acta medica

ORIGINAL ARTICLE

# Evaluation of Prognostic Markers in Cancer-associated Fibroblast Based Sub-groups of Colorectal Cancer

Seçil Demirkol Canlı
ORCID: 0000-0003-0200-7962

∼◌◌◑ A B S T R A C T ◐◌◌∼

Objective: Stromal cells in the tumor microenvironment (TME) are among the main players of carcinogenesis. Cancer-associated fibroblasts (CAFs) residing in tumor stroma are involved in cancer progression through various mechanisms, supporting tumor growth, cellular motility and invasiveness. The discovery of markers predicting recurrence risk in colorectal cancer (CRC) have led to the generation of several gene panels, including Coloprint. This study aimed to understand the impact of a CAF-rich and a CAF-poor TME on the performance of the prognostic markers in Coloprint.

Materials and Methods: Publicly available transcriptomic data of CRC tumors were used to generate tumor sub-groups based on CAF specific gene expression. Subsequently, prognostic relationships of Coloprint genes were assessed within these subgroups.

Results and Conclusion: Our data revealed that prognostic performance of Coloprint genes differed dramatically between CAF stratified subgroups compared to non-stratified analysis. We have found that multiple genes lost their prognostic significance and several genes showed an association in the opposite direction. 9 out of 17 genes were differentially expressed in at least one of the CAF-specific subgroups and majority of the genes predicted prognosis independent of CAF levels. These findings showed that the performance of the prognostic markers can vary significantly among CAF-poor and CAF-rich groups. Therefore testing potential biomarkers within such biological sub-groups may contribute to the development of more specific gene panels.

Keywords: Colorectal cancer, Coloprint, prognosis, biomarker, fibroblast, cancer- associated fibroblast

Corresponding Author: Seçil Demirkol Canlı
Molecular Pathology Application and Research Center, Hacettepe University, Ankara, Turkey.
E-mail: secil.demirkol@hacettepe.edu.tr

## INTRODUCTION

Colorectal cancer (CRC) is the third most common cancer type and cause of cancer related death in both men and women in the USA [1]. According to the statistics of the Health Ministry of the Turkish Republic, in Turkey colorectal cancer is the fourth and second most common cancer type in men and women, respectively [2]. The 5-year survival of CRC patients is 64%, ranging from 90% for localized disease to 14% for advanced-stage disease [3]. Clinical factors defining poor prognosis for CRC include obstruction and perforation of colon at diagnosis, low number of assessed lymph nodes (<12), T4 stage, high grade, vascular, lymphatic or perineural invasion and residual tumor [4]. It is well known that diseases at stages II and III can show varying clinical outcomes upon treatment, therefore gene expression based molecular tests have been also developed to assess patients at risk of recurrence or patients who are unlikely to benefit from conventional therapy. To address this issue, Salazar et al. evaluated microarray based expression profile of 188 CRC patients at various

stages (I to IV), and identified a set of 18 genes that were associated with metastasis-free survival [5]. This signature, called Coloprint, was then validated using an independent set of 206 samples from patients with stage I, II, and III CRC [5]. Multivariate analyses showed that Coloprint was a strong independent prognostic factor in both stage II and III diseases, and was superior to the ASCO criteria for the evaluation of recurrence risk in stage II patients (HR=3.34; p=0.017) [5].

The crosstalk between various cell types in TME has a major impact on tumor progression [6]. Numerous studies have shown that a sub-population of fibroblasts, called cancer-associated fibroblasts (CAFs) or 'activated fibroblasts' in the tumor stroma are prominent promoters of tumor growth and progression [7]. During tumor progression, fibroblasts are activated by *TGFβ*, monocyte chemotactic protein 1 (*CCL2*), and extracellular matrix (ECM) degrading agents such as matrix metalloproteinases (MMPs) [8]. The activated CAFs in turn affect cellular motility by secreting various growth factors and cytokines. They are also an important source of MMPs, which degrade and remodel ECM, thus enhancing tumor growth, invasion, angiogenesis, recruitment of inflammatory cells, and metastasis [7,9,10]. In CRC, elevated CAF signature has been associated with poor disease-free survival in patients who did not receive adjuvant chemotherapy [11]. A CAF index was found to be even more powerful than an epithelial-mesenchymal transition (EMT) score in predicting survival outcomes in a pan-cancer cohort [12]. In line with that, transcriptomic based fibroblast scores were higher in the consensus molecular subtype 4 (CMS4) of CRC tumors, which have the worst prognosis [13].

CAF-rich and CAF-poor tumors show differences in molecular dynamics, prognosis and aggressiveness, however the evaluation of prognostic markers has not been previously addressed in CAF level stratified tumors. The lack of this knowledge prompted us to re-evaluate a panel of prognostic markers in CRC tumors with different levels of CAFs. In order to do that, hierarchical clustering analysis based on expression of six CAF specific markers was performed in colorectal tumors. Then, a previously published and validated prognostic gene panel, Coloprint, was re-tested in stage II and III CRC separately in the aforementioned CAF based sub-groups. Majority of these markers showed loss of significance and significance in opposite directions in prognostic analyses when stratified by CAF levels. Overall, our results suggest that gene panels developed for risk prediction may lose their power if tumors are further subdivided into subgroups with different biological features.

## MATERIALS AND METHODS

### Study cohorts and microarray data processing
CEL files of colorectal tumors within GSE39582 [14], GSE17536 [15] and GSE14333 [16] datasets were downloaded from GEO database (https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi) and RMA normalized using "affy" package in R Bioconductor [17]. Clinical data was obtained from Array Express (http://www.ebi.ac.uk/arrayexpress). Consensus molecular subtype information was downloaded from Synapse platform (www.synapse.org) for samples in GSE39582.

### MCP-counter and ESTIMATE
MCP-counter R package was downloaded from Zenodo (https://doi.org/10.5281/zenodo.61372). Based on probeset level expression data as input, scores for fibroblasts were obtained separately for GSE39582, GSE17536 and GSE14333. Stromal scores were obtained via ESTIMATE R package (https://bioinformatics.mdanderson.org/estimate/rpackage.html) for GSE39582 [18]. Samples were sorted from lowest to highest based on stromal scores and divided into three groups including 188, 189 and 189 samples with low, intermediate and high score, respectively.

### Hierarchical Clustering
Cluster 3.0 and Treeview programs were used for hierarchical clustering and visualization of heatmaps (http://bonsai.hgc.jp/~mdehoon/software/cluster/software.htm#ctv). Gene expression data was standardized to mean of zero and standard deviation of 1 for each gene. This data was then used as input for Cluster 3.0 software. Hierarchical clustering was performed using euclidian distance as similarity metric and complete linkage as clustering method. The output in ".cdt" format was used as input for Treeview software for visualization and heatmap generation.

## Survival Analyses

Univariate cox regression analysis was performed to evaluate prognostic relationships using continuous gene expression data. For genes in Coloprint with multiple probesets, one probeset was selected to be used throughout the study using the following criteria: 1. For each probeset log-rank tests were performed at all possible cut-offs within 10-90 percentiles using "survival" package in R Bioconductor 2. The lowest p value for each probeset was noted, which was named "best cut-off p value" 3. The probeset with the lowest "best cut-off p value" was selected and used for each gene.

For categorical analyses of survival, expression based cut-off with the lowest log-rank p value within 25-75 percentiles was preferred in order to define low and high expression groups. If there was no significance in any of the cut-offs in 25-75%, then the cut-off resulting in the lowest p value within 10-90% range was used. Log-rank p values smaller than 0.05 were considered significant. If no significant p value was obtained at any of the cut-offs tested within 10-90 percentiles, the gene was considered not significantly associated with prognosis. Patients with survival time "0" were excluded from all survival analyses.
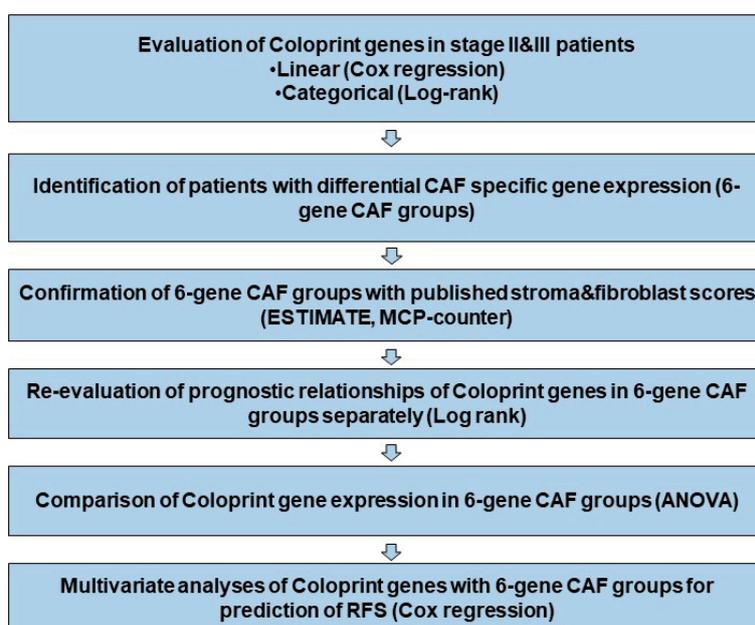
## Statistical Analyses

Gene expression plots comparing CAF groups were generated using "ggplot2" package [19], and ANOVA was performed using "oneway.test" function in R Bioconductor [20]. Pearson correlation analysis was performed using Microsoft Excel (2013) for the evaluation of intergenic correlations among CAF markers. Kaplan Meier graphs were generated using GraphPad Prism version 6 for Windows (Graphpad Software, San Diego, CA, USA) and cox regression analyses were performed using IBM SPSS Statistics for Windows, Version 23.0 (Armonk, NY, IBM Corp).

## RESULTS

### Relationship of gene expression and prognosis in stage II and III CRC

To assess the prognostic value of Coloprint genes in CAF-specific biological sub-groups, a transcriptomic based workflow was applied as summarized in Figure 1. Univariate cox regression analyses were performed for the genes included in the signature with recurrence-free survival (RFS) using microarray data from stage II and III CRC tumors in GSE39582 dataset (n=253 and n=200 with available recurrence-free survival data, respectively). 17 out of 18 genes in Coloprint were available in this dataset. Cox regression analyses showed that 7 and 2 genes were significantly associated with RFS in stage II and III patients, respectively (Table 1). However, expression values used in a continuous fashion did not show a significant prognostic relationship for the majority of the genes in both stage II and III disease (Table 1); therefore, the prognostic values were evaluated



**Figure 1.** A schematic of workflow of the study.

**Table 1.** Univariate cox regression analyses of 17 genes in stage II & III tumors (GSE39582, RFS).

| | Stage II | | | | Stage III | | | |
|---|---|---|---|---|---|---|---|---|
| | | | 95% CI for HR | | | | 95% CI for HR | |
| | p | HR | lower | upper | p | HR | lower | upper |
| CTSC | **<0.001** | 1.932 | 1.35 | 2.766 | 0.449 | 0.877 | 0.625 | 1.231 |
| CYFIP2 | 0.365 | 0.849 | 0.595 | 1.21 | 0.914 | 0.982 | 0.712 | 1.355 |
| EDEM1 | 0.344 | 0.747 | 0.409 | 1.365 | 0.43 | 0.825 | 0.512 | 1.33 |
| HSD3B1 | 0.079 | 3.376 | 0.868 | 13.127 | 0.131 | 2.556 | 0.757 | 8.63 |
| IL2RA | **0.001** | 2.126 | 1.386 | 3.261 | 0.202 | 0.76 | 0.498 | 1.159 |
| IL2RB | **0.032** | 1.439 | 1.032 | 2.006 | **0.002** | 0.58 | 0.407 | 0.825 |
| LAMA3 | 0.083 | 2.49 | 0.888 | 6.982 | 0.128 | 1.862 | 0.835 | 4.151 |
| LIF | **0.012** | 1.596 | 1.11 | 2.295 | 0.594 | 1.09 | 0.793 | 1.498 |
| MCTP1 | 0.211 | 0.82 | 0.6 | 1.12 | 0.377 | 0.88 | 0.662 | 1.169 |
| PIM3 | 0.401 | 1.288 | 0.713 | 2.327 | 0.816 | 1.064 | 0.632 | 1.789 |
| PLIN3 | 0.9 | 0.972 | 0.62 | 1.522 | 0.675 | 1.091 | 0.726 | 1.638 |
| PPARA | **0.001** | 0.515 | 0.352 | 0.754 | 0.356 | 0.837 | 0.574 | 1.221 |
| PYROXD1 | **0.011** | 0.755 | 0.609 | 0.937 | 0.063 | 0.786 | 0.611 | 1.013 |
| SLC6A11 | 0.072 | 2.313 | 0.927 | 5.774 | **0.044** | 2.137 | 1.021 | 4.475 |
| THNSL2 | **0.016** | 1.693 | 1.105 | 2.592 | 0.489 | 0.883 | 0.622 | 1.255 |
| ZBED4 | 0.1 | 0.664 | 0.408 | 1.081 | 0.184 | 0.778 | 0.537 | 1.127 |
| ZNF697 | 0.888 | 0.914 | 0.258 | 3.232 | 0.557 | 1.313 | 0.53 | 3.255 |

using expression data categorically.

When the prognostic relationship of each gene was re-evaluated by comparing prognosis of the tumor groups with high and low expression (see Methods) for stage II and stage III patients separately, a significant relationship with RFS was observed for 13 genes in stage II and 12 genes in stage III disease (Table 2). In a pooled analysis of 453 stage II and III patients, high expression of 7 genes (*CTSC, HSD3B1, IL2RA, LAMA3, LIF, SLC6A11, THNSL2*) and 6 genes (*CYFIP2, EDEM1, MCTP1, PPARA, PYROXD1, ZBED4*) were associated with shorter and longer RFS, respectively (Table 2). Significant relationships in the opposite directions were noted for *CTSC, IL2RA, IL2RB* and *THNSL2* genes. The expression of these genes were associated with shorter RFS in either stage II or pooled analyses, while they were associated with longer RFS in stage III disease (Table 2). In brief, multiple Coloprint genes did not show consistent prognostic associations in stage II and III patients in GSE39582.

## Sub-grouping method based on CAF specific gene expression

To define sub-groups based on CAF specific gene expression, six known CAF markers were used, *ATL1* [21], *PDGFRA, PDGFRB, FAP, ACTA2, S100A4* [22].

Based on intergenic correlations of all probesets of 6 CAF markers in 566 colorectal tumors from GSE39582 dataset (Table 3), the probeset with the highest mean Pearson r value was used in further analyses for genes with multiple probesets. Hierarchical clustering analyses performed separately for tumors in GSE39582 (n=566), GSE17536 (n=177) and GSE14333 (n=290) datasets showed 3 sub-groups with clear high, intermediate and low expression of markers consistently in 3 datasets (Figure 2, Supplementary Table 1). This classification was named "6-gene CAF groups" and used accordingly throughout the study. In line with these findings, mean expression of the six markers was significantly different among groups (Supplementary Figure 1).

In order to confirm that the groups represent enrichment of CAFs in the tumor microenvironment, a previously published algorithm called "MCP-counter" was used. This algorithm predicts the abundance of various cell types in the tumor microenvironment based on transcriptomic profiles, including fibroblasts [23]. Based on the MCP-counter algorithm, tumors were significantly infiltrated by fibroblasts (p<0.05) in GSE39582 (100% of all samples), GSE17536 (100% of all samples), and GSE14333 (99.3% of all samples) datasets. Among 6-gene CAF groups, CAF high

**Table 2.** Univariate cox regression analyses of 17 genes in stage II & III tumors (GSE39582, RFS).

| Probeset | Gene | Stage II&III (n=453) | | Stage II (n=253) | | Stage III (n=200) | |
|---|---|---|---|---|---|---|---|
| | | HR | P | HR | P | HR | P |
| 225646_at | CTSC | 1.5394 | 0.0118 | 2.5089 | 0.0006 | 0.5446 | 0.0118 |
| 215785_s_at | CYFIP2 | 0.4712 | 0.0255 | ng | ns | ng | ns |
| 203279_at | EDEM1 | 0.5615 | 0.0281 | ng | ns | 0.4707 | 0.0388 |
| 241111_at | HSD3B1 | 1.7473 | 0.0054 | 1.9253 | 0.0138 | 1.6385 | 0.0256 |
| 211269_s_at | IL2RA | 1.6919 | 0.0355 | 2.8776 | 0.0005 | 0.2973 | 0.0117 |
| 205291_at | IL2RB | ng | ns | 2.2127 | 0.0194 | 0.4095 | 0.0000 |
| 1568879_a_at | LAMA3 | 1.9139 | 0.0002 | 1.9588 | 0.0147 | 1.9082 | 0.0047 |
| 205266_at | LIF | 2.0601 | 0.0009 | 2.3729 | 0.0022 | 1.6982 | 0.0264 |
| 235740_at | MCTP1 | 0.6115 | 0.0299 | 0.3025 | 0.0326 | ng | ns |
| 224739_at | PIM3 | ng | ns | 2.0521 | 0.0176 | ng | ns |
| 202122_s_at | PLIN3 | ng | ns | ng | ns | ng | ns |
| 226978_at | PPARA | 0.4146 | 0.0000 | 0.3833 | 0.0002 | 0.5071 | 0.0225 |
| 213878_at | PYROXD1 | 0.4996 | 0.0000 | 0.3708 | 0.0001 | 0.5492 | 0.0076 |
| 230286_at | SLC6A11 | 2.3649 | 0.0002 | 1.9061 | 0.0167 | 1.9048 | 0.0068 |
| 219044_at | THNSL2 | 1.6922 | 0.0222 | 1.8598 | 0.0187 | 0.6014 | 0.0304 |
| 204799_at | ZBED4 | 0.5282 | 0.0011 | 0.5673 | 0.0483 | 0.6128 | 0.0414 |
| 1553702_at | ZNF697 | ng | ns | ng | ns | ng | ns |

ns: not significant
ng: HR was not given for nonsignificant relationships
Yellow and blue colors indicate relationships with poor (HR>1) and good (HR<1) prognosis, respectively.

**Table 3.** Intergenic correlation of CAF marker expression to select representative probesets (GSE39582, n=566).

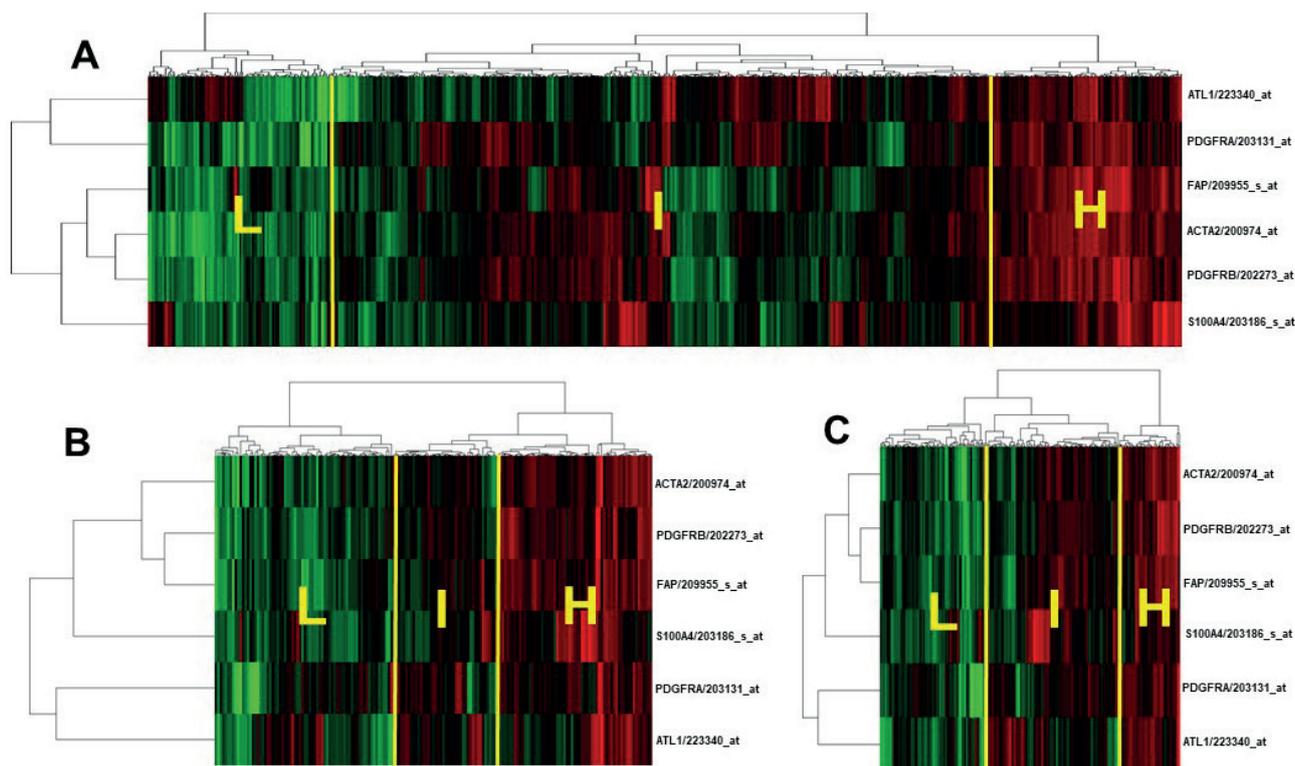| | | ACTA2 200974_at* | ACTA2 215787_at | ACTA2 243140_at | ATL1 223340_at* | FAP 209955_s_at* | PDGFRA 203131_at* | PDGFRA 211533_at | PDGFRA 215305_at | PDGFRA 1554828_at | PDGFRB 202273_at* | S100A4 203186_s_at* |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ACTA2 | **200974_at** | 1.000 | 0.058 | 0.567 | 0.380 | 0.634 | 0.628 | -0.090 | 0.284 | 0.167 | 0.806 | 0.482 |
| ACTA2 | 215787_at | 0.058 | 1.000 | 0.174 | -0.014 | 0.012 | 0.015 | -0.009 | 0.076 | 0.137 | 0.026 | -0.036 |
| ACTA2 | 243140_at | 0.567 | 0.174 | 1.000 | 0.209 | 0.417 | 0.373 | -0.036 | 0.204 | 0.161 | 0.486 | 0.294 |
| ATL1 | **223340_at** | 0.380 | -0.014 | 0.209 | 1.000 | 0.267 | 0.414 | -0.093 | -0.099 | -0.061 | 0.178 | 0.257 |
| FAP | **209955_s_at** | 0.634 | 0.012 | 0.417 | 0.267 | 1.000 | 0.419 | -0.157 | 0.225 | 0.030 | 0.700 | 0.493 |
| PDGFRA | **203131_at** | 0.628 | 0.015 | 0.373 | 0.414 | 0.419 | 1.000 | -0.127 | 0.337 | 0.045 | 0.569 | 0.249 |
| PDGFRA | 211533_at | -0.090 | -0.009 | -0.036 | -0.093 | -0.157 | -0.127 | 1.000 | -0.077 | 0.163 | -0.095 | -0.037 |
| PDGFRA | 215305_at | 0.284 | 0.076 | 0.204 | -0.099 | 0.225 | 0.337 | -0.077 | 1.000 | 0.140 | 0.389 | 0.187 |
| PDGFRA | 1554828_at | 0.167 | 0.137 | 0.161 | -0.061 | 0.030 | 0.045 | 0.163 | 0.140 | 1.000 | 0.165 | 0.132 |
| PDGFRB | **202273_at** | 0.806 | 0.026 | 0.486 | 0.178 | 0.700 | 0.569 | -0.095 | 0.389 | 0.165 | 1.000 | 0.491 |
| S100A4 | **203186_s_at** | 0.482 | -0.036 | 0.294 | 0.257 | 0.493 | 0.249 | -0.037 | 0.187 | 0.132 | 0.491 | 1.000 |
| | Mean r † | 0.447 | 0.131 | 0.350 | 0.222 | 0.367 | 0.357 | 0.040 | 0.242 | 0.189 | 0.429 | 0.319 |

r values are colored in red, white and green from highest to lowest, respectively
*Probesets used in further analyses. Probeset with the highest mean r was selected for genes with multiple probesets.
†Columnwise average r value

group had the highest MCP-counter fibroblast score and the score gradually and significantly decreased in intermediate and low groups (Supplementary Figure 2). Therefore, MCP-counter fibroblast scores were highly consistent with the 6-gene CAF groups. As seen in Table 4, tumor groups defined by another method, ESTIMATE algorithm, designed to predict presence of infiltrating stromal cells [18], overlaps to a large extent with 6-gene CAF groups. 91.3% of CAF high tumors had high stroma and 89.8% of CAF low tumors had low stroma scores by

ESTIMATE algorithm. These findings suggest that the 6-gene CAF groups identified in the current study were in line with tumor sub-groups defined by previously published CAF and stromal scoring methods and can clearly enable stratification of CRC tumors based on the level of CAFs in the tumor microenvironment. The 6-gene CAF groups also stratified patients with significantly different RFS, when evaluated in a pooled analysis of all stages and in stage II & III disease (Figure 3).

**Figure 2.** Hierarchical clustering analyses of colorectal tumors based on CAF markers. 566, 177 and 290 tumors were included in GSE39582 (A), GSE17536 (B), GSE14333 (C), respectively. «L», «I» and «H» letters indicate groups with low, intermediate and high CAF marker expression.

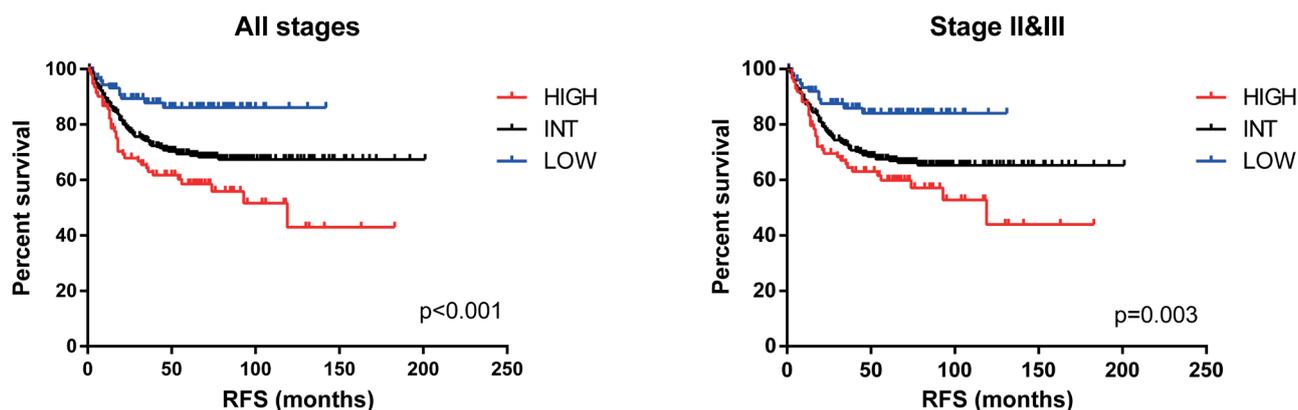**Table 4.** Distribution of ESTIMATE (stroma) and 6-gene CAF groups (GSE39582, n=566).

| | | 6 gene CAF groups | | | TOTAL |
| --- | --- | --- | --- | --- | --- |
| | | LOW | INT | HIGH | |
| ESTIMATE groups | LOW | 88 (89.8%) | 100 (27.4%) | 0 (0%) | 188 (33.2%) |
| | INT | 9 (9.2%) | 171 (46.8%) | 9 (8.7%) | 189 (33.4%) |
| | HIGH | 1 (1%) | 94 (25.8%) | 94 (91.3%) | 189 (33.4%) |
| | TOTAL | 98 (100%) | 305 (100%) | 103 (100%) | 566 (100%) |

Analysis of the distribution of CMS subtypes in CAF groups showed that 83.3% of the CAF high samples were CMS4 type (Supplementary Table 2), which was characterized as the mesenchymal subtype harboring prominent *TGFβ* activation, stromal invasion and angiogenesis [24]. The CAF intermediate group was heterogeneous and consisted of 20.2%, 54.7%, 12.4%, and 12.7% of the samples in CMS1, CMS2, CMS3 and CMS4 types, respectively. The CAF low group included only one CMS4 patient and the rest of the samples were distributed as 24.2% CMS1, 44.2% CMS2, and 30.5% CMS3. Overall, the CAF intermediate and low groups showed a heterogeneous distribution of CMS types, whereas CAF high group highly overlapped with tumors of CMS4 phenotype.

**Evaluation of prognostic genes in 6-gene CAF groups**

To re-assess the prognostic value of Coloprint genes within each 6-gene CAF groups, log-rank tests were performed based on the expression of each gene in a categorical fashion (see Methods) in stage II and III patients. As individual evaluation of prognostic relationships in each sub-group would reduce the sample sizes dramatically, this analysis was restricted to GSE39582 dataset, which has the highest number of samples with available survival data. Only three genes, *PPARA*, *PYROXD1* and *SLC6A11*, were significantly associated with RFS in all 6-gene CAF groups (Table 5). Hazard ratios (HR) indicated that high expressions of *PYROXD1* and *PPARA* were associated with longer RFS, and high expression of *SLC6A11* was associated with shorter RFS in all groups tested. Three genes, *CTSC, CYFIP2* and *ZNF697* were significant markers of RFS

## All stages



## Stage II&III



**Figure 3.** 6 gene CAF groups can predict RFS in GSE39582. Low, intermediate (INT) and high groups were assigned as defined in Figure 1. Log-rank p values are indicated.

in only the CAF high group. *EDEM1, MCTP1* genes were significantly related to longer RFS and high *THNSL2* expression was related to shorter RFS in CAF low and CAF intermediate groups, but not in CAF high, suggesting that high CAF content of the microenvironment might interfere with the prognostic role of these genes. The genes *ILR2A, LIF* and *ZBED4* had prognostic value in CAF low and CAF high groups, but not in CAF intermediate group, whereas *LAMA3* was significant in CAF intermediate and CAF high groups. *PIM3* was associated with unfavorable RFS in only CAF low group. *PLIN3* was the only gene that is not significantly related to

prognosis in any of the groups and in the pooled analyses. These data overall indicate that the prognostic relationships of most of these validated markers were highly heterogeneous when the tumors were stratified by CAF levels.

Interestingly, two genes, *HSD3B1* and *IL2RB* had contradictory relationships when evaluated in the 6-gene CAF groups separately. *HSNDB1* expression was significantly related to poor RFS in pooled analyses and the CAF intermediate group, whereas it was associated with good RFS in the CAF high group (Table 2&5). *IL2RB* was related to shorter RFS

**Table 5.** Log-rank based analyses of prognosis in 6-gene CAF groups (GSE39582, stage II&III, RFS).

| Gene | Low (n=77) | | Int (n=290) | | High (n=86) | |
|---|---|---|---|---|---|---|
| | HR | P | HR | P | HR | P |
| CTSC | ng | ns | ng | ns | 3.9437 | 0.0140 |
| CYFIP2 | ng | ns | ng | ns | 0.2539 | 0.0003 |
| EDEM1 | 0.1878 | 0.0169 | 0.574659 | 0.047502 | ng | ns |
| HSD3B1 | ng | ns | 1.6814 | 0.0189 | 0.3673 | 0.0034 |
| IL2RA | 6.7471 | 0.0352 | ng | ns | 2.2592 | 0.0491 |
| IL2RB | 3.3433 | 0.0416 | 0.5384 | 0.0040 | ng | ns |
| LAMA3 | ng | ns | 2.2358 | 0.0001 | 2.3384 | 0.0133 |
| LIF | 5.4958 | 0.0024 | ng | ns | 2.5422 | 0.0193 |
| MCTP1 | na* | 0.0374 | 0.626602 | 0.025783 | ng | ns |
| PIM3 | 3.4145 | 0.0328 | ng | ns | ng | ns |
| PLIN3 | ng | ns | ng | ns | ng | ns |
| PPARA | 0.3077 | 0.0409 | 0.5658 | 0.0475 | 0.4039 | 0.0123 |
| PYROXD1 | 0.2821 | 0.0330 | 0.5876 | 0.0131 | 0.3738 | 0.0028 |
| SLC6A11 | 3.7744 | 0.0193 | 2.0680 | 0.0106 | 4.3548 | 0.0078 |
| THNSL2 | 13.0820 | 0.0015 | 1.6135 | 0.0299 | ng | ns |
| ZBED4 | 0.1639 | 0.0011 | ng | ns | 0.4421 | 0.0168 |
| ZNF697 | ng | ns | ng | ns | 0.4241 | 0.0312 |

*Not available. Cox model resulted in an unrealistic HR due to lack of event in one of the groups
ns: not significant
ng: HR was not given for nonsignificant relationships
Yellow and blue colors indicate relationships with poor (HR>1) and good (HR<1) prognosis, respectively.
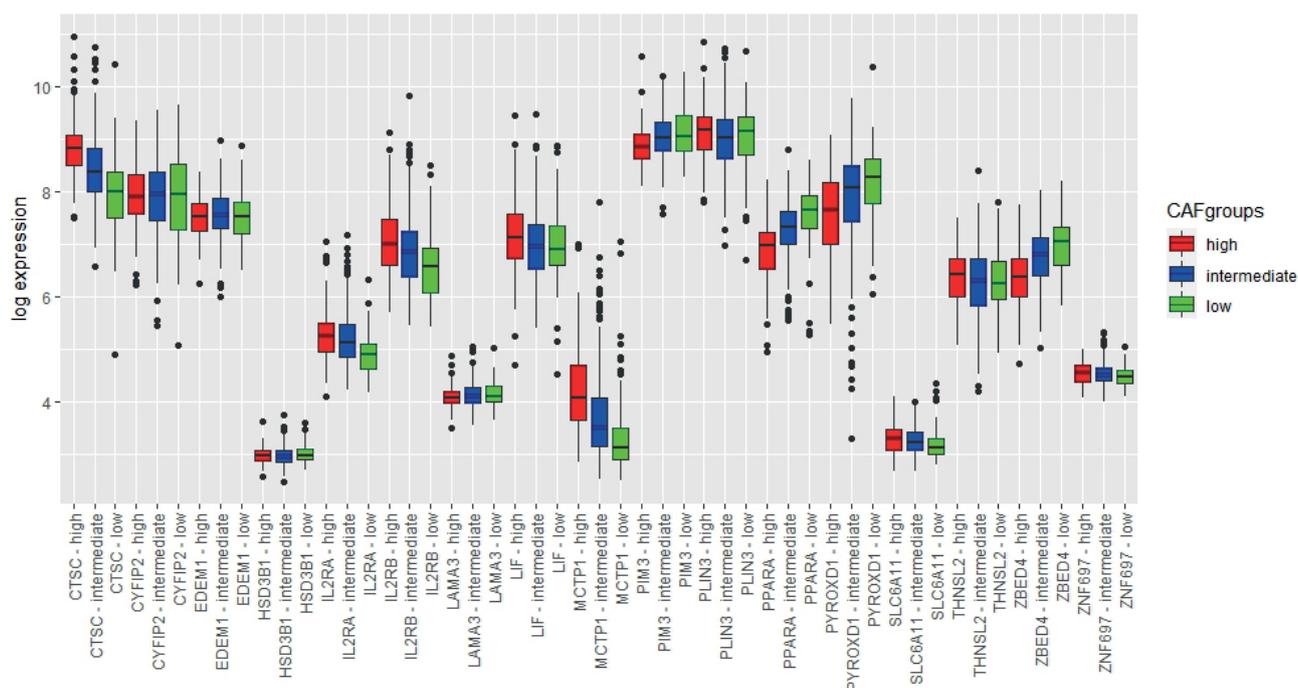
in CAF low, but longer RFS in the CAF intermediate group (Table 5). This type of opposite pattern was also observed in stage stratified analyses. High expression of *IL2RB* was associated with bad and good prognosis in stage II and stage III, respectively (Table 2). These findings were quite striking, since the prognostic relationships could show significant, but opposite patterns within 6-gene CAF groups, suggesting that CAF related changes in tumor microenvironment may have an effect on gene expression based prognostic predictions.

When expression of these genes between CAF groups were compared in stage II & III patients, 9 out of 17 genes (*CTSC, PPARA, ZBED4, MCTP1, IL2RA, IL2RB, PYROXD1, PIM3, SLC6A11*) showed significantly different expression in at least one CAF group (ANOVA p<0.05). Mean expression of *PPARA, ZBED4, PYROXD1, PIM3* decreased while mean expression of *CTSC, MCTP1, IL2RA, IL2RB, SLC6A11* increased gradually with increasing CAF level (Figure 4). Next, multivariate cox regression analyses (MVA) of 6-gene CAF groups were performed with each of the 13 genes which were significantly related to prognosis in a pooled analyses of stage II and III patients (Table 2). Our results indicated that 10 genes (*EDEM1, HSD3B1, LAMA3, LIF, MCTP1, PPARA, PYROXD1, SLC6A11, THNSL2, ZBED4*) were related to RFS independent

of 6-gene CAF groups while 3 genes (*CTSC, CYFIP2, IL2RA*) were not (Supplementary Table 3). *CTSC* and *IL2RA* genes were the only two genes, expression of which were elevated with increasing CAF levels and which lost significance in the multivariate cox model, indicating that the prognostic groups identified by these genes and 6-gene CAF groups were highly overlapping.

## DISCUSSION

Although it is known that CAFs are involved in tumor progression through secretion of various oncogenic signals and ECM-degrading proteases in the tumor microenvironment [25], how CAF involvement might affect the performance of putative prognostic or predictive biomarkers has not been elucidated so far. In this study, publicly available microarray data of CRC tumors were utilized and tumor sub-groups with low, intermediate and high CAF marker expressions were generated. This method enabled the identification of clear CAF sub-groups in three independent microarray datasets and these CAF groups were then confirmed with two previously published scoring methods for fibroblasts and tumor stroma. Therefore it is a fast, robust and practical way of obtaining CAF



**Figure 4.** Expression of Coloprint genes in 6-gene CAF groups. Boxes extend from 25th to 75th percentiles. The upper and lower whiskers extend 1.5 times the interquartile range above the upper quartile and below the lower quartile (Q1 - 1.5 * IQR or Q3 + 1.5 * IQR). Data beyond the end of the whiskers, outliers, are plotted individually.

related sub-groups based on transcriptomic data of CRC tumors. The significant gradual increase observed in mean expression of the 6 markers in 6-gene CAF groups, suggests that calculating mean expression of these markers could be noted as a practical and alternative approach to generate similar sub-groups. Our classification method may be also useful for transcriptomic data obtained by other technologies, such as RNA-seq and qRT-PCR, however further studies are needed to confirm the applicability.

Evaluation of gene expression-prognosis relationships showed that the biomarkers exhibit a heterogeneous pattern of significance in CAF groups. Only 3 genes *PPARA, PYROXD1* and *SLC6A11* out of 17 were significantly associated with clinical outcome in all 6-gene CAF groups. Furthermore, two genes *HSD3B1* and *IL1R3B* were significantly associated with RFS in the opposite directions in different CAF groups. The data further showed that CAFs are indeed a significant contributor to consensus molecular subtypes of colorectal cancer, as 83.3% of CAF high tumors were CMS4 type which is associated with an EMT phenotype and activation of matrix remodeling, angiogenesis and a gene expression profile compatible with stromal infiltration [24]. It is also known that CMS4 tumors have active *TGFβ* signaling [24], further supporting a CAF-rich microenvironment since *TGFβ*, released by cancer cells, is one of the key mediators of fibroblast activation [8]. Overall, these findings suggest that CAF levels may or should be considered as an important factor while evaluating putative markers as they contribute to significant changes in the tumor microenvironment that may affect the performance of biomarkers in prognostic panels.

In this study, the prognostic relationships were analyzed in a categorical way, via comparing high/low expression groups at all possible cut-offs. This approach enabled the identification of relationships that were weak to be significant in an analyses applied with the continuous log expression values. GSE39582 dataset was used for assessment of prognostic relationships in each 6-gene CAF group. There were 90, 336 and 93 patients with nonzero RFS and status information in CAF low, intermediate and high groups respectively. Therefore high sample size in this dataset enabled further dividing each 6-gene CAF group into two

groups based on expression of individual genes for prognostic comparison. Although the expression patterns of the six CAF markers were confirmed in two other independent datasets (GSE17536, GSE14333), these datasets were not utilized for prognostic analyses within CAF groups. GSE14333 had 44 samples in CAF high group with available clinical outcome data, and GSE17536 included 38 in CAF intermediate group. These sample sizes can be considered relatively low for prognostic comparisons, as categorical evaluations based on gene expression in these groups will lead to the comparison of data from only 15-20 patients to others. Thus, the categorical prognostic evaluations within each 6-gene CAF group were not performed in these datasets.

Upon evaluation of prognostic relationships of Coloprint genes, we noted clear changes in HR and p values for multiple genes when analyzed separately in 6-gene CAF groups. As Coloprint was developed based on Agilent oligonucleotide arrays [5], the platform-based differences such as the hybridization of probes to different transcript variants might have altered the direction and significance of prognostic relationships. In addition, differences in cohort-specific clinical characteristics may have had an effect on these inconsistencies.

The molecular function of genes in Coloprint included roles in cell proliferation, immune response, metabolism and cell invasion [5]. Among the genes involved in this signature, several genes have been previously linked to CAFs and CAF related molecular mechanisms. Laminin-332, an extracellular matrix (ECM) component composed of *LAMA3, LAMB3*, and *LAMC2* chains, was highly expressed in the tumor-normal interface. It was suggested that this may be a product of a paracrine intercellular reaction between invasive tumor cells in the tumor core and myofibroblasts in the tumor-normal interface to provide a suitable microenvironment for invasion in breast cancer [26]. Therefore, the fact that *LAMA3* was a significant predictor of prognosis in CAF high and CAF intermediate groups but not CAF low group, further supports that the role of this gene in prognostic prediction may rely on the presence of myofibroblasts in the microenvironment. Expression of *CTSC*, which positively correlated with CAF levels in our study, is expressed by fibroblasts and immune cells that mediates angiogenesis and growth of transplantable tumors [27]. In the current

study, *CTSC* was associated with RFS in only the CAF high group suggesting that *CTSC* expressed by CAFs may be relevant to its prognostic role. This is also in line with its lack of significance in a cox model including CAF levels in MVA. The expression of Leukemia inhibitory factor (*LIF*), which is secreted by both fibroblasts and neoplastic cells, is triggered by *TGFβ*. *LIF* is also involved in pro-invasive activation of stromal fibroblasts [28]. Although not significant, *LIF* expression was relatively higher in CAF high group, therefore it's likely that its role in the activation of fibroblasts might contribute to its prognostic associations.

Contradictory results in the direction and significance of multiple prognostic relationships upon individual evaluation of CAF groups, suggests that other prognostic markers or gene panels proposed in the literature may also show divergence in performance in CAF rich and CAF poor tumor microenvironments. Therefore it would be useful to take the involvement of CAFs into account while evaluating potential biomarkers and tumor sub-grouping gene panels.

## REFERENCES

[1] Siegel RL, Miller KD, Fuchs HE, Jemal A. Cancer Statistics, 2021. CA Cancer J Clin. 2021;71(1):7-33. https://doi.org/10.3322/caac.21654

[2] Akarsu N. AN, Aksoy M., Arslanhan S., Atabey E., Güllü İ., et al. Türkiye'de kanser kontrolü. In: Tuncer A, ed. Ankara: T.C. Sağlık Bakanlığı Kanserle Savaş Dairesi Başkanlığı; 2009:45-46.

[3] Siegel RL, Miller KD, Goding Sauer A, et al. Colorectal cancer statistics, 2020. CA Cancer J Clin. 2020;70(3):145-164. https://doi.org/10.3322/caac.21601

[4] Benson AB, 3rd, Schrag D, Somerfield MR, et al. American Society of Clinical Oncology recommendations on adjuvant chemotherapy for stage II colon cancer. J Clin Oncol. 2004;22(16):3408-3419. https://doi.org/10.1200/JCO.2004.05.063

[5] Salazar R, Roepman P, Capella G, et al. Gene expression signature to improve prognosis prediction of stage II and III colorectal cancer. J Clin Oncol. 2011;29(1):17-24. https://doi.org/10.1200/JCO.2010.30.1077

[6] Liotta LA, Kohn EC. The microenvironment of the tumour-host interface. Nature. 2001;411(6835):375-379. https://doi.org/10.1038/35077241

[7] Mueller MM, Fusenig NE. Friends or foes - bipolar effects of the tumour stroma in cancer. Nat Rev Cancer. 2004;4(11):839-849. https://doi.org/10.1038/nrc1477

[8] Bremnes RM, Donnem T, Al-Saad S, et al. The role of tumor stroma in cancer progression and prognosis: emphasis on carcinoma-associated fibroblasts and non-small cell lung cancer. J Thorac Oncol. 2011;6(1):209-217. https://doi.org/10.1097/JTO.0b013e3181f8a1bd

[9] Kessenbrock K, Plaks V, Werb Z. Matrix metalloproteinases: regulators of the tumor microenvironment. Cell. 2010;141(1):52-67. https://doi.org/10.1016/j.cell.2010.03.015

[10] Boire A, Covic L, Agarwal A, Jacques S, Sherifi S, Kuliopulos A. PAR1 is a matrix metalloprotease-1 receptor that promotes invasion and tumorigenesis of breast cancer cells. Cell. 2005;120(3):303-313. https://doi.org/10.1016/j.cell.2004.12.018

[11] Isella C, Terrasi A, Bellomo SE, et al. Stromal contribution to the colorectal cancer transcriptome. Nat Genet. 2015;47(4):312-319. https://doi.org/10.1038/ng.3224

[12] Ko YC, Lai TY, Hsu SC, et al. Index of Cancer-Associated Fibroblasts Is Superior to the Epithelial-Mesenchymal Transition Score in Prognosis Prediction. Cancers (Basel). 2020;12(7). https://doi.org/10.3390/cancers12071718

[13] Becht E, de Reynies A, Giraldo NA, et al. Immune and Stromal Classification of Colorectal Cancer Is Associated with Molecular Subtypes and Relevant for Precision Immunotherapy. Clin Cancer Res. 2016;22(16):4057-4066. https://doi.org/10.1158/1078-0432.CCR-15-2879

[14] Marisa L, de Reynies A, Duval A, et al. Gene expression classification of colon cancer into molecular subtypes: characterization, validation, and prognostic value. PLoS Med. 2013;10(5):e1001453. https://doi.org/10.1371/journal.pmed.1001453

[15] Smith JJ, Deane NG, Wu F, et al. Experimentally derived metastasis gene expression profile predicts recurrence and death in patients with colon cancer. Gastroenterology. 2010;138(3):958-968. https://doi.org/10.1053/j.gastro.2009.11.005

[16] Jorissen RN, Gibbs P, Christie M, et al. Metastasis-Associated Gene Expression Changes Predict Poor Outcomes in Patients with Dukes Stage B and C Colorectal Cancer. Clin Cancer Res. 2009;15(24):7642-7651. https://doi.org/10.1158/1078-0432.CCR-09-1431

[17] Gautier L, Cope L, Bolstad BM, Irizarry RA. affy--analysis of Affymetrix GeneChip data at the probe level. Bioinformatics. 2004;20(3):307-315. https://doi.org/10.1093/bioinformatics/btg405

[18] Yoshihara K, Shahmoradgoli M, Martinez E, et al. Inferring tumour purity and stromal and immune cell admixture from expression data. Nature Communications. 2013;4:2612. https://doi.org/10.1038/ncomms3612

[19] ggplot2: Elegant Graphics for Data Analysis [computer program]. New York: Springer-Verlag; 2016.

[20] R: A language and environment for statistical computing [computer program]. Vienna, Austria: R Foundation for Statistical Computing; 2019.

[21] Son GM, Kwon MS, Shin DH, Shin N, Ryu D, Kang CD. Comparisons of cancer-associated fibroblasts in the intratumoral stroma and invasive front in colorectal cancer. Medicine (Baltimore). 2019;98(18):e15164. https://doi.org/10.1097/MD.0000000000015164

[22] Liu T, Han C, Wang S, et al. Cancer-associated fibroblasts: an emerging target of anti-cancer immunotherapy. J Hematol Oncol. 2019;12(1):86. https://doi.org/10.1186/s13045-019-0770-1

[23] Becht E, Giraldo NA, Lacroix L, et al. Estimating the population abundance of tissue-infiltrating immune and stromal cell populations using gene expression. Genome Biol. 2016;17(1):218. https://doi.org/10.1186/s13059-016-1070-5

[24] Guinney J, Dienstmann R, Wang X, et al. The consensus molecular subtypes of colorectal cancer. Nat Med. 2015;21(11):1350-1356. https://doi.org/10.1038/nm.3967

[25] Kalluri R, Zeisberg M. Fibroblasts in cancer. Nat Rev Cancer. 2006;6(5):392-401. https://doi.org/10.1038/nrc1877

[26] Kim BG, An HJ, Kang S, et al. Laminin-332-rich tumor microenvironment for tumor invasion in the interface zone of breast cancer. Am J Pathol. 2011;178(1):373-381. https://doi.org/10.1016/j.ajpath.2010.11.028

[27] Ruffell B, Affara NI, Cottone L, et al. Cathepsin C is a tissue-specific regulator of squamous carcinogenesis. Genes Dev. 2013;27(19):2086-2098. https://doi.org/10.1101/gad.224899.113

[28] Albrengues J, Bourget I, Pons C, et al. LIF mediates proinvasive activation of stromal fibroblasts in cancer. Cell Rep. 2014;7(5):1664-1678. https://doi.org/10.1016/j.celrep.2014.04.036