing with the false discovery rate method. log2-transformed expression values are presented as medians with interquartile ranges (IQR) and ranges. The table is available as a supplementary Excel file.*

*Supplementary Table S6: Clinical characteristic of the collagen clusters. Numeric variables are presented as medians with interquartile ranges (IQR) and ranges. Nominal variables are presented as percentages and counts of categories within the cluster.*

Cohort	Variable <sup>a</sup>	Collagen low	Collagen high	Significance <sup>b</sup>	Effect size <sup>b</sup>
GSE54460	PSA at diagnosis	7.1 [IQR: 5.3 - 13] range: 1.8 - 40 n = 57	7.2 [IQR: 5.6 - 11] range: 1.8 - 73 n = 46	ns (p = 1)	r = 0.0085
	Pathological tumor stage	T1: 1.7% (n = 1) T2: 81% (n = 48) T3: 15% (n = 9) T4: 1.7% (n = 1) n = 59	T1: 28% (n = 13) T2: 54% (n = 25) T3: 17% (n = 8) T4: 0% (n = 0) n = 46	p = 0.0025	V = 0.41
	Gleason score	5 - 6: 12% (n = 7) 7: 78% (n = 47) 8+: 10% (n = 6) n = 60	5 - 6: 8.7% (n = 4) 7: 72% (n = 33) 8+: 20% (n = 9) n = 46	ns (p = 0.72)	V = 0.14
	Surgical margins	negative: 61% (n = 34) positive: 39% (n = 22) n = 56	negative: 60% (n = 27) positive: 40% (n = 18) n = 45	ns (p = 1)	V = 0.0073
GSE70768	Age at diagnosis, years	62 [IQR: 56 - 65] range: 42 - 73 n = 62	63 [IQR: 58 - 69] range: 41 - 93 n = 62	ns (p = 0.15)	r = 0.17
	PSA at diagnosis	7.1 [IQR: 5.8 - 11] range: 4 - 19 n = 62	8.7 [IQR: 7 - 14] range: 3.2 - 280 n = 61	ns (p = 0.15)	r = 0.19
	Clinical tumor stage	T1: 57% (n = 35) T2: 30% (n = 18) T3: 13% (n = 8) n = 61	T1: 54% (n = 27) T2: 30% (n = 15) T3: 16% (n = 8) n = 50	ns (p = 0.9)	V = 0.044
	Pathological tumor stage	T2: 34% (n = 21) T3: 66% (n = 40) T4: 0% (n = 0) n = 61	T2: 26% (n = 13) T3: 72% (n = 36) T4: 2% (n = 1) n = 50	ns (p = 0.64)	V = 0.14
	Pathological node stage	N0: 89% (n = 42) N1: 11% (n = 5) n = 47	N0: 93% (n = 40) N1: 7% (n = 3) n = 43	ns (p = 0.9)	V = 0.064
	Gleason score	5 - 6: 9.7% (n = 6) 7: 81% (n = 50) 8+: 9.7% (n = 6)	5 - 6: 22% (n = 13) 7: 62% (n = 37) 8+: 17% (n = 10)	ns (p = 0.15)	V = 0.21

Cohort	Variable <sup>a</sup>	Collagen low	Collagen high	Significance <sup>b</sup>	Effect size <sup>b</sup>
GSE70769		n = 62	n = 60		
	Surgical margins	negative: 81% (n = 50) positive: 19% (n = 12) n = 62	negative: 75% (n = 43) positive: 25% (n = 14) n = 57	ns (p = 0.9)	V = 0.063
	PSA at diagnosis	6.9 [IQR: 5.1 - 11] range: 2.2 - 35 n = 42	8.6 [IQR: 6.4 - 12] range: 1.5 - 120 n = 48	ns (p = 0.27)	r = 0.19
	Clinical tumor stage	T1: 50% (n = 21) T2: 36% (n = 15) T3: 14% (n = 6) n = 42	T1: 43% (n = 20) T2: 51% (n = 24) T3: 6.4% (n = 3) n = 47	ns (p = 0.34)	V = 0.18
	Pathological tumor stage	T2: 61% (n = 27) T3: 39% (n = 17) n = 44	T2: 46% (n = 21) T3: 54% (n = 25) n = 46	ns (p = 0.34)	V = 0.16
	Pathological node stage	N0: 100% (n = 7) n = 7	N0: 100% (n = 11) n = 11	ns (p = 0.4)	V = Inf
	Gleason score	5 - 6: 32% (n = 14) 7: 59% (n = 26) 8+: 9.1% (n = 4) n = 44	5 - 6: 15% (n = 7) 7: 64% (n = 30) 8+: 21% (n = 10) n = 47	ns (p = 0.27)	V = 0.24
GSE220095	Surgical margins	negative: 57% (n = 25) positive: 43% (n = 19) n = 44	negative: 53% (n = 26) positive: 47% (n = 23) n = 49	ns (p = 0.88)	V = 0.038
	Extracapsular extension	39% (n = 17) n = 44	53% (n = 25) n = 47	ns (p = 0.34)	V = 0.15
	PSA at diagnosis	8.4 [IQR: 5.7 - 14] range: 1 - 77 n = 78	8.1 [IQR: 5.4 - 14] range: 2.4 - 120 n = 98	ns (p = 0.83)	r = 0.016
	Pathological tumor stage	T2: 65% (n = 51) T3: 24% (n = 19) T4: 10% (n = 8) n = 78	T2: 67% (n = 66) T3: 28% (n = 27) T4: 5.1% (n = 5) n = 98	ns (p = 0.55)	V = 0.1
	Pathological node stage	N0: 91% (n = 68) N1: 9.3% (n = 7) n = 75	N0: 84% (n = 78) N1: 16% (n = 15) n = 93	ns (p = 0.55)	V = 0.1
	Gleason score	5 - 6: 13% (n = 10) 7: 79% (n = 62)	5 - 6: 27% (n = 26) 7: 59% (n = 58)	ns (p = 0.064)	V = 0.22

Cohort	Variable <sup>a</sup>	Collagen low	Collagen high	Significance <sup>b</sup>	Effect size <sup>b</sup>
TCGA		8+: 7.7% (n = 6) n = 78	8+: 14% (n = 14) n = 98		
	Age at diagnosis, years	61 [IQR: 56 - 66] range: 44 - 78 n = 236	62 [IQR: 57 - 67] range: 41 - 77 n = 257	ns (p = 0.16)	r = 0.075
	PSA at diagnosis	0.1 [IQR: 0.03 - 0.11] range: 0 - 320 n = 206	0.1 [IQR: 0.03 - 0.12] range: 0 - 37 n = 230	ns (p = 0.88)	r = 0.018
	Pathological tumor stage	T2: 44% (n = 102) T3: 54% (n = 125) T4: 2.6% (n = 6) n = 233	T2: 33% (n = 84) T3: 65% (n = 165) T4: 1.6% (n = 4) n = 253	ns (p = 0.081)	V = 0.12
	Pathological node stage	N0: 82% (n = 161) N1: 18% (n = 36) n = 197	N0: 81% (n = 181) N1: 19% (n = 42) n = 223	ns (p = 0.98)	V = 0.0072
DKFZ	Gleason score	5 - 6: 11% (n = 26) 7: 55% (n = 129) 8+: 34% (n = 81) n = 236	5 - 6: 7.4% (n = 19) 7: 45% (n = 116) 8+: 47% (n = 122) n = 257	ns (p = 0.051)	V = 0.14
	Age at diagnosis, years	48 [IQR: 45 - 49] range: 38 - 52 n = 61	48 [IQR: 46 - 49] range: 32 - 52 n = 57	ns (p = 0.54)	r = 0.056
	PSA at diagnosis	8.4 [IQR: 5.8 - 16] range: 1.9 - 740 n = 60	7.7 [IQR: 6 - 39] range: 3.2 - 150 n = 56	ns (p = 0.54)	r = 0.065
	Pathological tumor stage	T2: 68% (n = 41) T3: 32% (n = 19) T4: 0% (n = 0) n = 60	T2: 59% (n = 33) T3: 29% (n = 16) T4: 12% (n = 7) n = 56	ns (p = 0.073)	V = 0.26
	Gleason score	5 - 6: 9.8% (n = 6) 7: 80% (n = 49) 8+: 9.8% (n = 6) n = 61	5 - 6: 12% (n = 7) 7: 67% (n = 38) 8+: 21% (n = 12) n = 57	ns (p = 0.38)	V = 0.17

<sup>a</sup>PSA: prostate-specific antigen.

<sup>b</sup>Qualitative variables:  $\chi^2$  test with Cramer V effect size statistic. Numeric variables: Mann-Whitney test with r effect size statistic. P values corrected for multiple testing with the false discovery rate.

*Supplementary Table S7: Non-malignant cell numbers predicted for the collagen clusters by the MCP Counter and xCell algorithms. Statistical significance was assessed by Mann-Whitney test with r effect size statistic. P values were corrected for multiple testing with the false discovery method. The table is available as a supplementary Excel file.*

*Supplementary Table S8: Gene set variation analysis with the Reactome pathway gene signatures. Differences in ssGSEA scores between collagen high and collagen low cancers were investigated by two-tailed T test with Cohen's d effect size statistic. Results for signatures significantly regulated with at least weak effect size (d at least 0.2) in at least five cohorts are presented. P values were corrected for multiple testing with the false discovery rate method (FDR). The table is available as a supplementary Excel file.*

*Supplementary Table S9: Genes differentially expressed in the collagen high cluster as compared with collagen low cancers were identified by two-tailed T test with the 1.25-fold regulation cutoff and the Cohen's d effect size statistic of 0.2 P values were corrected for multiple testing with the false discovery rate method (FDR). The table is available as a supplementary Excel file.*

*Supplementary Table S10: Biological process gene ontology (GO) term enrichment within genes differentially regulated in the collagen clusters. The enrichment analysis was performed with goana tool, enrichment p values were corrected for multiple testing with the false discovery rate (FDR) method. Significant enrichment was defined by  $pFDR < 0.05$  and odds ratio (OR) for enrichment within differentially regulated genes of at least 1.44. OR for enrichment within genes upregulated in the collagen high and collagen low clusters are presented for significant GO terms shared by at least five cohorts. The table is available as a supplementary Excel file.*

*Supplementary Table S11: Activity of transcriptional regulons in the collagen high cluster as compared with the collagen low cluster predicted by the collecTRI model. Regulon activity was estimated with uni-parameter linear modeling with whole-transcriptome effect sizes of differential gene expression, p values were corrected for multiple testing with the false discovery rate (FDR) method. Linear model scores are presented for regulons significantly activated or inhibited in at least five cohorts. The table is available as a supplementary Excel file.*

*Supplementary Table S12: Activity of signaling pathways in the collagen high cluster as compared with the collagen low cluster predicted by the PROGENy model. Pathway activity was estimated with multi-parameter linear modeling with whole-transcriptome effect sizes of differential gene expression. P values were corrected for multiple testing with the false discovery rate (FDR) method, linear model scores serve as measures of pathway activity. The table is available as a supplementary Excel file.*

*Supplementary Table S13: Biochemical reactions predicted to be significantly activated in collagen high as compared with collagen low cancers. Statistical significance was determined by a Monte Carlo simulation. P values were corrected for multiple testing with the false discovery rate (FDR) method. The table is available as a supplementary Excel file.*

*Table 14: Results of enrichment analysis for significantly activated and inhibited biochemical reactions within the Recon metabolism subsystem. Statistical significance was determined by random sampling from the entire reaction pool p values were corrected for multiple testing with the false discovery rate (FDR) method. Effect size of enrichment of the subsystem in significantly activated or inhibited reactions was measured by odds ratio (OR) statistic. The table is available as a supplementary Excel file.*

---

---

---



*Supplementary Table S15: Transcriptomic collagen score was established by Elastic Net Cox modeling of biochemical relapse-free survival in the TCGA cohort. Hazard ratios (HR) of the member genes are presented.*

<b>Variable</b>	<b>HR</b>
<i>ALDH18A1</i>	0.942
COL11A2	1.430
COL14A1	0.983
COL1A1	1.080
COL4A6	0.901
COL6A3	0.863
COL7A1	1.030
COL9A1	1.000
CTSS	1.060
LAMA3	1.040
LOX	0.980
MMP13	1.010
P4HB	0.932
PCOLCE2	0.950
PLOD3	1.040
SERPINH1	1.050
<i>ALDH18A1</i> <sup>2</sup>	0.962
<i>COL11A2</i> <sup>2</sup>	1.010
<i>COL1A1</i> <sup>2</sup>	1.190
<i>COL4A3</i> <sup>2</sup>	0.997
<i>COL4A6</i> <sup>2</sup>	0.905
<i>COL6A3</i> <sup>2</sup>	0.996
<i>COL7A1</i> <sup>2</sup>	1.050

Variable	HR
COL9A1 <sup>2</sup>	1.070
COL9A3 <sup>2</sup>	0.947
CTSS <sup>2</sup>	1.010
DST <sup>2</sup>	0.999
LAMA3 <sup>2</sup>	1.000
LAMC2 <sup>2</sup>	0.991
MMP7 <sup>2</sup>	0.890
P4HA1 <sup>2</sup>	1.070
P4HB <sup>2</sup>	0.949
PCOLCE2 <sup>2</sup>	0.975
PLOD3 <sup>2</sup>	1.030
SERPINH1 <sup>2</sup>	1.030

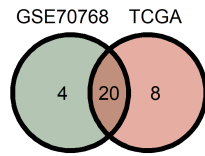
*Supplementary Table S16: Performance of the transcriptomic collagen score at prediction of biochemical relapse-free survival in the TCGA training cohort and validation collectives.*

<b>Data set type</b>	<b>Cohort</b>	<b>Complete cases</b>	<b>Cases with biochemical relapse</b>	<b>Concordance index</b>	<b>Integrated Brier score</b>
training	TCGA	493	93	0.764	0.176
	DKFZ	105	24	0.734	0.133
	GSE220095	176	75	0.572	0.188
test	GSE54460	106	55	0.597	0.233
	GSE70768	111	19	0.667	0.121
	GSE70769	92	45	0.618	0.210

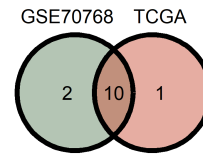
## Supplementary Figures

**A**

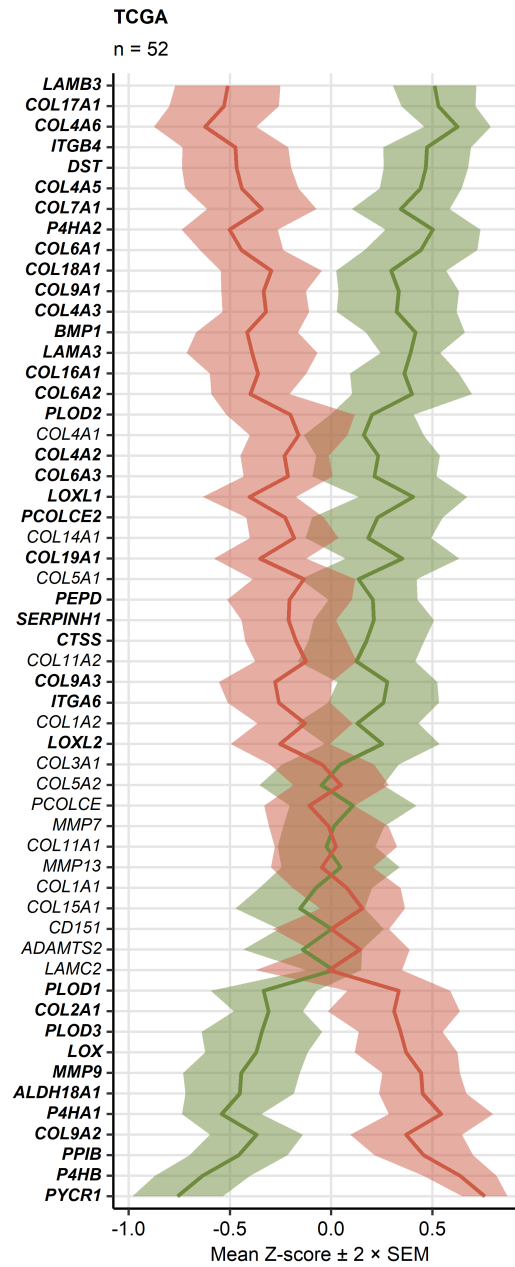
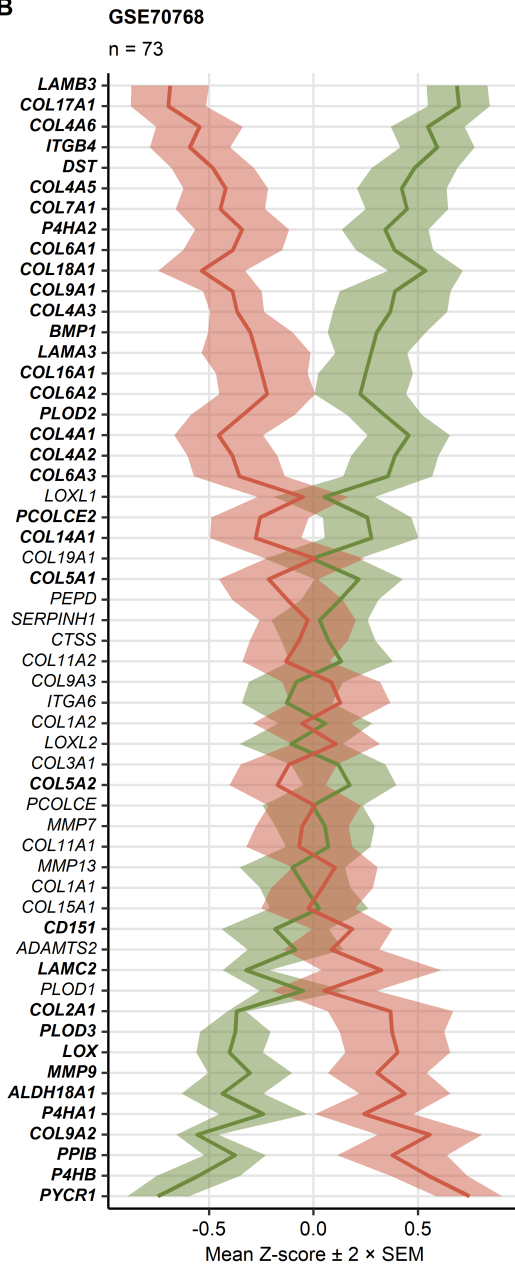
**Upregulated transcripts**  
Tumor vs benign



**Downregulated transcripts**  
Tumor vs benign



**B**



**Supplementary Figure S1. Expression of collagen-related genes in the normal prostate and prostate cancer tissue.**

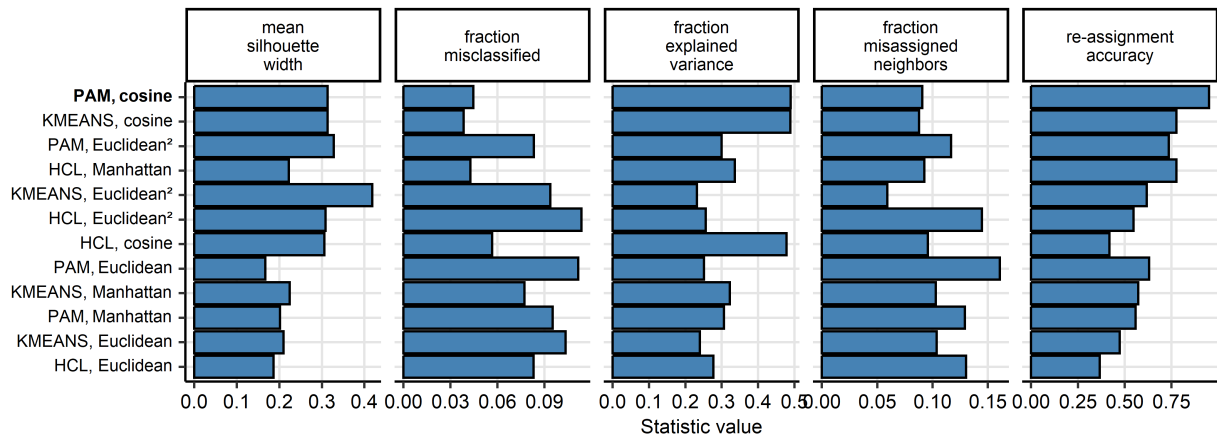
*Differences in  $\log_2$ -transformed expression of 55 collagen-related genes between donor-matched pairs of the prostate cancer and benign tissue were assessed by paired T test with Cohen's effect size statistic in the GSE70768 and TCGA cohorts. P values were corrected for multiple testing with the false discovery rate (FDR) method. Full analysis results are listed in Supplementary Table S3*

*(A) Numbers of significantly up- and downregulated genes in the tumor tissue as compared with the benign tissue in the investigated cohorts presented in Venn plots.*

*(B) Mean normalized  $\log_2$  expression of the collagen-related genes presented as lines. Tinted regions represent the  $2 \times SEM$  (standard error of the mean) intervals. Significant effects are highlighted with bold font in the Y axis. Numbers of tissue pairs are displayed in the plot captions.*

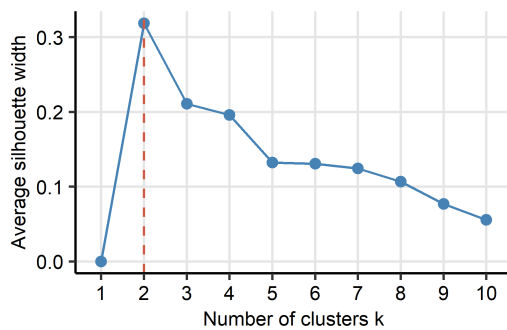
**A**

**Choice of clustering algorithm, TCGA, cross-validation**



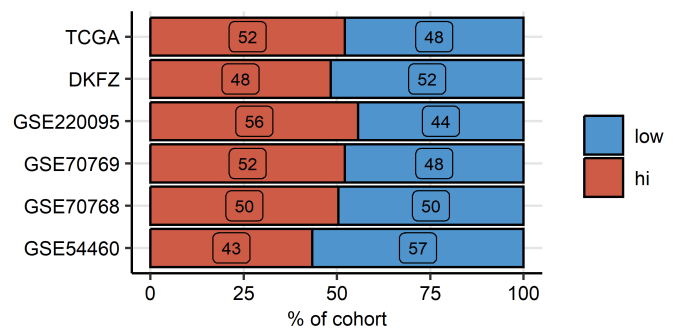
**B**

**Choice of cluster number, TCGA**



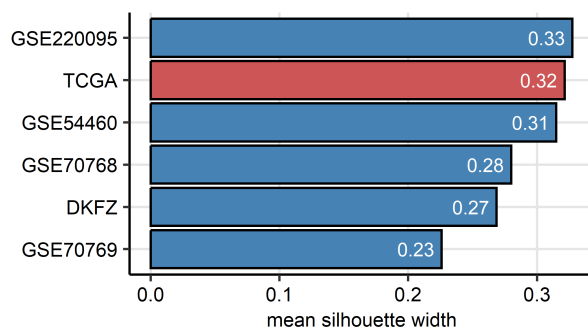
**C**

**Collagen cluster size**

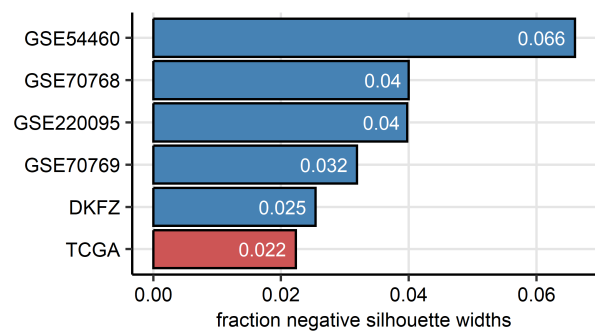


**D**

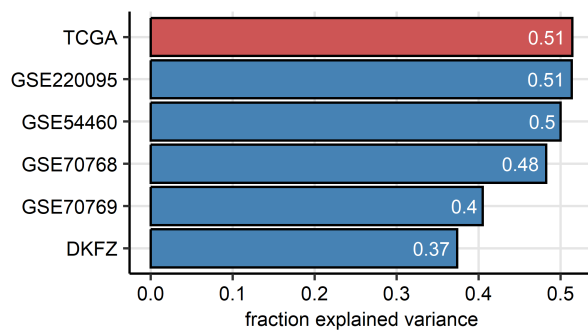
**Cluster separation**



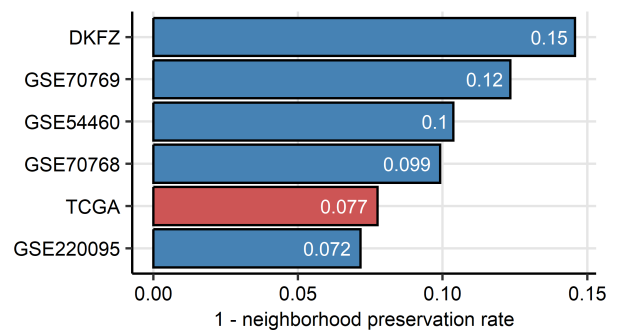
**Misclassification rate**



**Explained variance**



**Misassigned neighbors**





## **Supplementary Figure S2. Semi-supervised clustering of prostate cancer samples in respect to expression of collagen-related genes.**

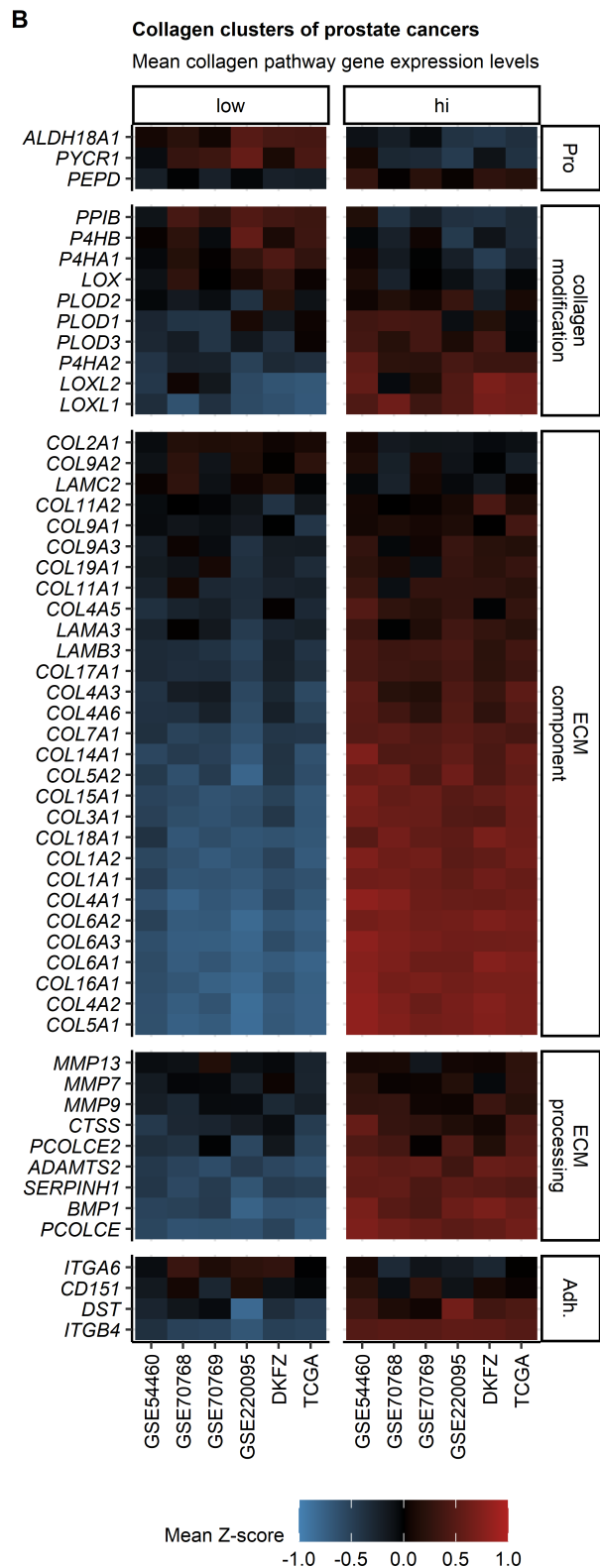
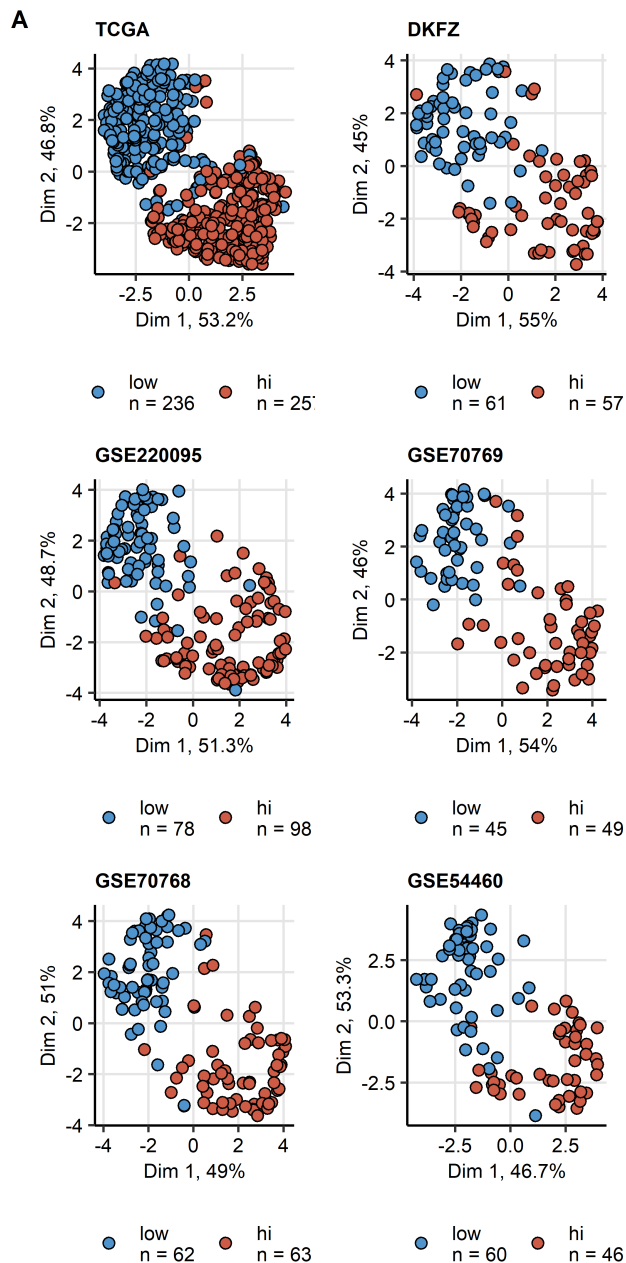
*Cancer samples in the TCGA training cohort (n = 493) were clustered in respect to normalized,  $\log_2$ -transformed expression levels of the collagen-related genes of interest with the PAM (partition around medoids) algorithm with cosine distance metric. Two clusters were defined: collagen low and collagen high. Assignment of the tumor samples in the training GSE54460 (n = 106), GSE70768 (n = 125), GSE70769 (n = 94), GSE220095 (n = 176), and DKFZ collective (n = 118) to the collagen clusters was accomplished by an inverse distance-weighted k-nearest neighbor classifier.*

*(A) Comparison of performance of several clustering algorithms (PAM: partition around medoids, HCL: hierarchical clustering/Ward D2 and KMEANS) and distance metrics (Euclidean, Manhattan, squared Euclidean and cosine) in 5-fold cross-validation of the TCGA data set gauged by mean silhouette width as a measure of cluster separation, fraction of observations with negative silhouette widths indicative of possible misclassification, explained clustering variance and fraction of misclassified 5-nearest neighbors, cluster assignment accuracy in the cross-validation folds. Note superior performance of the PAM/cosine algorithm used for definition of the collagen clusters.*

*(B) Choice of cluster number (k = 2) for PAM/cosine clustering of cancer samples in the TCGA cohort motivated by the peak mean silhouette width statistic.*

*(C) Distribution of sizes of the collagen clusters in the TCGA training cohort and the validation collectives. Percentages of samples assigned to the collagen clusters are presented in a stack plot.*

*(D) Quality of collagen clusters in the training cohort (TCGA) and validation collectives. Cluster separation was assessed by mean silhouette width, misclassification rate is expressed as fraction of observations with negative silhouette widths. Explanatory performance was gauged by fraction of explained clustering variance. Neighborhood preservation was assessed by mean fraction of 5-nearest neighbors placed in different clusters.*

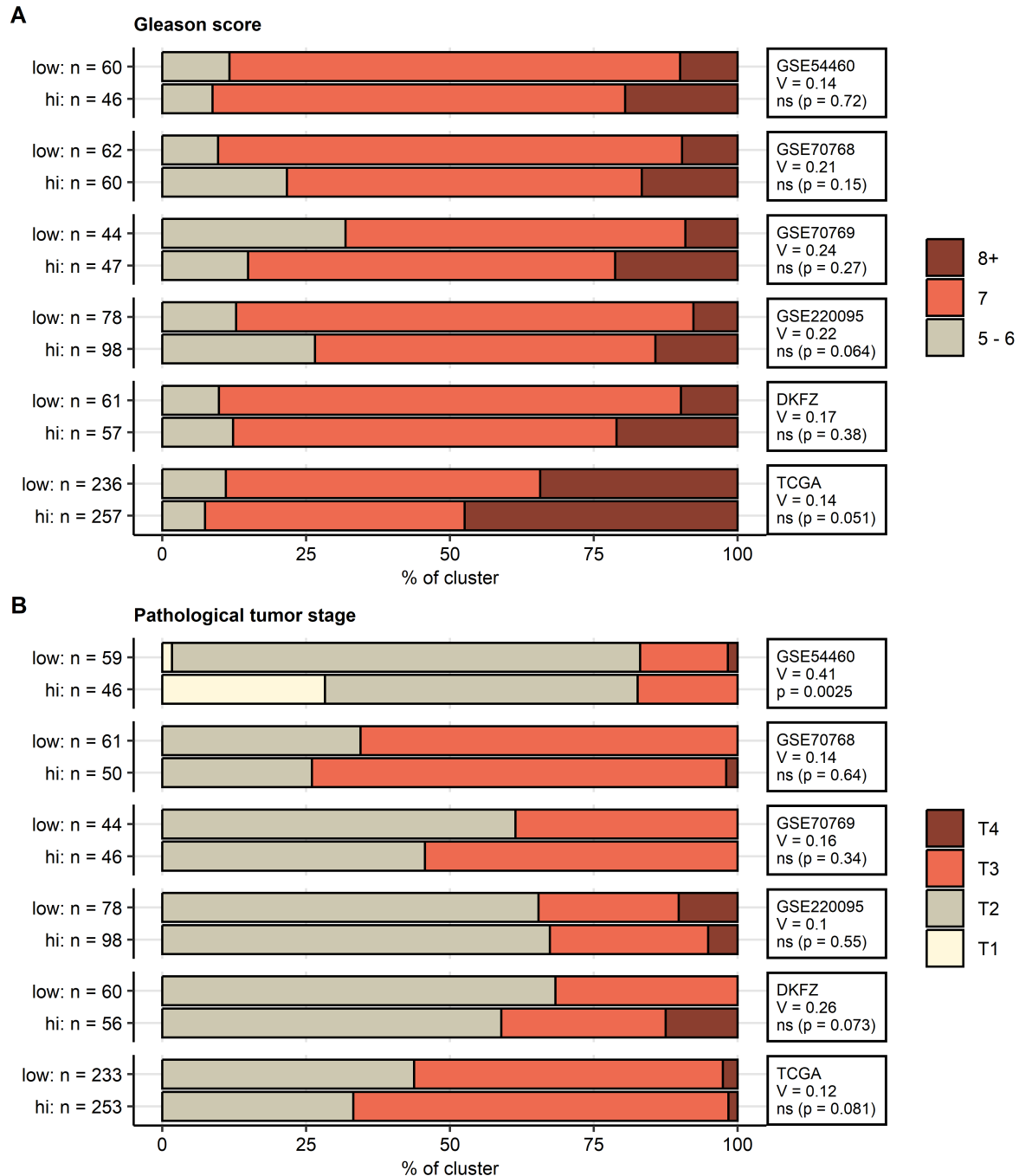


**Supplementary Figure S3. Expression of the collagen-related genes in the collagen clusters of prostate cancers.**

*The collagen clusters of prostate cancer samples were developed by semi-supervised PAM clustering as presented in Supplementary Figure S2. Results of comparison of levels of the cluster-defining factors between the collagen clusters are listed in Supplementary Table S5.*

*(A) Uniform manifold approximation and projection (UMAP) layout of normalized  $\log_2$ -transformed expression values of the collagen-related genes utilized for definition of the collagen clusters. Each point represents a single cancer sample, cluster assignment is color-coded. Numbers of cancer samples in the collagen clusters are indicated in the figure legends.*

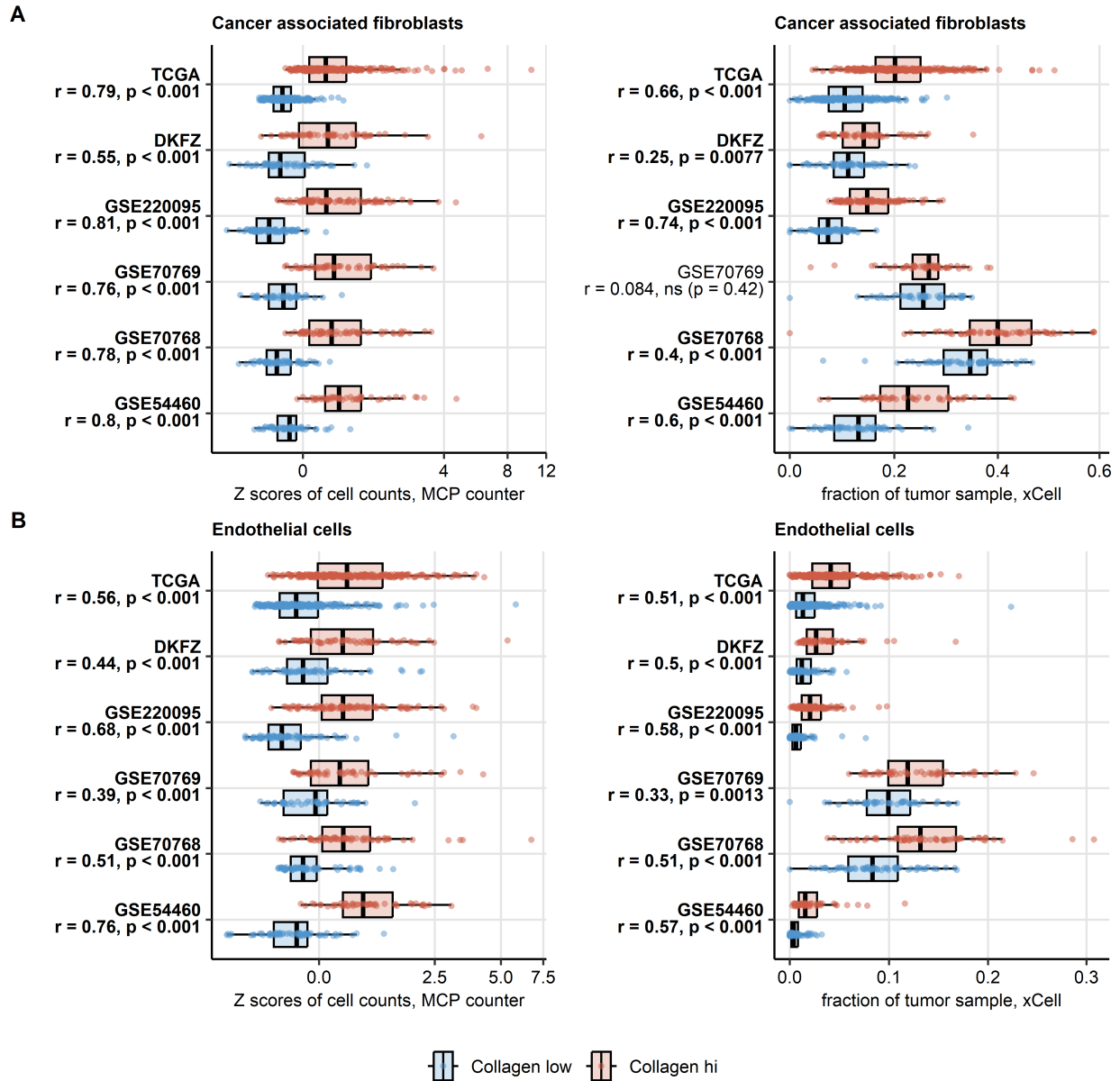
*(B) Mean normalized  $\log_2$ -transformed expression values of the cluster-defining collagen-related genes visualized in a heat map. The cluster-defining genes were classified by their biological function (Pro: proline metabolism, ECM: extracellular matrix, Adh.: adhesion).*



**Supplementary Figure S4. Gleason score and pathological tumor stage in the collagen clusters.**

Differences in distribution of Gleason scores (A, 5 - 6, 7 and  $\geq 8$ ) and pathological tumor stages (B) between the collagen clusters were investigated by  $\chi^2$  test with Cramer V effect size statistic. P values were corrected for multiple testing with the false discovery rate method.

*Percentages of samples assigned to the Gleason score strata or tumor stages within each cluster are shown in stack plots. Effect sizes and p values are displayed in the plot facets. Numbers of complete observations in the clusters are indicated in the Y axes. Full clinical and pathological characteristic of the collagen clusters is shown in Supplementary Table S6.*

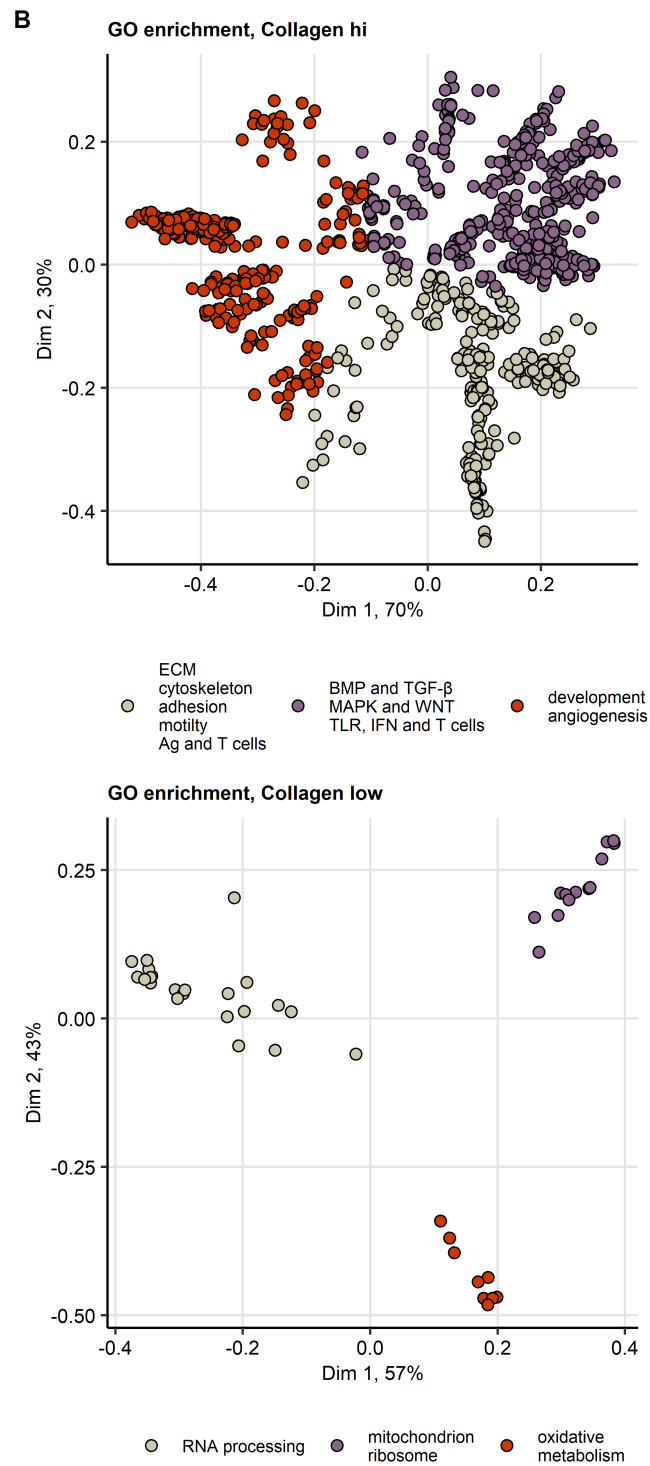
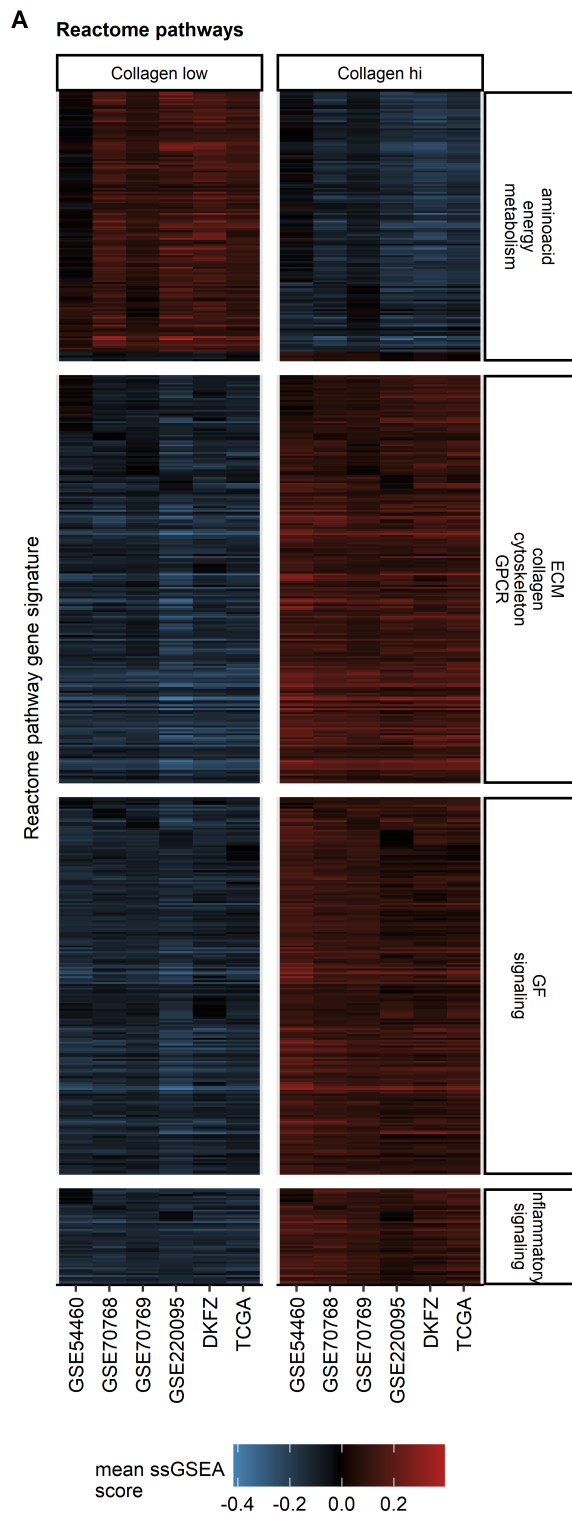


### Supplementary Figure S5. Cancer-associated fibroblasts and endothelial cell infiltration in the collagen clusters.

Levels of cancer-associated fibroblasts (A) and endothelial cells (B) in the collagen clusters were predicted by the MCP Counter and xCell algorithms. Statistical significance of differences between the clusters was assessed by Mann-Whitney test with  $r$  effect size statistic.  $P$  values were corrected for multiple testing with the false discovery rate method. Median infiltration levels with interquartile ranges are visualized as boxes, whiskers span over 150% of the interquartile ranges. Points represent single cancer samples. Effect sizes and  $p$  values are

*displayed in the Y axes, significant effects are highlighted in bold. GSE54460: Collagen low: n = 60, Collagen hi: n = 46, GSE70768: Collagen low: n = 62, Collagen hi: n = 63, GSE70769: Collagen low: n = 45, Collagen hi: n = 49, GSE220095: Collagen low: n = 78, Collagen hi: n = 98, DKFZ: Collagen low: n = 61, Collagen hi: n = 57, TCGA: Collagen low: n = 236, Collagen hi: n = 257. Results of the analysis for other cell populations are listed in Supplementary Table S7.*

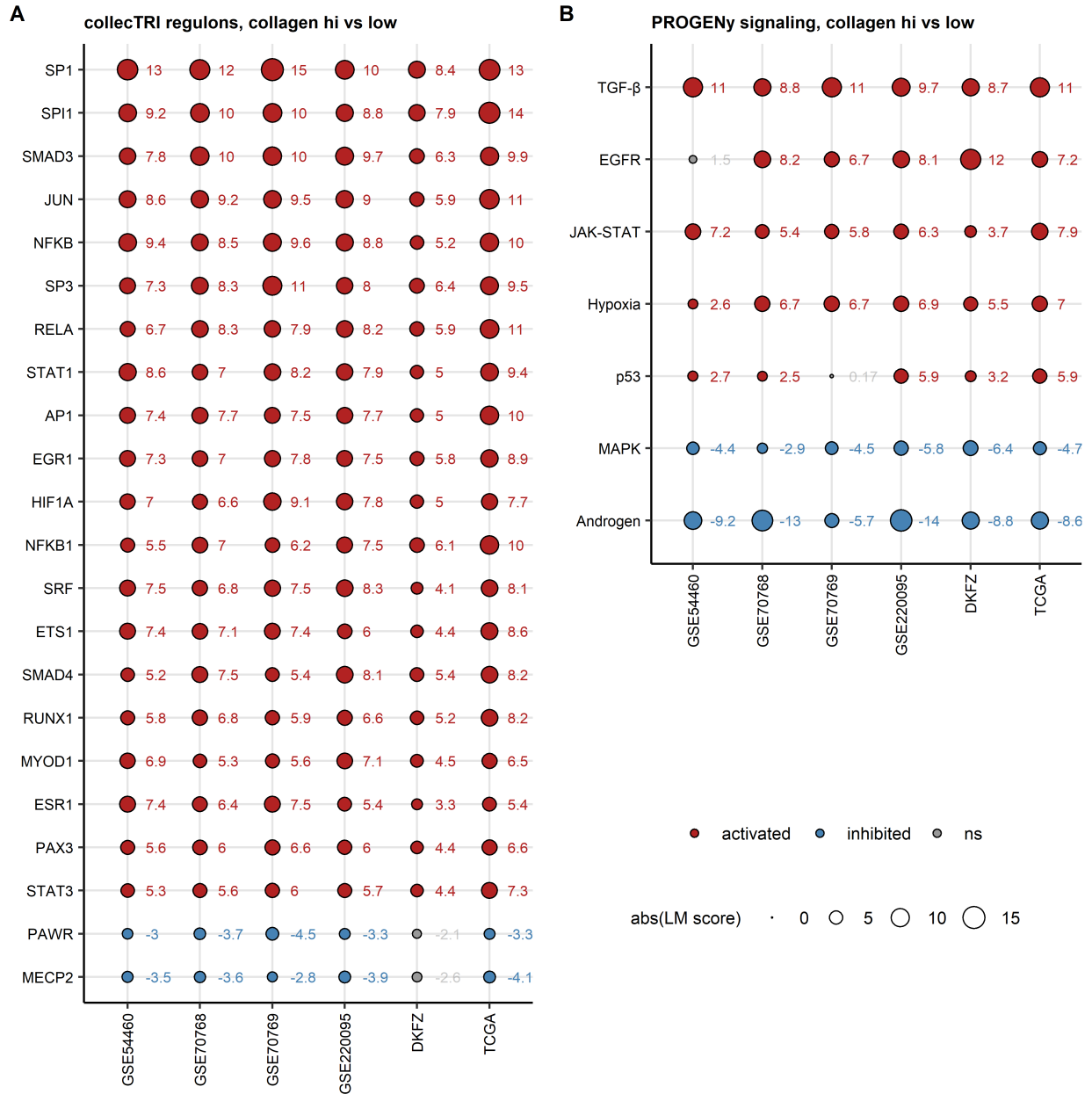




**Supplementary Figure S6. Gene set variation analysis of Reactome pathway gene signatures and gene ontology term enrichment analysis for the collagen clusters.**

*(A) Gene set variation analysis (GSVA) with Reactome pathway gene signatures. Statistical significance for differences in signature single sample gene set enrichment analysis scores (ssGSEA) between the collagen clusters was determined by two-tailed T test with Cohen's d effect size statistic. P values were corrected for multiple testing with the false discovery rate method. Mean ssGSEA scores in the clusters and cohorts for signatures significantly regulated in at least five cohorts with at least weak effect size ( $d \geq 0.2$ ) are presented in a heat map. Full analysis results are presented in Supplementary Table S8. The common regulated gene signatures were classified in respect their co-expression patterns by KMEANS clustering; signature classification is indicated in the Y axis facets of the heat map (ECM: extracellular matrix, GF: growth factor, GPCR: G protein-coupled receptor).*

*(B) Genes differentially expressed between the collagen clusters were identified by false discovery rate (FDR) corrected two-tailed T test with Cohen's d effect size statistic. Differentially expressed genes were defined by the  $pFDR < 0.05$  and  $d \geq 0.2$  cutoffs (Supplementary Table 9). Biological process gene ontology (GO) enrichment within genes significantly unregulated in the collagen high and the collagen low clusters was investigated with the goana algorithm. Enrichment p values were adjusted for multiple testing with the FDR method, odds ratio of enrichment over the whole-transcriptome GO frequency served as an effect size statistic. GO terms significantly enriched with at least weak effect size ( $OR \geq 1.44$ ) in at least five cohorts were subjected to multi-dimensional scaling and unsupervised KMEANS clustering in respect to their semantic similarity. Multi-dimensional scaling layouts of Wang distances between the common enriched GO terms are presented. Each point represents a single GO term, GO cluster assignment is color coded. Complete results of the GO enrichment analysis are listed in Supplementary Table S10.*



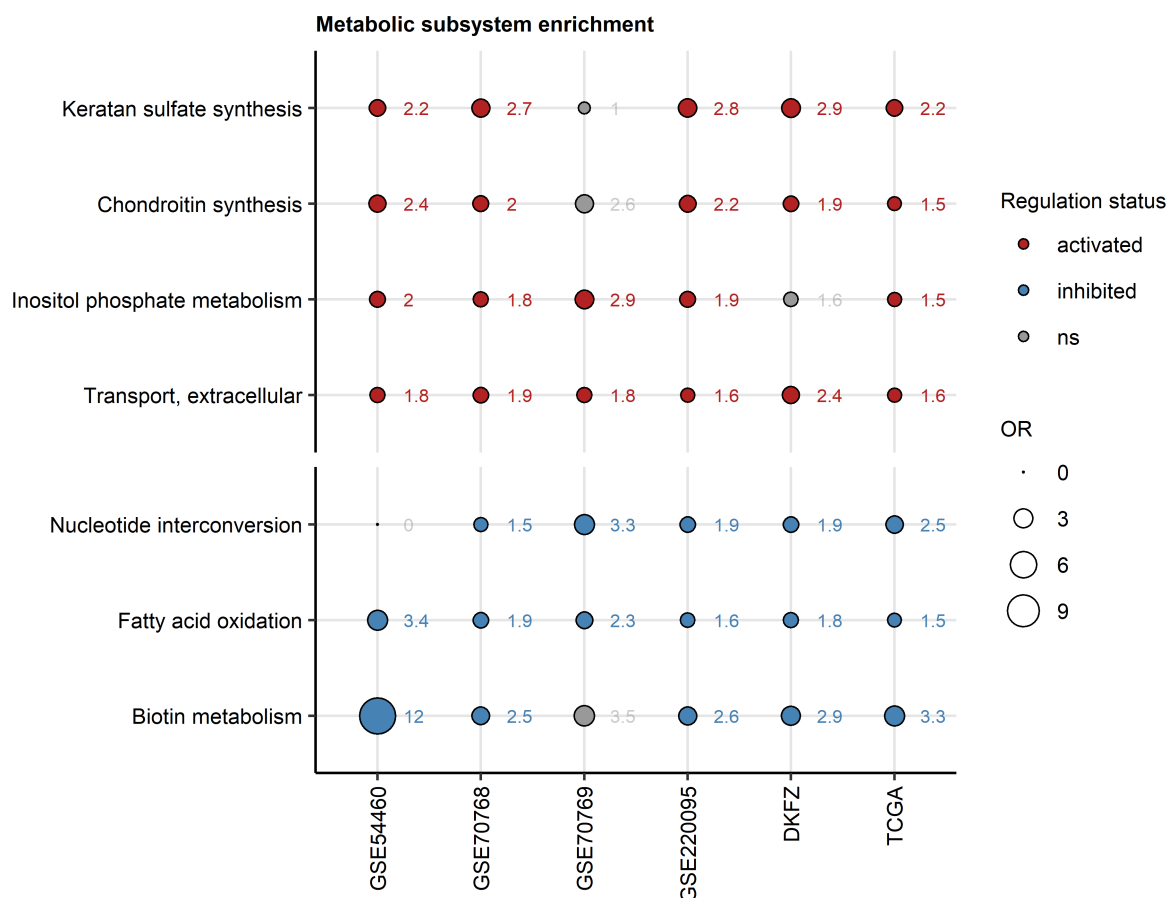
**Supplementary Figure S7. Transcriptional regulons and signaling pathway activity in the collagen clusters.**

*Differences in activity of transcriptional regulons (gene sets controlled by a common transcription factor, A) and signaling pathways (B) in the collagen high cluster as compared with the collagen low cluster were assessed by linear modeling with the collectTRI and PROGENy databases implemented by the decoupler algorithm. Linear model score (LM score) served as a measure of regulon or pathway activity, p values for significant modulation of*

*activity (LM score  $\neq 0$ ) were corrected for multiple testing with the false discovery rate method. Full analysis results are listed in Supplementary Tables S11 and S12.*

*(A) LM scores of top modulated transcriptional regulons found to be significantly activated or inhibited in at least five cohorts visualized as bubble plots. Point color codes for the regulation sign, point size codes for absolute value of the LM score. LM scores for regulon activity in particular cohorts are presented next to the data points.*

*(B) LM scores of signaling pathways found to be significantly activated or inhibited in at least five cohorts visualized as bubble plots. Point color codes for the regulation sign, point size codes for absolute value of the LM score. LM scores for regulon activity in particular cohorts are presented next to the data points.*



### Supplementary Figure S8. Metabolism in the collagen clusters.

*Modulation of Recon2 model metabolic reactions in collagen high tumors as compared with collagen low cancers was predicted by Monte Carlo simulation implemented by the BiGGR and biggrExtra algorithms based on whole-genome differential gene expression estimates (Supplementary Table 13). Enrichment of significantly activated and significantly inhibited reactions in the Recon metabolic subsystems was assessed by comparing the observed frequency of the subsystem's significantly regulated reactions with 10000 random draws from the entire reaction pool. Enrichment p values were corrected for multiple testing with the false discovery rate method. Enrichment odds ratios (OR) for significantly enriched metabolic subsystems with at least weak effect size ( $OR \geq 1.44$ ) shared by at least five cohorts are presented in a bubble plot. Point color codes for reaction modulation status, point sizes represent absolute values of OR. Points are labeled with their OR values. Full analysis results are presented in Supplementary Table S14.*

## References

1. Wickham H, Averick M, Bryan J, Chang W, McGowan L, François R, Grolemund G, Hayes A, Henry L, Hester J, et al. Welcome to the Tidyverse. *Journal of Open Source Software* (2019) 4:1686. doi: [10.21105/joss.01686](https://doi.org/10.21105/joss.01686)
2. Henry L, Wickham Hadley. rlang: Functions for Base Types and Core R and 'Tidyverse' Features. (2022) <https://cran.r-project.org/web/packages/rlang/index.html>
3. Gagolewski M, Tartanus B. Package 'stringi'. (2021) <https://cran.r-project.org/web/packages/stringi/index.html>  
<http://cran.ism.ac.jp/web/packages/stringi/stringi.pdf>
4. Sean D, Meltzer PS. GEOquery: a bridge between the Gene Expression Omnibus (GEO) and BioConductor. *Bioinformatics (Oxford, England)* (2007) 23:1846–1847. doi: [10.1093/BIOINFORMATICS/BTM254](https://doi.org/10.1093/BIOINFORMATICS/BTM254)
5. Pagès H, Carlson M, Falcon S, Li N. AnnotationDbi: Manipulation of SQLite-based annotations in Bioconductor. (2022) doi: [10.18129/B9.bioc.AnnotationDbi](https://doi.org/10.18129/B9.bioc.AnnotationDbi)
6. Carlson M. org.Hs.eg.db: Genome wide annotation for Human. (2022) doi: [10.18129/B9.bioc.org.Hs.eg.db](https://doi.org/10.18129/B9.bioc.org.Hs.eg.db)
7. Sturm G, Finotello F, List M. Immunedeconv: An R Package for Unified Access to Computational Methods for Estimating Immune Cell Fractions from Bulk RNA-Sequencing Data. *Methods in molecular biology (Clifton, NJ)* (2020) 2120:223–232. doi: [10.1007/978-1-0716-0327-7\\_16](https://doi.org/10.1007/978-1-0716-0327-7_16)
8. Aran D, Hu Z, Butte AJ. xCell: Digitally portraying the tissue cellular heterogeneity landscape. *Genome Biology* (2017) 18:220. doi: [10.1186/s13059-017-1349-1](https://doi.org/10.1186/s13059-017-1349-1)
9. Becht E, Giraldo NA, Lacroix L, Buttard B, Elarouci N, Petitprez F, Selves J, Laurent-Puig P, Sautès-Fridman C, Fridman WH, et al. Estimating the population abundance of tissue-infiltrating immune and stromal cell populations using gene expression. *Genome Biology* (2016) 17:218. doi: [10.1186/s13059-016-1070-5](https://doi.org/10.1186/s13059-016-1070-5)
10. Hänzelmann S, Castelo R, Guinney J. GSEA: Gene set variation analysis for microarray and RNA-Seq data. *BMC Bioinformatics* (2013) 14:7. doi: [10.1186/1471-2105-14-7](https://doi.org/10.1186/1471-2105-14-7)
11. Kassambara A, Mundt F. factoextra: Extract and Visualize the Results of Multivariate Data Analyses. (2020) <https://cran.r-project.org/web/packages/factoextra/index.html>

12. Schubert E, Rousseeuw PJ. Faster k-Medoids Clustering: Improving the PAM, CLARA, and CLARANS Algorithms. *Lecture notes in computer science (including subseries lecture notes in artificial intelligence and lecture notes in bioinformatics)*. Springer (2019). p. 171–187 doi: [10.1007/978-3-030-32047-8\\_16](https://doi.org/10.1007/978-3-030-32047-8_16)
13. Drost H-G. Philentropy: Information Theory and Distance Quantification with R. *Journal of Open Source Software* (2018) 3:765. doi: [10.21105/joss.00765](https://doi.org/10.21105/joss.00765)
14. Konopka T. umap: Uniform Manifold Approximation and Projection. (2022) <https://cran.r-project.org/web/packages/umap/index.html>
15. Kassambara A. rstatix: Pipe-Friendly Framework for Basic Statistical Tests. (2021) <https://cran.r-project.org/package=rstatix>
16. Young MD, Wakefield MJ, Smyth GK, Oshlack A. Gene ontology analysis for RNA-seq: accounting for selection bias. *Genome Biology* (2010) 11:R14. doi: [10.1186/gb-2010-11-2-r14](https://doi.org/10.1186/gb-2010-11-2-r14)
17. Yu G, Li F, Qin Y, Bo X, Wu Y, Wang S. GOSemSim: an R package for measuring semantic similarity among GO terms and gene products. *Bioinformatics* (2010) 26:976–978. doi: [10.1093/BIOINFORMATICS/BTQ064](https://doi.org/10.1093/BIOINFORMATICS/BTQ064)
18. Müller-Dott S, Tsirvouli E, Vazquez M, Ramirez Flores RO, Badia-i-Mompel P, Fallegger R, Türei D, Lægreid A, Saez-Rodriguez J. Expanding the coverage of regulons from high-confidence prior knowledge for accurate estimation of transcription factor activities. *Nucleic acids research* (2023) 51:10934–10949. doi: [10.1093/NAR/GKAD841](https://doi.org/10.1093/NAR/GKAD841)
19. Schubert M, Klinger B, Klünemann M, Sieber A, Uhlitz F, Sauer S, Garnett MJ, Blüthgen N, Saez-Rodriguez J. Perturbation-response genes reveal signaling footprints in cancer gene expression. *Nature Communications* 2017 9:1 (2018) 9:1–11. doi: [10.1038/s41467-017-02391-6](https://doi.org/10.1038/s41467-017-02391-6)
20. Badia-I-Mompel P, Vélez Santiago J, Braunger J, Geiss C, Dimitrov D, Müller-Dott S, Taus P, Dugourd A, Holland CH, Ramirez Flores RO, et al. decoupleR: ensemble of computational methods to infer biological activities from omics data. *Bioinformatics Advances* (2022) 2: doi: [10.1093/BIOADV/VBAC016](https://doi.org/10.1093/BIOADV/VBAC016)
21. Gavai AK, Supandi F, Hettling H, Murrell P, Leunissen JAM, Van Beek JHGM. Using Bioconductor Package BiGGR for Metabolic Flux Estimation Based on Gene Expression Changes in Brain. *PLOS ONE* (2015) 10:e0119016. doi: [10.1371/JOURNAL.PONE.0119016](https://doi.org/10.1371/JOURNAL.PONE.0119016)
22. King ZA, Lu J, Dräger A, Miller P, Federowicz S, Lerman JA, Ebrahim A, Palsson BO, Lewis NE. BiGG Models: A platform for integrating, standardizing and sharing genome-scale models. *Nucleic Acids Research* (2016) 44:D515–D522. doi: [10.1093/NAR/GKV1049](https://doi.org/10.1093/NAR/GKV1049)

23. Friedman J, Hastie T, Tibshirani R. Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software* (2010) 33:1–22. doi: [10.18637/jss.v033.i01](https://doi.org/10.18637/jss.v033.i01)
24. Therneau TM, Grambsch PM. *Modeling Survival Data: Extending the Cox Model*. 1st ed. New York: Springer Verlag (2000).
25. Kassambara A, Kosinski M, Biecek P. survminer: Drawing Survival Curves using 'ggplot2'. (2016) <https://cran.r-project.org/package=survminer>
26. Wilke CO. *Fundamentals of Data Visualization: A Primer on Making Informative and Compelling Figures*. 1st ed. Sebastopol: O'Reilly Media (2019).
27. Gohel D. flextable: Functions for Tabular Reporting. (2022) <https://cran.r-project.org/web/packages/flextable/index.html>
28. Liu J, Lichtenberg T, Hoadley KA, Poisson LM, Lazar AJ, Cherniack AD, Kovatich AJ, Benz CC, Levine DA, Lee AV, et al. An Integrated TCGA Pan-Cancer Clinical Data Resource to Drive High-Quality Survival Outcome Analytics. *Cell* (2018) 173:400–416.e11. doi: [10.1016/J.CELL.2018.02.052](https://doi.org/10.1016/J.CELL.2018.02.052)
29. Abeshouse A, Ahn J, Akbani R, Ally A, Amin S, Andry CD, Annala M, Aprikian A, Armenia J, Arora A, et al. The Molecular Taxonomy of Primary Prostate Cancer. *Cell* (2015) 163:1011–1025. doi: [10.1016/j.cell.2015.10.025](https://doi.org/10.1016/j.cell.2015.10.025)
30. Gerhauser C, Favero F, Risch T, Simon R, Feuerbach L, Assenov Y, Heckmann D, Sidiropoulos N, Waszak SM, Hübschmann D, et al. Molecular Evolution of Early-Onset Prostate Cancer Identifies Molecular Risk Markers and Clinical Trajectories. *Cancer cell* (2018) 34:996–1011.e8. doi: [10.1016/J.CCELL.2018.10.016](https://doi.org/10.1016/J.CCELL.2018.10.016)
31. Long Q, Xu J, Osunkoya AO, Sannigrahi S, Johnson BA, Zhou W, Gillespie T, Park JY, Nam RK, Sugar L, et al. Global transcriptome analysis of formalin-fixed prostate cancer specimens identifies biomarkers of disease recurrence. *Cancer research* (2014) 74:3228–3237. doi: [10.1158/0008-5472.CAN-13-2699](https://doi.org/10.1158/0008-5472.CAN-13-2699)
32. Ross-Adams H, Lamb A, Dunning M, Halim S, Lindberg J, Massie C, Egevad L, Russell R, Ramos-Montoya A, Vowler S, et al. Integration of copy number and transcriptomics provides risk stratification in prostate cancer: A discovery and validation cohort study. *EBioMedicine* (2015) 2:1133–1144. doi: [10.1016/j.ebiom.2015.07.017](https://doi.org/10.1016/j.ebiom.2015.07.017)
33. Schimmelpfennig C, Rade M, Füssel S, Löffler D, Blumert C, Bertram C, Borkowetz A, Otto DJ, Puppel SH, Hönscheid P, et al. Characterization and evaluation of gene fusions as a measure of genetic instability and disease prognosis in prostate cancer. *BMC cancer* (2023) 23: doi: [10.1186/S12885-023-11019-6](https://doi.org/10.1186/S12885-023-11019-6)



34. Kocher F, Tymoszuk P, Amann A, Sprung S, Salcher S, Daum S, Haybaeck J, Rinnerthaler G, Huemer F, Kauffmann-Guerrero D, et al. Deregulated glutamate to pro-collagen conversion is associated with adverse outcome in lung cancer and may be targeted by renin-angiotensin-aldosterone system (RAS) inhibition. *Lung Cancer* (2021) 159:84–95. doi: [10.1016/j.lungcan.2021.06.020](https://doi.org/10.1016/j.lungcan.2021.06.020)
35. Cohen J. Statistical Power Analysis for the Behavioral Sciences. *Statistical Power Analysis for the Behavioral Sciences* (2013) doi: [10.4324/9780203771587](https://doi.org/10.4324/9780203771587)
36. Benjamini Y, Hochberg Y. Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing. *Journal of the Royal Statistical Society: Series B (Methodological)* (1995) 57:289–300. doi: [10.1111/j.2517-6161.1995.tb02031.x](https://doi.org/10.1111/j.2517-6161.1995.tb02031.x)
37. Rousseeuw PJ. Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *Journal of Computational and Applied Mathematics* (1987) 20:53–65. doi: [10.1016/0377-0427\(87\)90125-7](https://doi.org/10.1016/0377-0427(87)90125-7)
38. Venna J, Kaski S. Neighborhood preservation in nonlinear projection methods: An experimental study. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)* (2001) 2130:485–491. doi: [10.1007/3-540-44668-0\\_68](https://doi.org/10.1007/3-540-44668-0_68)
39. Lange T, Roth V, Braun ML, Buhmann JM. Stability-based validation of clustering solutions. *Neural Computation* (2004) 16:1299–1323. doi: [10.1162/089976604773717621](https://doi.org/10.1162/089976604773717621)
40. Leng M, Wang J, Cheng J, Zhou H, Chen X. Adaptive semi-supervised clustering algorithm with label propagation. *Journal of Software Engineering* (2014) 8:14–22. doi: [10.3923/jse.2014.14.22](https://doi.org/10.3923/jse.2014.14.22)
41. Pichler R, Siska PJ, Tymoszuk P, Martowicz A, Untergasser G, Mayr R, Weber F, Seeber A, Kocher F, Barth DA, et al. A chemokine network of T cell exhaustion and metabolic reprogramming in renal cell carcinoma. *Frontiers in Immunology* (2023) 14:1208. doi: [10.3389/FIMMU.2023.1095195/BIBTEX](https://doi.org/10.3389/FIMMU.2023.1095195/BIBTEX)
42. 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. doi: [10.18637/JSS.V039.I05](https://doi.org/10.18637/JSS.V039.I05)
43. Zou H, Hastie T. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society Series B: Statistical Methodology* (2005) 67:301–320. doi: [10.1111/j.1467-9868.2005.00503.x](https://doi.org/10.1111/j.1467-9868.2005.00503.x)
44. Grambsch PM, Therneau TM. Proportional Hazards Tests and Diagnostics Based on Weighted Residuals. *Biometrika* (1994) 81:515. doi: [10.2307/2337123](https://doi.org/10.2307/2337123)

45. Harrell FE, Lee KL, Mark DB. Multivariable prognostic models: Issues in developing models, evaluating assumptions and adequacy, and measuring and reducing errors. *Statistics in Medicine* (1996) 15:361–387. doi: [10.1002/\(SICI\)1097-0258\(19960229\)15:4<361::AID-SIM168>3.0.CO;2-4](https://doi.org/10.1002/(SICI)1097-0258(19960229)15:4<361::AID-SIM168>3.0.CO;2-4)
46. Graf E, Schmoor C, Sauerbrei W, Schumacher M. Assessment and comparison of prognostic classification schemes for survival data. *Statistics in Medicine* (1999) 18:2529–2545. doi: [10.1002/\(sici\)1097-0258\(19990915/30\)18:17/18<2529::aid-sim274>3.0.co;2-5](https://doi.org/10.1002/(sici)1097-0258(19990915/30)18:17/18<2529::aid-sim274>3.0.co;2-5)