

# A Digital Pathology-Based Shotgun-Proteomics Approach to Biomarker Discovery in Colorectal Cancer

Stefan Zahnd<sup>1</sup>, Sophie Braga-Lagache<sup>2</sup>, Natasha Buchs<sup>2</sup>, Alessandro Lugli<sup>1</sup>, Heather Dawson<sup>1</sup>, Manfred Heller<sup>2</sup>, Inti Zlobec<sup>1</sup>

<sup>1</sup>Institute of Pathology, University of Bern, Bern, Switzerland, <sup>2</sup>Department for BioMedical Research, Proteomics and Mass Spectrometry Core Facility, University of Bern, Bern, Switzerland

Received: 31 August 2018

Accepted: 21 February 2019

Published: 12 December 2019

## Abstract

**Background:** Biomarkers in colorectal cancer are scarce, especially for patients with Stage 2 disease. The aim of our study was to identify potential prognostic biomarkers from colorectal cancers using a novel combination of approaches, whereby digital pathology is coupled to shotgun proteomics followed by validation of candidates by immunohistochemistry (IHC) using digital image analysis (DIA). **Methods and Results:** Tissue cores were punched from formalin-fixed paraffin-embedded colorectal cancers from patients with Stage 2 and 3 disease ( $n = 26$ , each). Protein extraction and liquid chromatography-mass spectrometry (MS) followed by analysis using three different methods were performed. Fold changes were evaluated. The candidate biomarker was validated by IHC on a series of 413 colorectal cancers from surgically treated patients using a next-generation tissue microarray. DIA was performed by using a pan-cytokeratin serial alignment and quantifying staining within the tumor and normal tissue epithelium. Analysis was done in QuPath and Brightness\_Max scores were used for statistical analysis and clinicopathological associations. MS identified 1947 proteins with at least two unique peptides. To reinforce the validity of the biomarker candidates, only proteins showing a significant ( $P < 0.05$ ) fold-change using all three analysis methods were considered. Eight were identified, and of these, cathepsin B was selected for further validation. DIA revealed strong associations between higher cathepsin B expression and less aggressive tumor features, including tumor node metastasis stage and lymphatic vessel and venous vessel invasion ( $P < 0.001$ , all). Cathepsin B was associated with more favorable survival in univariate analysis only. **Conclusions:** Our results present a novel approach to biomarker discovery that includes MS and digital pathology. Cathepsin B expression analyzed by DIA within the tumor epithelial compartment was identified as a strong feature of less aggressive tumor behavior and favorable outcome, a finding that should be further investigated on a more functional level.

**Keywords:** Biomarker discovery, colorectal cancer, digital image analysis, digital pathology, mass spectrometry

## BACKGROUND

An important issue in biomarker research is the intra- and inter-tumor heterogeneity of the composition of tumors.<sup>[1,2]</sup> Tumors arise and grow in a very complex environment of different cell types, including primary tumor cells, stromal cells, and immune cells, among others. All of these different cell populations contribute to varying degrees to any protein, DNA, or RNA signal obtained as a result of studies working on patient material. This heterogeneity has been proposed as a significant challenge to successful cancer biomarker research.<sup>[3]</sup> Consequently, multiple authors have demonstrated the importance of enrichment to obtain specific signals for unbiased biomarker discovery.<sup>[4-6]</sup>

The lack of biomarkers remains an urgent issue in colorectal cancer. To this day, the tumor node metastasis (TNM) staging

system remains the “gold standard” for tumor classification in this disease, which assigns patients into one of four categories with unique treatment regimens. The lack of biomarkers is most prominent in Stage 2 and 3 patients. While Stage 2 patients are usually surgically treated, Stage 3 patients are readily administered adjuvant therapies. These differences in

**Address for correspondence:** Prof. Inti Zlobec,  
Institute of Pathology, University of Bern, Murtenstrasse 31,  
3008 Bern, Switzerland.  
E-mail: [inti.zlobec@pathology.unibe.ch](mailto:inti.zlobec@pathology.unibe.ch)

This is an open access journal, and articles are distributed under the terms of the Creative Commons Attribution-NonCommercial-ShareAlike 4.0 License, which allows others to remix, tweak, and build upon the work non-commercially, as long as appropriate credit is given and the new creations are licensed under the identical terms.

**For reprints contact:** [reprints@medknow.com](mailto:reprints@medknow.com)

**How to cite this article:** Zahnd S, Braga-Lagache S, Buchs N, Lugli A, Dawson H, Heller M, *et al.* A digital pathology-based shotgun-proteomics approach to biomarker discovery in colorectal cancer. *J Pathol Inform* 2019;10:40.

Available FREE in open access from: <http://www.jpathinformatics.org/text.asp?2019/10/1/40/272776>

### Access this article online

#### Quick Response Code:



**Website:**  
[www.jpathinformatics.org](http://www.jpathinformatics.org)

**DOI:**  
10.4103/jpi.jpi\_65\_18

treatment regimen are based predominantly on the absence or presence, respectively, of lymph node metastases. The need for additional biomarkers for improved patient stratification and management is found in the conflicting finding that subgroups of Stage 3 patients (Stage 3A) are characterized by more favorable prognosis than subgroups of Stage 2 patients (Stage 2B and 2C).

Among the available methods for biomarker discovery, mass spectrometry (MS) certainly ranks among the methods with the highest accuracy and throughput. The typical application for mass spectrometry in protein biomarkers discovery refers to an untargeted, or “shotgun,” approach. This procedure involves extraction and denaturation of proteins from frozen or formalin-fixed paraffin-embedded (FFPE) tissue into solution, followed by enzymatic digestion with specific proteases such as trypsin. The resulting peptide solution is chromatographically separated, ionized, and injected for analysis into an MS instrument in an approach known as liquid chromatography-based MS [LC-MS]. Because of the principle of injecting a solution of fragmented proteins into an MS instrument, this approach has also been coined “shotgun proteomics.”

MS represents a well-suited method for biomarker discovery, but additional methods are required for the subsequent validation of identified candidates. In that respect, the combination of digital image analysis (DIA) and immunohistochemistry (IHC) has recently shown very promising results. For example, the Ki67 marker has been used for years as a proliferative marker in breast cancer but is limited by significant interobserver variability. Recently, however, Stålhammar *et al.*<sup>[7]</sup> have used DIA for more reliable and objective scoring for Ki67. The potential of DIA can be further increased in combination with a next-generation tissue microarray<sup>®</sup> (ngTMA) approach that allows for targeted regions of tissue to be included in digital analysis.<sup>[8-10]</sup>

The aim of this study is the discovery of novel potential biomarkers for patients with Stage 2 colorectal cancer. This aim is achieved by developing a pipeline for biomarker discovery that begins by enriching tumor samples for MS analysis using ngTMA technology, followed by validation of candidate proteins using quantitative DIA. This combination of methodologies represents a new milestone for the identification and validation of biomarkers in cancer research.

## METHODS

### Patients

A total of 413 surgically treated patients with primary colorectal cancer were entered into this study. Patients were treated at the University Hospital of Bern between 2002 and 2013. Patient characteristics are found in Supplementary Table 1. All diagnostic slides were re-reviewed by expert gastrointestinal pathologists, and the following histopathological information was recorded: tumor grade, lymphatic and venous invasion, pathological T, N, and M stage, overall TNM stage, perineural

invasion, tumor border configuration, tumor budding score assessed according to the ITBCC guidelines,<sup>[11]</sup> microsatellite status, and tumor histology. Clinical data included age at diagnosis, gender, survival time, and therapy. No cases were preoperatively treated. Median overall survival time for the full cohort was 42.7 months.

### Next-generation tissue microarray construction

Blocks were retrieved from the archives of the Institute of Pathology, University of Bern, Switzerland. From each block, an H and E slide was created and scanned (digitalized) on a Panoramic 250 scanner (3DHitech, Budapest, Hungary). On each digital slide created this way, multiple circular digital annotations (i.e., regions of interest [ROI]) with a diameter of 0.6 and 1.0 mm were annotated in histological tumor regions of all patients in the cohort. Digital slides were aligned with their corresponding blocks, and annotated regions were extracted from the block and placed in a new, empty paraffin block. This resulted in ngTMA<sup>®</sup> featuring multiple cores per patient in the validation cohort [Supplementary Figure 1]. The validation cohort ngTMA<sup>®</sup> contained in a total of 796 cores. The use of clinical data and tissue samples was approved by the Ethics Committee of the Canton of Bern (KEK 2014/200).

### Protein extraction and mass spectrometry

From these 413 patients, 26 Stage 2 and 26 Stage 3 cases were selected for MS. Patients were selected to represent a balanced cohort without missing data for any feature including survival. Using the same technology as described above for TMAs, digital annotations were used to target regions of tumor epithelium. Blocks were aligned with the digital scans, and four cores from each patient were collected in 2 ml tubes. Extracted cores were re-embedded in a 7 mm × 7 mm paraffin block, 20 sections (15 μm each) from each block were cut, and protein was extracted according to the “Extraction of Total Protein from FFPE Tissue Sections” protocol of the Qproteome<sup>®</sup> FFPE Tissue Handbook (p. 18-19, Qiagen). 15 μl of the resulting protein extract was injected on an LC-MS/MS system (EASY nLC1000 liquid chromatograph coupled to a QExactive mass spectrometer, both from ThermoFisher Scientific). To reinforce the confidence in identified biomarkers, MS results were evaluated with one commercial (Proteome Discoverer version 2.2, ThermoFisher) and one freely available (MaxQuant version 1.5.8.3) software with the following settings: up to two missed cleavages in tryptic digest peptides, carbamidomethylation on cysteine as fixed modification and N-terminal acetylation and methionine oxidation as variable modifications. A mass tolerance of 10 ppm was used for precursor ions and 20 ppm for fragment ions. False discovery rate on protein level was set to 0.01 (1%), with a minimum of two unique peptides required for protein identification.

After result evaluation, a python script was specifically developed to extract intensities of all identified proteins from both software packages by summing up individual intensities of the three most abundant peptide ions. This approach is

known as “Top 3” and has been previously shown to be the most accurate approximation of abundance in label-free protein quantification.<sup>[12,13]</sup> MaxQuant-calculated “label-free quantification” (LFQ) scores were also included in the statistical evaluation. The three calculated quantification methods (LFQ and Top 3 scores from MaxQuant, Top 3 score from Proteome Discoverer) were imported into R version 3.4.2<sup>[14]</sup> and were used for calculating fold change and associated *P* values for all proteins identified with at least two unique peptides. All proteins exhibiting significant (*P* < 0.05) fold changes in all three scores were considered as true biomarkers. The cathepsin B protein was selected for subsequent validation on the previously constructed ngTMA<sup>®</sup> cohort.

### Immunohistochemistry

For IHC validation, two sequential slides of the validation ngTMA<sup>®</sup> were created and stained with cytokeratin (clone AE1/AE3, Cell Signaling) and cathepsin B (clone EPR4323, GeneTex), respectively, on the automated BOND RX<sup>®</sup> (Leica Biosystems, Newcastle, UK). Digital copies of slides were obtained at highest quality settings and ×40 resolution on a Panoramic 250 Scanner and stored on our institute servers. The intensity of cathepsin B was visually assessed and classified into one of two “visual intensity” categories, namely category 0 (low expression) or category 1 (high expression) based on reference intensities. These visual intensity categories were also combined with clinical, pathological data for inclusion in the statistical analysis.

### Digital image analysis using QuPath

For quantification with DIA, slides were imported into the QuPath software version 0.1.2,<sup>[15]</sup> core areas were detected using the built-in TMA dearrayer module, and the position of tissue cores was manually adjusted to ensure proper recognition and further processing of all cores on the slide. All cytokeratin-stained tissue cores were then evaluated for inclusion or exclusion in further data analysis. Tissue cores were excluded from analysis if (i) their shape was distorted such that the tissue region extended beyond the borders of the 0.85 mm circular annotation, (ii) they contained any amount of normal epithelial tissue, (iii) they showed clear indications that the staining protocol did not work, (iv) they contained overlapping tissue, or (v) they contained <25% of tissue. Cores for which observers were not in agreement were discussed after evaluation until agreement of both evaluators was reached.

To all tissue cores included in the analysis, QuPath’s simple linear iterative clustering superpixel segmentation (Gaussian sigma: 5 μm, superpixel spacing: 10 μm, iterations: 15, regularization: 0.1) was applied. This approach is very well suited for capturing the histological shape of tissue cores. QuPath intensity features (preferred pixel size: 2 μm, region: ROI, tile diameter: 25 μm, compute all features including Haralick features [Haralick *et al.*, 1973] with 32 bins) were calculated for all superpixels identified this way.

After calculation, annotations were placed on randomly chosen tissue cores of each cytokeratin slide in three

regions: epithelium, stroma, and whitespace. It was ensured that each annotation contained 3000 (±5%) superpixels. Annotated superpixels were used to build a Random Trees detection classifier for the entire TMA slide, and this approach ensured a balanced training set in the three selected regions. Superpixels belonging to the same class were then combined to a single annotation, resulting in three annotations (epithelium, stroma, and whitespace) on each tissue core of all slides.

After classification of superpixels, epithelial regions were manually transferred to corresponding cores of the sequential ngTMA<sup>®</sup> slide stained for cathepsin B. Corresponding cores were aligned, the epithelial annotation from the cytokeratin slide was manually transferred to the cathepsin B slide, and necrotic areas as well as intraglandular debris of the transferred annotation were manually removed before further analysis. Annotations were not transferred if the core on the cathepsin B slide (i) showed clear indications that the staining did not work, (ii) contained overlapping tissue, (iii) was distorted in a way that the core tissue region extended beyond the borders of the 0.85 mm circular annotation, or (iv) contained <25% of tissue. Upon successful transfer of all epithelial annotations to the candidate biomarker slide, the watershed cell detection method (detection image: optical density sum, requested pixel size: 0.5 μm, background radius: 8 μm, median filter radius: 0 μm, Sigma: 1.5 μm, minimum area: 10 μm<sup>2</sup>, maximum area: 400 μm<sup>2</sup>, threshold: 0.1, checked options: splitting by shape, include cell nucleus, smooth boundaries) was applied to all annotations on the candidate biomarker slide. Cell features were calculated for all cells using QuPath’s add intensity features option (preferred pixel size: 0.5 μm, region: ROI, tile diameter: 25 μm, compute all features including Haralick features with 32 bins). This approach calculates 163 parameters for each detected cell on the candidate biomarker slide, which were exported into a tab-delimited file using a QuPath script specifically developed for this purpose. Exported measurements from all cells on the biomarker candidate slide were combined with relevant clinical pathological features and imported into R version 3.4.2 for statistical analysis.

### Study design

Patients for biomarker discovery (26 patients from Stage 2 and 3, respectively) as well as for the construction of the validation cohort were chosen from an initial four hundred and thirteen CRC patients. The complete study design is shown in Figure 1.

### Statistical analysis

All statistical analysis was performed in R version 3.4.2. Differences in survival curves were assessed using the log-rank test. Biomarkers in the three scoring approaches were identified using the *t*-test. Only proteins with *P* < 0.05 in all three comparisons were considered significant. The association of visual intensity categories and Brightness\_Max with clinical-pathological features was performed using the Spearman correlation coefficient, Chi-square, or Wilcoxon’s rank-sum test.

## RESULTS

### Differentially expressed proteins between Stage 2 and 3 patients

Differential protein expression between Stage 2 and 3 patients revealed 1947 proteins identified with at least two unique peptides. The three types of quantification scores (LFQ scores from MaxQuant and Top 3 scores from MaxQuant and Proteome Discoverer) were combined using a python script developed specifically for this purpose and analyzed in R. To reinforce the validity of biomarker candidates, only proteins with a significant *P* value ( $P < 0.05$ ) fold change in all three scores were considered as true biomarkers [Figure 2a]. This resulted in eight candidates for subsequent validation [Figure 2b]. Based on the body of available literature, the cathepsin B protein was chosen for further validation on the validation cohort using DIA.

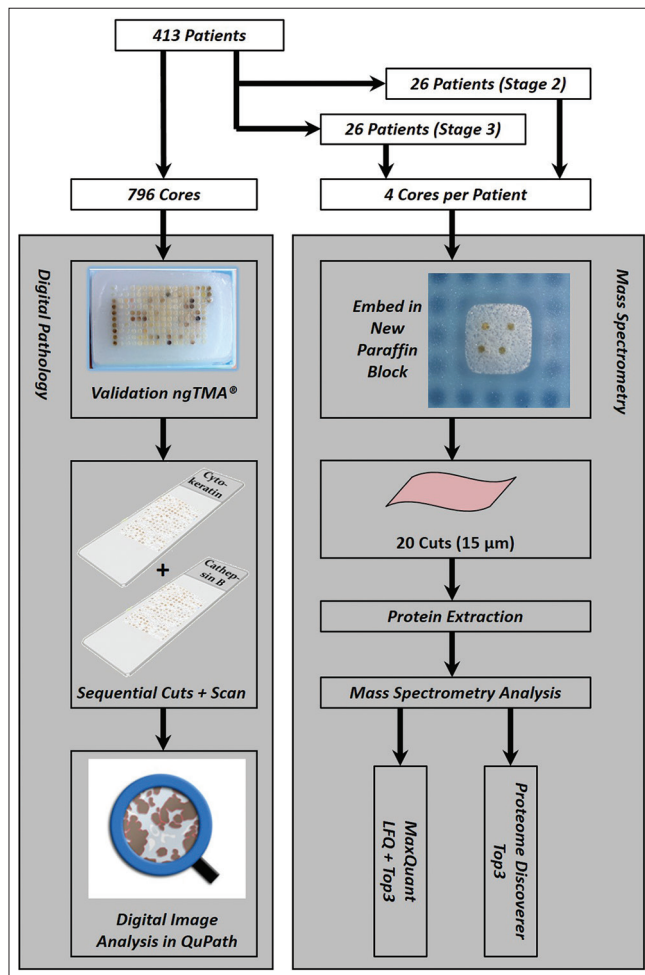
### Assessment of cathepsin B by visual assessment of next-generation tissue microarray

In a first step, the expression of cathepsin B in all patients in the validation cohort was visually assessed and assigned one

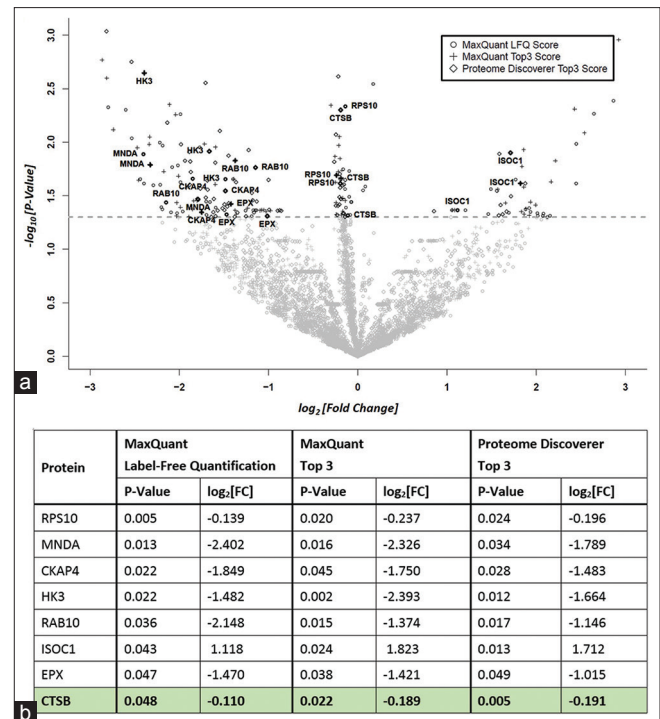
of two categories “low/no cathepsin B expression” or “high cathepsin B expression” based on reference images [Figure 3]. All patients in the cohort were then stratified for overall survival by this visual intensity category [Figure 4a]. However, these results did not indicate any significant stratification of patients (log-rank test,  $P > 0.05$ ). In addition, this categorization into expression categories did not reveal any significant associations with clinical-pathological features [Table 1]. Based on these findings, subgroup analysis for available Stage 2 patients was performed. Stratification of Stage 2 patients by visual intensity is shown in Figure 4b. Although not statistically significant, a trend in the stratification of patients can be observed.

### Digital image analysis of cathepsin B

Our categorization of cathepsin B by visual assessment did not reveal any significant stratification of patients. However, we hypothesized that DIA, which allows analysis of IHC expression on a continuous rather than a categorical scale, would improve the associations with clinical-pathological features and potentially yield a “digital” surrogate for the expression of cathepsin B. Therefore, to identify which of the measured QuPath variables would most closely represent the visually assessed intensity, each variable



**Figure 1:** Workflow for digital pathology (left) and mass spectrometry (right) analysis in the presented study

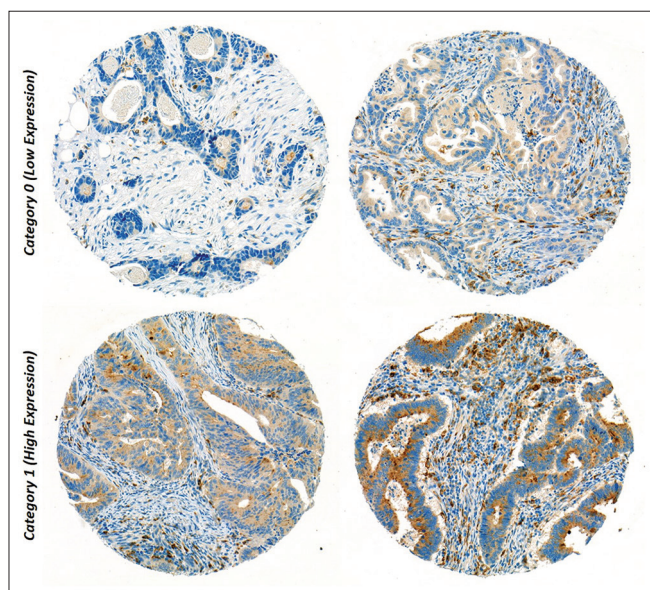


**Figure 2:** (a) *P* values and associated fold change between Stage 2 and Stage 3 patients for all proteins identified in our mass spectrometry results for each of the three quantification scores (◊: “Top3” scores from ProteomeDiscoverer, ◯: label-free quantification scores from MaxQuant, +: “Top 3” scores from MaxQuant). Note that a negative fold change indicates higher expression in Stage 2. Eight proteins with  $P < 0.05$  in all three quantification scores (indicated in bold) were considered as biomarker candidates. (b) Additional information for each of the eight biomarker candidates

**Table 1: Association of cathepsin B staining by visual assessment and digital image analysis and clinicopathological features**

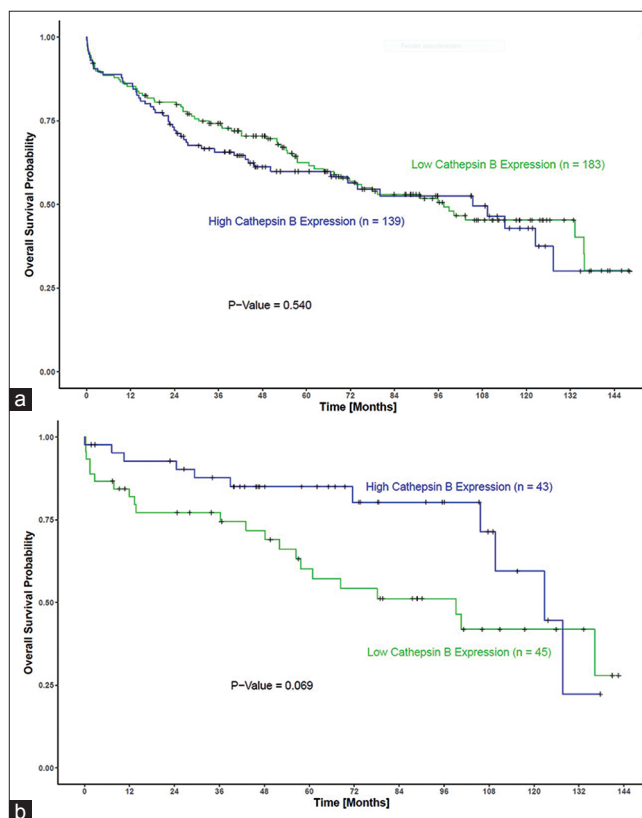
Clinical-pathological variable	Visual intensity (P)	Brightness_Max in Stage 2 (P)	Brightness_Max in full cohort (P)
Age	0.960	0.089	0.747
Gender	0.296	0.517	0.341
Tumor histology	0.338	0.797	0.488
Tumor location	0.392	0.030	0.583
pT	0.917	0.143	<0.001
pN	0.831	-	0.004
cM	1.000	-	0.001
TNM stage	0.831	-	<0.001
Grade	0.511	0.343	0.001
L	1.000	0.089	<0.001
V	0.321	0.650	0.008
Pn	0.430	0.296	<0.001
Number of buds (ITBCC)	1.000	0.042	0.007
Budding category	0.833	0.086	0.042
Tumor border configuration	0.615	0.117	0.001
MMR status	1.000	0.090	0.024
Overall survival	0.612	0.044	0.042

TNM: Tumor, node, metastasis



**Figure 3:** Reference tissue cores stained for cathepsin B used for assessment of visual intensity. The top row showing cores which were categorized as “low cathepsin B expression,” the bottom row showing cores categorized as “high cathepsin B expression”

was independently analyzed in a Cox proportional hazards model containing Stage 2 patients in the validation cohort. The limited size of the cohort precluded the application of a multivariate model. Univariate Cox analysis resulted in three QuPath variables with significant ( $P < 0.05$ ) hazard ratios: Red\_Max, Residual\_Min, and Brightness\_Max. We considered Brightness\_Max, which is representative of the maximum median pixel intensity of all cells detected on particular patient core in the tumor epithelium region, as the best surrogate for visual intensity. Since visual



**Figure 4:** Overall survival in the full validation cohort (a) and in Stage 2 patients of the validation cohort (b), stratified by visual intensity category. Indicated P values are the result of a log-rank test

intensity is manifested by higher staining intensity (i.e., darker staining), we had an *a priori* expectation that Brightness\_Max would be inversely correlated with visual intensity. Indeed, our data showed significantly higher

values of Brightness\_Max values in the “low cathepsin B expression” category, which confirmed our *a priori* hypothesis [Supplementary Figure 2].

### Association of cathepsin B Brightness\_Max scores with clinical pathological features

Our DIA approach incorporated the alignment of cytokeratin-stained and cathepsin B-stained cores, followed by cell detection and digital assessment of staining intensity. Based on the specificity of cytokeratin to epithelial regions, this approach ensured that cathepsin B expression would be digitally assessed only in the tumor epithelial regions, thereby removing any potential bias from stromal signals [Figure 5]. Based on our DIA data, we successfully established QuPath’s Brightness\_Max measure as a digital surrogate of visual intensity. It also enabled an investigation of the association with important clinical pathological features of the validation cohort. This analysis was performed for both Stage 2 patients and for all patients in the validation cohort [Table 1]. Interestingly, we found lower values of Brightness\_Max (i.e., higher expression of cathepsin B) significantly associated with less aggressive pathological features of CRC in the full cohort, such as the absence of lymphatic ( $P < 0.001$ ) and perineural ( $P < 0.001$ ) invasion, and a lower number of tumor buds ( $P = 0.007$ ). We could also demonstrate a significant association of Brightness\_Max with overall survival in the full cohort ( $P < 0.05$ ), although this was only found in univariate analysis. Unfortunately, Brightness\_Max does not represent an independent prognostic factor in multivariate analysis when adjusting for other confounders.

## CONCLUSIONS

The lack of prognostic and predictive biomarkers remains an important issue in proper patient classification and management for CRC. In this study, we investigated whether the introduction of a novel approach coupling DIA and state-of-the-art tissue enrichment methodology to mass spectrometry for biomarker discovery would allow for the identification of clinically relevant novel biomarkers, with a particular focus on markers in Stage 2 patients. The presented study identified the human cysteine cathepsin B as a promising biomarker and proposed validation guidelines using ngTMA<sup>®</sup> coupled to DIA on a previously established validation cohort of 413 patients. Although the cathepsin B protein is not as interesting in terms of differential protein expression compared to other candidates such as Hexokinase 3, it was chosen for validation based on a number of different reasons. First, the role of cathepsin B in cancer malignancy and progression is supported by a large body of literature: cathepsin B expression has been directly associated with a more aggressive phenotype and positive expression in tumor buds in CRC,<sup>[16]</sup> and the cathepsin family of proteins is generally attributed an important promoting role in a multitude of cancers.<sup>[17-20]</sup> Second, the investigation of tumor budding as the main focus of our research group aligned well with the previously reported function of cathepsin B. Finally, the availability of a commercial antibody for cathepsin B and its excellent quality of staining was an inevitable prerequisite for validation in our study. Taken together, this renders cathepsin B the most obvious choice for validation in our list of biomarker candidates.

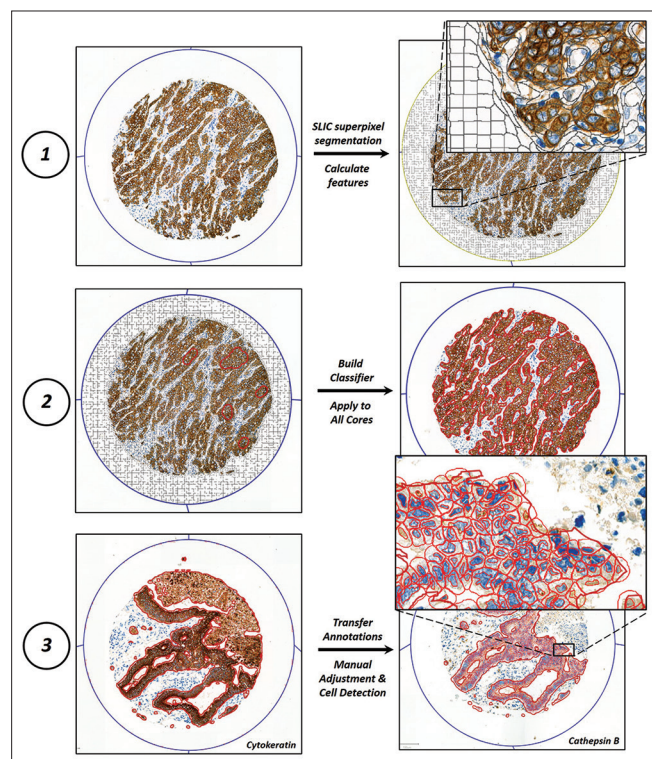


Figure 5: Digital image analysis workflow

While our results indicated that the visually inspected cathepsin B category was not a suitable prognostic factor on the validation cohort, the analysis specific to Stage 2 patients revealed its prognostic potential. Stratification of patients by visual intensity category revealed two prognostic subgroups, with a trend toward statistical significance. Interestingly, the better prognostic subgroup is characterized by increased expression of cathepsin B, indicating a “protective” effect of this protein.

The role of cathepsin B, one of the eleven human cysteine cathepsins (B, C, F, H, L, K, O, S, V, W, X/Z), has been well described.<sup>[21]</sup> The protein is constitutively expressed as a precursor on the rough endoplasmatic reticulum and transported to lysosomes, where it is converted into its active form. A dual role has been reported for this protein in cancer: On the one hand, it has been shown to promote tumor progression via increased secretion into the extracellular matrix and degradation of the basement membrane.<sup>[22-24]</sup> On the other hand, cathepsin B has also been reported as an important proapoptotic component of apoptosis by cleaving anti-apoptotic members of the Bcl-2 family of proteins.<sup>[25,26]</sup> Taken together, evidence in the literature suggests opposing roles for cathepsin B in malignancy. Our data indicate that the “protective” effect of cathepsin B is more pronounced in our validation cohort since the favorable prognostic subgroup of

Stage 2 patients is characterized by increased expression of cathepsin B.

With respect to DIA, our study revealed important points for consideration. We have proposed a method to successfully identify a digital “surrogate” measure based on the framework of visually assessed intensity levels. In this context, our approach outlines the benefits and challenges of incorporating DIA into biomarker discovery studies. On the one hand, DIA provides a continuous measure of visual intensity, which bears obvious benefits in terms of statistical evaluation compared to the classical categorical assessment of visual intensity. The benefits of a digital, continuous measure of intensity for biomarker validation are shown in Table 1, which outlines significant associations with relevant clinical pathological features in both Stage 2 patients and in all patients of the validation cohort. A comparison of the same features with categorical visual intensity fails to reach statistical significance.

On the other hand, our approach also outlines the challenges associated with DIA-based approaches. The amount of data produced by QuPath (and DIA software in general) is orders of magnitude larger than any type of visual assessment. This offers many more possibilities for data analysis, such as an examination of quantitative cell-specific expression data or cell–cell interactions. However, the resulting data analysis, processing, and storage are more challenging and require advanced expertise in bioinformatics, which may not be readily available in research institutions.

Our study may be limited in several aspects. One very important point is the mode of synthesis of cathepsin B. Proteins synthesized in the form of a precursor which is only activated at a later stage form a major challenge for IHC, since antibodies specific to either the precursor or its active form are difficult to obtain. In most cases, the IHC staining does not differentiate between the two forms, including the presented work. Therefore, our results cannot be used to distinguish whether the more favorable identified prognostic subgroup in Stage 2 is truly characterized by increased expression of cathepsin B, its precursor form, or a combination of the two. In addition, our experience has shown that while the incorporation of DIA approaches may increase the objectivity of results, image analysis software remains prone to imperfections. These imperfections may be manifested in the form of tissue artifacts, mast cells which intensely pick up a wide range of staining, or two closely adjacent cells which are falsely identified as one much larger cell. Such issues are avoided in visual assessment of staining intensities since they are “obvious” to the observer. This addresses an important aspect that will require further refinement in the comparison of visual versus digital pathology.

In conclusion, cathepsin B is a marker for more indolent tumor behavior and favorable prognosis in colorectal cancer. Our work has shown that inclusion of DIA increases objectivity compared to visual assessment. However, we have also outlined some of the challenges associated, including the amount and analysis of validation data. We estimate that with

increasing popularity of various DIA software packages, this approach will be the focus of many future studies which will help to clarify current challenges in the comparison between visual and digital pathology.

### Acknowledgments

This project was supported by the Werner and Hedy Berger-Janser Stiftung. The funders had involvement neither in the study design, collection, analysis, and interpretation of the data, nor in writing of the report and decision to submit the paper for publication. The authors would also like to thank Prof. Dr. Karl-Friedrich Becker and his team for support and assistance in preparations of FFPE samples for mass spectrometry analysis.

### Financial support and sponsorship

This project was supported by the Werner and Hedy Berger-Janser Stiftung.

### Conflicts of interest

There are no conflicts of interest.

### REFERENCES

1. Lawrence MS, Stojanov P, Mermel CH, Robinson JT, Garraway LA, Golub TR, *et al.* Discovery and saturation analysis of cancer genes across 21 tumour types. *Nature* 2014;505:495-501.
2. Andor N, Graham TA, Jansen M, Xia LC, Aktipis CA, Petritsch C, *et al.* Pan-cancer analysis of the extent and consequences of intratumor heterogeneity. *Nat Med* 2016;22:105-13.
3. Cyll K, Ersv er E, Vlatkovic L, Pradhan M, Kildal W, Avranden Kj er M, *et al.* Tumour heterogeneity poses a significant challenge to cancer biomarker research. *Br J Cancer* 2017;117:367-75.
4. Calon A, Lonardo E, Berenguer-Llargo A, Espinet E, Hernando-Mombona X, Iglesias M, *et al.* Stromal gene expression defines poor-prognosis subtypes in colorectal cancer. *Nat Genet* 2015;47:320-9.
5. Becht E, de Reyni es A, Giraldo NA, Pilati C, Buttard B, Lacroix L, *et al.* Immune and stromal classification of colorectal cancer is associated with molecular subtypes and relevant for precision immunotherapy. *Clin Cancer Res* 2016;22:4057-66.
6. Dunne PD, McArt DG, Bradley CA, O'Reilly PG, Barrett HL, Cummins R, *et al.* Challenging the cancer molecular stratification dogma: Intratumoral heterogeneity undermines consensus molecular subtypes and potential diagnostic value in colorectal cancer. *Clin Cancer Res* 2016;22:4095-104.
7. St lhammar G, Robertson S, Wedlund L, Lippert M, Rantalainen M, Bergh J, *et al.* Digital image analysis of Ki67 in hot spots is superior to both manual Ki67 and mitotic counts in breast cancer. *Histopathology* 2018;72:974-89.
8. Nolte S, Zlobec I, Lugli A, Hohenberger W, Croner R, Merkel S, *et al.* Construction and analysis of tissue microarrays in the era of digital pathology: A pilot study targeting CDX1 and CDX2 in a colon cancer cohort of 612 patients. *J Pathol Clin Res* 2017;3:58-70.
9. Koelzer VH, Sokol L, Zahnd S, Christe L, Dawson H, Berger MD, *et al.* Digital analysis and epigenetic regulation of the signature of rejection in colorectal cancer. *Oncoimmunology* 2017;6:e1288330.
10. Sokol L, Koelzer VH, Rau TT, Karamitopoulou E, Zlobec I, Lugli A, *et al.* Loss of tapasin correlates with diminished CD8(+) T-cell immunity and prognosis in colorectal cancer. *J Transl Med* 2015;13:279.
11. Lugli A, Kirsch R, Ajioka Y, Bosman F, Cathomas G, Dawson H, *et al.* Recommendations for reporting tumor budding in colorectal cancer based on the international tumor budding consensus conference (ITBCC) 2016. *Mod Pathol* 2017;30:1299-311.
12. Ahrm e E, Molzahn L, Glatter T, Schmidt A. Critical assessment of proteomics-wide label-free absolute abundance estimation strategies. *Proteomics* 2013;13:2567-78.

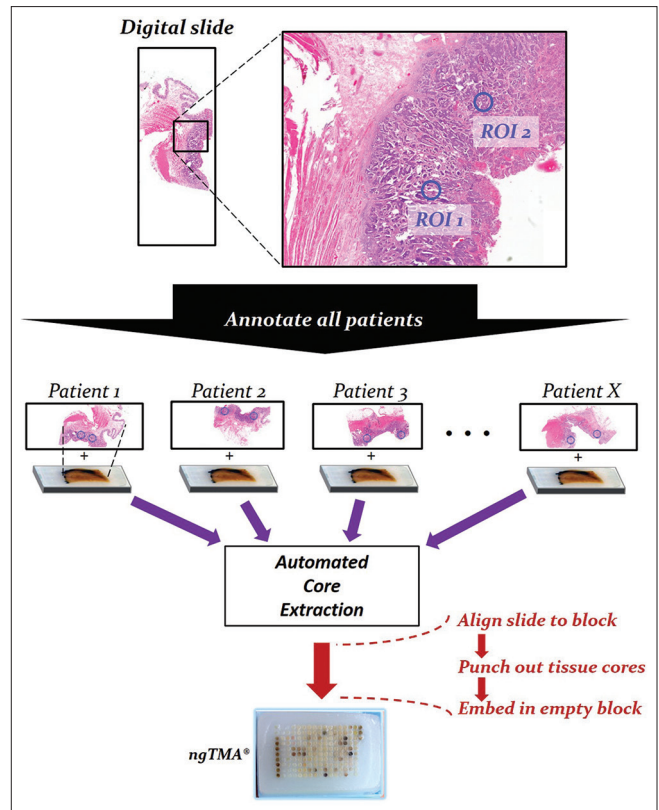
13. Grossmann J, Roschitzki B, Panse C, Fortes C, Barkow-Oesterreicher S, Rutishauser D, *et al.* Implementation and evaluation of relative and absolute quantification in shotgun proteomics with label-free methods. *J Proteomics* 2010;73:1740-6.
14. R Core Team. R: A Language and Environment for Statistical Computing; 2014. Available from: <https://www.r-project.org/>.
15. Bankhead P, Loughrey MB, Fernández JA, Dombrowski Y, McArt DG, Dunne PD, *et al.* QuPath: Open source software for digital pathology image analysis. *Sci Rep* 2017;7:16878.
16. Guzińska-Ustymowicz K. MMP-9 and cathepsin B expression in tumor budding as an indicator of a more aggressive phenotype of colorectal cancer (CRC). *Anticancer Res* 2006;26:1589-94.
17. Koblinski JE, Ahram M, Sloane BF. Unraveling the role of proteases in cancer. *Clin Chim Acta* 2000;291:113-35.
18. Jedeszko C, Sloane BF. Cysteine cathepsins in human cancer. *Biol Chem* 2004;385:1017-27.
19. Mohamed MM, Sloane BF. Cysteine cathepsins: Multifunctional enzymes in cancer. *Nat Rev Cancer* 2006;6:764-75.
20. Tan GJ, Peng ZK, Lu JP, Tang FQ. Cathepsins mediate tumor metastasis. *World J Biol Chem* 2013;4:91-101.
21. Turk V, Stoka V, Vasiljeva O, Renko M, Sun T, Turk B, *et al.* Cysteine cathepsins: From structure, function and regulation to new frontiers. *Biochim Biophys Acta* 2012;1824:68-88.
22. Steffan JJ, Snider JL, Skalli O, Welbourne T, Cardelli JA. Na<sup>+</sup>/H<sup>+</sup> exchangers and rhoA regulate acidic extracellular pH-induced lysosome trafficking in prostate cancer cells. *Traffic* 2009;10:737-53.
23. Rozhin J, Sameni M, Ziegler G, Sloane BF. Pericellular pH affects distribution and secretion of cathepsin B in malignant cells. *Cancer Res* 1994;54:6517-25.
24. Glunde K, Guggino SE, Solaiyappan M, Pathak AP, Ichikawa Y, Bhujwala ZM, *et al.* Extracellular acidification alters lysosomal trafficking in human breast cancer cells. *Neoplasia* 2003;5:533-45.
25. Cirman T, Oresić K, Mazovec GD, Turk V, Reed JC, Myers RM, *et al.* Selective disruption of lysosomes in HeLa cells triggers apoptosis mediated by cleavage of bid by multiple papain-like lysosomal cathepsins. *J Biol Chem* 2004;279:3578-87.
26. Droga-Mazovec G, Bojic L, Petelin A, Ivanova S, Romih R, Repnik U, *et al.* Cysteine cathepsins trigger caspase-dependent cell death through cleavage of bid and antiapoptotic Bcl-2 homologues. *J Biol Chem* 2008;283:19140-50.



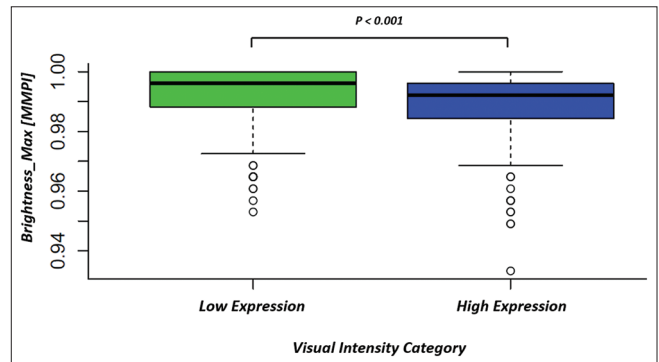
### Supplementary Table 1: Clinical-pathological features of the patient cohort

Feature	Frequency (%)
Age (median) (minimum-maximum)	70.6 (19.1-92.1)
Sex	
Female	161 (39.2)
Male	250 (60.8)
Grade	
1	10 (2.9)
2	262 (75.3)
3	76 (21.8)
L	
0	126 (34.7)
1	237 (65.3)
V	
0	168 (45.8)
1	199 (54.2)
pT	
0	22 (5.7)
1	9 (2.3)
2	66 (17.0)
3	210 (54.0)
4	82 (21.0)
pN	
0	199 (48.5)
1	124 (30.2)
2	87 (21.2)
pM	
0	296 (72.0)
1	115 (28.0)
TNM stage	
0	19 (4.6)
1	47 (11.4)
2	106 (25.8)
3	124 (30.2)
4	115 (28.0)
Pn	
0	281 (78.7)
1	76 (21.3)
Preoperative TX	
No	360 (87.6)
Yes	51 (12.4)
Tumor border configuration (percentage expanding, median) (minimum-maximum)	50 (0-100)
Budding (ITBCC, median) (minimum-maximum)	3 (0-195)
Microsatellite status	
MSI	30 (14.7)
MSS	174 (85.3)
Tumor histology	
Adenocarcinoma	335 (87.5)
Mucinous	32 (8.4)
Other	16 (4.2)
Overall survival (median) (minimum-maximum)	42.7 (0-182.9)
Disease-free survival (median) (minimum-maximum)	38.8 (0-161.5)

TNM: Tumor, node, metastasis, MSI: Microsatellite instability, MSS: Microsatellite stable



Supplementary Figure 1: Next-generation Tissue Microarray Construction



Supplementary Figure 2: Significantly higher values of Brightness\_Max scores in the "low" compared to "high" Cathepsin B expressing category