



AperTO - Archivio Istituzionale Open Access dell'Università di Torino

Tumour-educated circulating monocytes are powerful candidate biomarkers for diagnosis and disease follow-up of colorectal cancer

This is the author's manuscript								
Original Citation:								
Availability:								
This version is available http://hdl.handle.net/2318/1724034	since 2020-01-20T14:43:24Z							
Published version:								
DOI:10.1136/gutjnl-2014-308988								
Terms of use:								
Open Access								
Anyone can freely access the full text of works made available as "Open Access". Works made available under a Creative Commons license can be used according to the terms and conditions of said license. Use of all other works requires consent of the right holder (author or publisher) if not exempted from copyright								

(Article begins on next page)

protection by the applicable law.

Alexander Hamm^{1,2#}, Hans Prenen^{3#}, Wouter Van Delm^{4#}, Mario Di Matteo^{1,2}, Mathias Wenes^{1,2}, Estelle Delamarre^{1,2}, Thomas Schmidt⁵, Jürgen Weitz^{5,6}, Roberta Sarmiento⁷, Angelo Dezi⁷, Giampietro Gasparini⁷, Françoise Rothé⁸, Robin Schmitz⁵, André D'Hoore⁹, Hannes Iserentant¹⁰, Alain Hendlisz⁸ & Massimiliano Mazzone^{1,2}

¹Lab of Molecular Oncology and Angiogenesis, Vesalius Research Center, VIB, Leuven, Belgium ²Lab of Molecular Oncology and Angiogenesis, Vesalius Research Center, Department of Oncology, KU Leuven, Leuven, Belgium

³Digestive Oncology, University Hospitals Leuven and Department of Oncology, KU Leuven, Leuven, Belgium

⁴Nucleomics Core, VIB, Leuven, Belgium

⁵Department of General, Visceral, and Transplantation Surgery, University of Heidelberg, Heidelberg, Germany

⁶Department of Visceral, Thoracic, and Vascular Surgery, University Hospital Carl Gustav Carus,

Technical University Dresden, Dresden, Germany

⁷Department of Oncology, San Filippo Neri, Rome, Italy

⁸Medical Oncology Clinic, Institut Jules Bordet, Brussels, Belgium

⁹Department of Abdominal Surgery, University Hospitals Leuven, KU Leuven, Leuven, Belgium

¹⁰VIB, Zwijnaarde, Belgium

[#]contributed equally to this study

Correspondence:

Massimiliano Mazzone, <u>massimiliano.mazzone@vib-kuleuven.be</u>, Tel: +32-16-373213, Fax +32-16-372585 VIB Vesalius Research Center, KU Leuven, Herestraat 49, Bus 912, 3000 Leuven, Belgium

Hans Prenen, hans.prenen@uzleuven.be, Tel: +32-16-340238

Word Count: 4127

Key Words: Monocytes, colorectal cancer, screening, inflammation

https://mc.manuscriptcentral.com/gut

LIST OF ABBREVIATIONS

- AUC area under the curve
- BER balanced error rate

1 2 3

4 5

6

7

8

9

10

11

12

13

14

15

16

17 18

19

20

21

22

23

24

25

26 27

28

29 30

31 32

33

34

35

36

37 38

39

- CEA carcino-embryonic antigen
- CRC colorectal cancer
- ENS ensemble method
- FIT fecal immunochemical test
- FOBT fecal occult blood test
- MACS magnet-associated cell sorting
- MCCV Monte Carlo cross validation
- NSAID non-steroid anti-inflammatory drugs
- PBM peripheral blood monocytes
- peripheral blood mononuclear cells PBMC
- qPCR quantitative RT-PCR
- RF random forest
- ROC receiver operating characteristics
- RT-PCR reverse-transcription polymerase chain reaction
- Se sensitivity
- SGMV single gene majority vote
 - Sp specificity
- SVM support vector machine
- UICC Union internationale contre le cancer

Labels of patient groups:

HV	healthy volunteer
Р	non-metastatic CRC patient

- ents P. PM non-metastatic and metastatic CRC patients
- PC pancreatic cancer patient
- PG gastric cancer patient
 - PGT gastritis patient
 - ΡM metastatic CRC patient
 - PR patient in remission from CRC

ABSTRACT

Objective: Cancer immunology is a growing field of research whose aim is to develop innovative therapies and diagnostic tests. Starting from the hypothesis that immune cells promptly respond to harmful stimuli, we utilized peripheral blood monocytes (PBM) in order to characterize a distinct gene expression profile and to evaluate its potential as a candidate diagnostic biomarker in colorectal cancer (CRC) patients, a still unmet clinical need.

Design: We performed a case-control study including 360 PBM samples from four European oncological centres and defined a gene expression profile specific to CRC. The robustness of the genetic profile and disease specificity, were assessed in an independent setting.

Results: This screen returned 43 putative diagnostic markers, which we refined and validated in the confirmative multicentric analysis to 23 genes with outstanding diagnostic accuracy (AUC=0.99 [0.99:1.00]. Se=100.0% [100.0%:100.0%]. Sp=92.9% [78.6%;100.0%] in multiple-gene ROC analysis). The diagnostic accuracy was robustly maintained in prospectively collected independent samples (AUC=0.95 [0.85;1.00], Se=92.6% [81.5%;100.0%], Sp=92.3% [76.9%;100.0%]. This monocyte signature was expressed at early disease onset, remained robust over the course of disease progression, and was specific for the monocytic fraction of mononuclear cells. The gene modulation was induced specifically by soluble factors derived from transformed colon epithelium in comparison to normal colon or other cancer histotypes. Moreover, expression changes were plastic and reversible, as they were abrogated upon withdrawal of these tumour-released factors. Consistently, the modified set of genes reverted to normal expression upon curative treatment and was specific for CRC.

<text><text><text><text>

SUMMARY BOX

What is already known on this subject?

- Early diagnosis of colorectal cancer is crucial for curative surgical treatment,
 highlighting the need for efficient screening tools.
- Colorectal cancer screening is a rapidly evolving field, as several strategies for supplementing the invasive colonoscopic screening are explored.
- Circulating cells of the immune system in the blood stream are easily accessible, yet understudied with regard to their precise role in tumour immunology.
- Tumour-associated macrophages deriving from circulating monocytes can display diverse phenotypes and affect tumour growth and metastasis by different means, depending on the cellular context.

What are the new findings?

- Monocytes are plastic cells that are modified by early occurrence of colorectal cancer, resulting in a highly specific genetic fingerprint, which is independent of tumour stage.
- The changes in monocyte expression profiles are reversible, highly specific to the tissue type and cancer histotype, and induced in response to soluble factors released by the cancer cells in the primary or metastatic site.
- The specific genetic fingerprint in circulating monocytes can be harnessed for diagnosis and disease follow-up of colorectal cancer.

How might it impact on clinical practice in the foreseeable future?

 μη

 μοτινα και.

 μοτινα ματά τη τη μοτιατίτη τη τη μοτιατίτη τη τη προστή τη τη μοτιατίτη τη τη προστή προστή τη προστή τη προστή τη προστή προστή τη προστή προστή τη > If the initiated prospective validation study supports our sound results, our

INTRODUCTION

Colorectal cancer (CRC) is the second leading cause of cancer-related deaths in the US¹. Its incidence and the difficulty in early-diagnosis make CRC a primary focus in the oncology community². Early CRC is symptomless, and, consequently, is frequently diagnosed when already advanced. Metastatic disease (found in 30 to 40% of CRC patients) is associated with a poor 5-year survival rate of less than 10%. In contrast, up to 80% of patients can be cured by early tumour resection, rendering timely diagnosis a crucial factor for proper disease management². Nevertheless, endoscopic screening as well as stool tests (fecal immunochemical test, FIT, or fecal occult blood test, FOBT⁵) are not widely accepted by the target population, while the socioeconomical burden of these procedures is high². Thus, there is urgent need to identify specific, non-invasive biomarkers for early CRC diagnosis and treatment monitoring to avoid disease progression to advanced stages that are difficult to cure⁶. Peripheral blood is one of the least invasive sample sources that can be intensively screened for CRC biomarkers. Within the blood stream, peripheral blood monocytes (PBM) represent a reservoir of inflammatory cells that contribute to disease progression by different means^{7 8}. These cells are recognized to be plastic and versatile cells, which can change their phenotype in response to microenvironmental stimuli, yielding either tumouricidal or pro-tumourigenic features depending on the stromal context or tumour type¹⁰¹¹. Interestingly, recent studies have suggested distinct expression profiles in circulating monocytes in several pathological conditions such as diabetes¹², atherosclerosis¹³, and dysmenorrhea¹⁴, though none have convincingly demonstrated a specific regulation of monocyte heterogeneity by malignantly transformed cells apart from descriptive studies in vitro on monocytic cell lines¹⁵.

Several novel accessible diagnostic tools share the major opportunity to make frequent screening more appealing to a greater number of patients, as a less invasive method is likely to increase compliance and allow for decreased screening intervals (recently comprehensively reviewed⁶). While conventional blood-based tumour markers (particularly carcino-embryonic antigen, CEA¹⁶) have been established as supplemental markers in treatment monitoring, they have failed to yield high diagnostic accuracy as primary screening tools. In addition to the established FIT or FOBT⁵, other potential diagnostic markers include serumassociated biomarkers (e.g. circulating tumour DNA¹⁷, micro-RNA¹⁸, methylation markers like SEPT9¹⁹), genetic marker sets in white blood cells²⁰⁻²³, and, most recently, fecal tumour DNA²⁴. However, all of these approaches display limited sensitivity and specificity⁶. In this study, we therefore assess the sensitivity and specificity of a novel gene signature in circulating monocytes for the diagnosis of CRC in comparison to healthy individuals or to other cancer types, and assess its robustness in prospectively obtained samples.

PATIENTS AND METHODS

Patients

We collected a total of 360 samples between January 13, 2010 and January 26, 2015, comprised of the following cohorts: cohort I (genome-wide screening in 27 patients with non-metastatic stage I, stage II, or stage III CRC (P), 28 patients with metastatic stage IV CRC (PM), and 38 healthy volunteers (HV) (without history or evidence of acute or chronic disease)), cohort II (multicentric validation in 73 patients and 61 healthy volunteers from four different oncological centres), cohort III (robustness assessment in 27 patients and 13 asymptomatic healthy individuals with colonoscopy-confirmed absence of disease), cohort IV (15 patients with gastric cancer (PG), 16 patients with pancreatic cancer (PC), 10 patients with gastritis (PGT), all treatment-naïve, and 13 HV), cohort V (15 curatively treated patients), and cohort VI (comparative expression analysis in PBM and PBMC in 17 patients and 7 healthy volunteers). See Figure 1 for allocation of collected samples to analyses. All participants gave written informed consent, and the study was approved by the respective institutional review boards. Details on inclusion and exclusion criteria, participating centers and ethical approval can be found in Supplementary Methods.

Identification of a gene signature

Genome-wide expression analysis was performed on the Illumina platform (Illumina) on RNA obtained from peripheral blood monocytes (PBM), isolated by a two-step procedure with density gradient centrifugation and positive selection for CD14 using the MACS system (Miltenyi). Details are reported in Supplementary Methods. Differential expression was assessed with the limma package of R²⁶. Putative candidate genes were confirmed on a random subset of cohort I and validated by

Gut

quantitative RT-PCR (qPCR) on the 7500Fast System (Applied Biosystems) using intron-spanning PrimeTime qPCR Assays (Integrated DNA Technologies) listed in Supplementary Table 1 as described in Supplementary Methods. For statistical analysis, we followed a three-step top-down approach to construct a gene signature for CRC, with details explained in Supplementary Methods.

Multicentric validation study

For validation of a diagnostic test, we used cohort II to train and validate a multi-gene classifier. Splits in training and test sets for validation were performed by stratified random sampling for centre of origin and class label as detailed out in Supplementary Methods. Samples with missing values for more than 25% of the genes were excluded from the analysis. We ruled out an effect of the class labeling on the percentage of missing values with Fisher's exact test (Supplementary Table 2).

The training dataset was used to build three types of classifiers: a support vector machine (SVM)²⁹ with linear kernel, a single-gene majority vote (SGMV) classifier, and a random forest classifier (RF³⁰). Subsequently, we applied an ensemble method³¹ that votes according to the majority of the three independent classifiers. Performance was validated both with ranking (AUC) and classification (balanced error rate, BER, Se, Sp) scores with 95% confidence intervals ([lower boundary; upper boundary]). We explicitly opted for relatively simple computational models in order to limit chances of over-fitting the training data and to maximize interpretability of the models' internal decision-making process. Model flexibility was further controlled through a Monte-Carlo cross-validation scheme (MCCV)³², before final estimation of the model parameters. Validation of the predictive models was done on

the test set of cohort II, which were not included during development of the models. Details on all classification methods are specified in Supplementary Methods. In order to avoid biased conclusions, the analysis of the 23 genes was complemented with a study by an independent team (DNAlytics, Belgium) that adopted a slightly modified analysis protocol (see Supplementary Methods). All complementary analyses were performed in R with scripts designed by DNAlytics,

In vitro model system

To study the effects of tumour-released soluble factors on gene expression in monocytes, we established an *in vitro* model system, where monocytes from healthy donors were challenged with tumour-released soluble factors and changes in gene expression profile were analyzed by qPCR. See Supplementary Methods for details.

fully independently from other analyses described in this paper.

RESULTS

Establishment of putative biomarkers by genome-wide expression analysis

<text><text><text> To obtain a set of putative biomarkers that might facilitate early diagnosis of CRC, we have performed a genome-wide expression analysis on PBMs from 55 untreated patients newly diagnosed with CRC and 38 healthy volunteers (cohort I). All relevant clinicopathological information on patient cohorts can be found in Table 1.

Cohort I P,PM HV			II P,PM HV									V	VI		
										P,PM	HV	HV P	P,PM HV		
			LEU ^a	HD⁵	SFN℃	IJB ^d	LEU ^a	HD⁵	SFN℃	IJB ^d					
Number of	55	38	39	19	10	5	20	12	14	15	27	13	15	17	7
samples															
Age															
median	67	55	66	69	72	59	49	55	47	49	66	62	69	78	42
range	44-87	42-79	47-78	42-76	50-85	52-82	42-69	46-75	40-63	42-62	44-90	43-74	45-81	62-89	42-57
Gender															
male	22	15	24	11	5	1	15	7	11	2	14	8	8	11	5
female	33	23	15	8	5	4	5	5	3	13	13	5	7	6	2
metastatic	28	1	16	3	2	2	/	1	/	1	16	/	0	6	/
non-metastatic	27	1	23	16	8	3	1	1	/	1	11	1	15	11	/
UICC stage															
- 1	3	/	7	2	1	1	1	/	/	1	2	/	4	2	/
2	12	1	8	8	2	0	1	1	/	1	3	/	7	7	/
3	12	1	8	6	5	2		1	/	1	6	/	4	2	/
4	28	1	16	3	2	2		1	/	1	16	/	0	6	/
Tumour localization															
Caecum	5	1	3	3	0	0	/		/	1	2	/	1	2	/
Ascendens	11	1	4	3	1	0	1	I	1	1	6	/	4	4	/
Transversum	0	1	4	3	0	0	/			1	2	/	1	0	/
Descendens	4	1	2	1	3	3	/	1		. 🔶 /	0	/	0	1	/
Sigmoid	28	1	15	3	2	0	/	1	1	1	10	/	5	6	/
Rectum	6	/	8	5	3	2	/	1			7	/	4	2	/
Double	1	1	3	1	1	0	1	1	1		0	/	0	2	1

^aLeuven, ^bHeidelberg, ^cRome, ^dBrussels. See Supplementary Methods for the detailed description of contributing centres

The purity of the monocyte fraction was >90%, as assessed by FACS analysis in the pilot phase (Supplementary Figure 1a) and verified by hemocytometric analysis for each individual sample (Supplementary Figure 1b). Both absolute and relative monocyte counts were not different between patients and healthy volunteers (Supplementary Figure 1c). We therefore investigated differentially expressed genes by genome-wide expression analysis using the Illumina HumanHT-12 v4 Expression BeadChip Kit. The data discussed in this publication have been deposited in NCBI's Gene Expression Omnibus³³ and are accessible through GEO Series accession number GSE47756

(http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?token=hvmpvoswuqaeybc&acc=GSE 47756). In first instance, we compared the average expression values of all CRC patients (P,PM), comprised of non-metastatic (P) and metastatic (PM) patients, to that of healthy volunteers (HV). The resulting gene signature of (P,PM) versus HV consisted of 36 upregulated and 4 downregulated probes (Figure 2a, b, Table 2). In second instance, we were interested if the gene signature in patients with synchronous metastases *i.e.*, at the time of diagnosis (PM, n=28) was different from that in non-metastatic patients (P, n=27). Interestingly, the number of up- and downregulated genes was comparable in both P and PM (in comparison to HV) (Table 2 and Supplementary Figure 2a, b), while there were no genes found to be differentially expressed between the two patient groups (Supplementary Figure 2a, b), indicating that the gene signature induced at early onset stays robust over disease progression. Indeed, when post-hoc assessing those samples from patients with early stages (Tis and T1), they clustered with the rest of the patient samples (data not shown). A power analysis revealed that, for the given number of genes, samples and observed variation, chances were very low $(<10^{-10})$ that truly differentially expressed genes with

https://mc.manuscriptcentral.com/gut

fold changes larger than 1.5 had been missed. Therefore, adding more samples would probably have changed little to the panel of candidate genes that our screen returned.

Confirmation of the gene signature in independently processed samples

To validate the genetic signature, we performed quantitative RT-PCR (qPCR) analysis on a random subset of PBM from 8 samples of each of the three groups (P, PM, and HV), normalizing to reference gene *B2M*, which was selected after an extensive screening procedure (Supplementary Note 1). To avoid bias in the confirmation procedure, we freshly extracted RNA from independently stored samples for confirmative expression analysis. In analyzing 43 putative marker genes with probes listed in Supplementary Table 1, 23 genes showed differential expression between (P, PM) and HV (Supplementary Figure 3b, Table 2, and Supplementary Table 4). Thus, we were able to confirm a subset of the previously established gene signature, independent of the RNA extraction and the platform used for expression analysis. Information on the annotated biological function of the genes of the diagnostic signature can be found in Supplementary Table 5 and Supplementary Note 2.

Confirmation of the gene signature in a multicentric validation set

For a rigorous validation of the gene signature, we collected an independent multicentric validation set (cohort II) from a total of 4 different European oncological centres with stratified training and test sets as described in Supplementary Methods. Using the panel of 23 genes confirmed previously, we found consistently differential expression between all patients and the healthy volunteers (Figure 2c and

Supplementary Figure 4). In line with the findings from the screening phase, there were no detectable differences in expression levels between P and PM (Supplementary Figure 5), while either patient group alone compared to HV was differentially expressed (data not shown).

In ROC analysis for single genes, we found that some, but not all of the genes that displayed significantly differential expression were able to discriminate patient samples from healthy individual samples with acceptable AUCs (Supplementary Figure 6 and data not shown). We therefore hypothesized that a marker panel consisting of multiple genes might yield better results in discriminating sample identity. To address this question, we decided to test three different classification algorithms on this data set, namely a support vector machine $(SVM)^{29}$ with linear kernel, a single-gene majority vote (SGMV) classifier, a random forest classifier (RF³⁰), and a combined classification by an ensemble method³¹, using the outcome of the three classification algorithms for a final diagnostic decision. To limit overestimation of the performance by the particular training and test set, we performed a MCCV as a conservative estimate with 1,000 cross-validations. Performance of all classification algorithms in cohorts II – VI, including the conservative estimate of the MCCV in cohort II, is given in detail in Table 2.

https://mc.manuscriptcentral.com/gut

TABLE 2: PERFORMANCE SCORES OF MULTIGENE CLASSIFIER

	SGMV	SVM	RF	ENS
Cohort II (Validation)				
AUC [95% CI]	0.99 [0.99;1.00] 3.6	1.00 [1.00;1.00] 3 3	0.99 [0.97;1.00] 3.6	0.99 [0.99;1.00] 3 6
Sensitivity [95% CI]	100 [100:100]	93.3 [80.0:100]	100 [100:100]	100 [100:100]
Specificity [95% CI]	92.9 [78.6;100]	100 [100;100]	92.9 [78.6;100]	92.9 [78.6;100]
Cohort II (MCCV)	[,]	,]		
AUC [95% ĆI]	0.94 [0.86;1.00]	0.92 [0.83;0.99]	0.93 [0.83;1.00]	0.86 [0.72;0.99]
BER	13.3	20.0	13.3	13.3
Sensitivity [95% CI]	80.0 [60.0;100]	66.7 [20.0;93.3]	86.7 [60.0;100]	80.0 [60.0;100]
Specificity [95% CI]	93.3 [66.7;100]	93.3 [80.0;100]	93.3 [73.3;100]	93.3 [80.0;100]
Cohort III				
AUC [95% CI]	0.96 [0.89;0.99]	0.91 [0.80;0.99]	0.93 [0.79;1.00]	0.95 [0.85;1.00]
BER	7.7	15.0	7.6	7.6
Sensitivity [95% CI]	100 [100;100]	77.8 [59.3;92.6]	92.6 [81.5;100]	92.6 [81.5;100]
Specificity [95% CI]	84.6 [61.5;100]	92.1 [76.9;100]	92.3 [76.9;100]	92.3 [76.9;100]
Cohort IV (gastric cancer)				
Sensitivity [95% CI]	33.3 [13.3;60.0]	26.7 [6.7;46.7]	20.0 [0.0;40.0]	20.0 [0.0;40.0]
Conort IV (pancreatic cance	r)	0.0.10.0.0.01	0.0.10.0.0.01	0.0.10.0.0.01
Sensitivity [95% CI]	0.0 [0.0;0.0]	0.0 [0.0;0.0]	0.0 [0.0;0.0]	0.0 [0.0;0.0]
Conort IV (gastritis)	10.0.10.0.20.01	10 0 [0 0.20 0]	10 0 [0 0.20 0]	10 0 [0 0.20 0]
Sensitivity [95% Ci]	10.0 [0.0,30.0]	10.0 [0.0,30.0]	10.0 [0.0,30.0]	10.0 [0.0,30.0]
Sonsitivity [05% CI]	50 0 120 0 20 01	10 0 0 0.30 01	20 0 [0 0.50 0]	20 0 [0 0.50 0]
Cobort VI (PBMC)	50.0 [20.0,80.0]	10.0 [0.0,30.0]	20.0 [0.0,30.0]	20.0 [0.0,30.0]
	0 51 [0 19.0 80]	0 44 [0 13.0 74]	0 64 [0 31.0 94]	0 44 [0 19:0 66]
BER	59.3	49.3	52 1	52 1
Sensitivity [95% CI]	10 0 [0 0:30 0]	30.0 [10.0.60.0]	10.0 [0.0:30.0]	10 0 [0 0.30 0]
Specificity [95% CI]	71.4 [28.6:100]	71.4 [42.5:100]	85.7 [57.1:100]	85.7 [57.1:100]
Cohort VI (PBM)	[,]			
AUC [95% CI]	1.00 [1.00;1.00]	0.79 [0.54;1.00]	1.00 [1.00;1.00]	1.00 [1.00;1.00]
BER	0.0	30.8	0.0	0.0
Sensitivity [95% CI]	100 [100;100]	38.5 [15.4;61.5]	100 [100;100]	100 [100;100]
Specificity [95% CI]	100 [100;100]	100 [100;100]	100 [100;100]	100 [100;100]

Listed are the performance scores of all multi-gene classifiers (SGMV, SVM, RF) and their combined ensemble method (ENS) of all different cohorts – please see methods for details. ^aSensitivity for labeling a gastric cancer sample as CRC

^bSensitivity for labeling a curatively treated patient in full remission as CRC

Strikingly, we achieved a remarkably high AUC of 0.99 [0.99;1.00] with a BER of 3.6% (Figure 2d and Table 2), translating into a sensitivity of 100.0% [100.0%;100.0%] and a specificity of 92.9% [78.6%;100.0%] (Table 2). Neither of the classification algorithms was capable of separating P from PM or detect differences dependent on tumour localization (Supplementary Note 3).

In order to assess whether the diagnostic gene signature is actually suitable for diagnosis of CRC in a screening setting, we have initiated a prospective sample collection in both patients and healthy individuals who are subjected to colonoscopy. In a pilot analysis in 27 patients (newly diagnosed with CRC by screening colonoscopy) and 13 healthy individuals negative to screening colonoscopy (cohort III), we found an AUC of 0.95 [0.85;1.00] with a BER of 7.6%, yielding a sensitivity of 92.6% [81.5%;100.0%] and a specificity of 92.3% [76.9%;100.0%] (Figure 2e and Table 2). The complementary data analysis (independently performed by DNAlytics, Belgium) on the same panel of 23 genes led to the matching conclusions in terms of performance. The first experiment consisted in cross-validating a model on Cohort II (BER: 8.4% [3.4%;13.4%]; AUC: 0.93 [0.88;0.98]). A second experiment consisted in learning the same type of model on Cohort II and having it make predictions on Cohort III (BER: 13.2%; AUC: 0.92).

Soluble factors released by colorectal cancer cells induce an early, tumour type-specific and reversible genetic fingerprint in monocytes

We hypothesized that tumour-released soluble factors are the key players in inducing the genetic signature in circulating monocytes. Thus, we established an *in vitro* model system where we cultured freshly isolated human monocytes from healthy donors in different conditions. In order to assess alterations in gene expression, we

first analyzed which of the 23 genes comprising the gene signature was up- or downregulated in culture after 72 hours without any additional stimulus and excluded these from the further *in vitro* studies (Supplementary Figure 7). Out of the remaining gene signature, the majority (7/9) was specifically upregulated when culturing naïve monocytes in medium conditioned by the CRC cell line HCT116, while expression levels were not affected by mock medium (Figure 3a). Moreover, in line with the coherent induction of the specific signature independent of the stage of the disease, the induction *in vitro* was independent of hypoxic cues, as HCT116-conditioned medium in hypoxia did not induce any different expression levels than medium obtained in normoxia (Figure 3b). Likewise, the changes in expression levels of all these genes occurred already 18 hours after stimulating monocytes with the conditioned medium, consistent with the fact that already early stages are detectable by the diagnostic signature.

To rule out an off-target effect of conditioned medium *i.e.*, unspecific cues from cell metabolites, apoptotic bodies, pH, etc., we assessed the expression levels of the genes upregulated by HCT116-conditioned medium in comparison to a benign colon epithelium cell line, CCD 841 CoN (CCD), which did not induce alterations in gene expression levels different from the Mock control (Figure 3a).

Prompted by this finding, we investigated if the induction of the genetic signature was a general effect of malignant transformation or might be specific to the histotype of cancer. To address this question, we conditioned medium with a gastric cancer cell line, MKN-45 (MKN), to compare CRC to another frequent gastrointestinal solid neoplasm. Remarkably, when comparing the expression levels in naïve monocytes upon stimulation with the different conditioned media, we found that MKN-45

conditioned medium did not induce the same upregulation of the genes of interest as HCT116 conditioned medium (Figure 3c).

As immune cells are highly versatile and plastic cells mirroring the microenvironment, where they are embedded, we reasoned that the genetic signature induced by CRC in monocytes might be dependent on the continuous presence of the stimulating agents and thus be reversible upon inversion of the conditions. We therefore incubated naïve monocytes first with HCT116-conditioned medium for 18 hours and then refreshed the medium with plain culture medium, thus withdrawing the tumour-released soluble factors. Strikingly, the previously elevated expression levels of a set of marker genes were almost entirely reverted to the original (and to the mock control) expression levels 72 hours after withdrawing the tumour-cell conditioned medium (Figure 3d), whereas they remained constantly overexpressed when the conditioned medium was maintained (data not shown).

The monocyte signature is specific for CRC and might serve as a candidate biomarker of disease follow-up

Based on the *in vitro* results showing that the genetic signature is specific to CRC, we sought to confirm these findings *in vivo*. We therefore assessed the diagnostic signature in patients with i. cancer of the stomach and gastro-esophageal junction (PG, n=15) and ii. pancreatic ductal adenocarcinoma (PC, n=16), two other frequent cancers of the gastrointestinal tract¹. In addition, we analysed iii. patients with gastritis (PGT, n=10) in order to compare the gene signature in CRC to a benign inflammatory condition of the gastrointestinal tract (cohort IV). In line with the *in vitro* results we saw that the vast majority of all genes were not significantly different between either of the patient groups and healthy volunteers, indicating the specificity

of this monocyte imprinting by colorectal cancer cells (Figure 4a and data not shown). Moreover, the classifier established to diagnose CRC could not separate patients with gastric cancer (AUC 0.63 [0.48;0.77]), pancreatic cancer (AUC 0.41 [0.27;0.50]), or gastritis (AUC 0.52 [0.35;0.68]) from healthy individuals (Figure 4c, d and Table 2).

The finding that the genetic signature is reverted upon withdrawal of the stimulating agents prompted us to investigate in a pilot phase the behaviour of the entire diagnostic signature in patients upon curative treatment *i.e.*, patients with surgically removed tumours without any evidence of residual disease. To this end, we isolated monocytes from 15 patients of stages I to III treated with curative intent (with or without adjuvant treatment) and presenting at follow-up without detectable residual disease (PR) (cohort V). Here, we found that virtually all of the previously upregulated genes were reverted to expression levels comparable to those of healthy volunteers (Figure 4b). Consequently, when applying the previously established classifier, we found that it was able to distinguish accurately between patients in remission and patients with tumour, while it could not detect differences between patients in remission and healthy volunteers (Table 2).

Finally, as the plasticity of the signature offers the perspective to use the gene signature for follow-up of treated patients, we became interested if the same signature could be used to diagnose relapse (frequently as metachronous metastases rather than local recurrence²). Although our dataset was not powered to address this question with sufficient significance, we post-hoc identified four patients from cohorts I and II included at presentation with metachronous metastases. All four clustered clearly in the group of patients, separately from the healthy volunteers (Figure 4e), suggesting that the signature might be used to detect disease relapse in

line with the previous results that show coherent expression over disease progression.

The gene signature is specific to monocytes in comparison to all peripheral blood mononuclear cells (PBMCs)

To rigorously assess if the genetic fingerprint identified in monocytes was specific to this cell type or an epiphenomenon of genetic shifts in the entire population of PBMCs, we isolated both monocytes and full PBMC fractions from 17 patients and 7 healthy volunteers for a comparative analysis (cohort VI). Interestingly, we found that while in the monocyte population, the diagnostic marker set of 23 genes was upregulated in all patients (both P and PM) in accordance with our previous results (Figure 5a), there were no significant differences in the expression levels of the analyzed genes in the full PBMC compartment when comparing patients to healthy volunteers (Figure 5a). Consistently, applying the previously established classifier with the defined cut-off values, it was impossible to separate the patient group from the healthy volunteer group in PBMC (Figure 5b and Table 2), while the classifier confirmed its accuracy in PBM (Figure 5c and Table 2). Thus, the differential regulation of the gene signature in PBM used for CRC diagnosis is specific to the monocytic lineage, reinforcing our initial working hypothesis that these cells are specifically affected by tumour-secreted factors.

DISCUSSION

The dismal prognosis of CRC can be effectively attenuated by an early and accurate diagnosis, which is however hampered by low compliance rates to the available screening strategies²⁶. With this study, we present a hypothesis-driven approach to screen for specific biomarkers for diagnosis of CRC, which exploits the canonical knowledge on tumour-stroma interactions¹⁰. By using genome-wide expression analysis, we show that a distinct gene signature is detectable in circulating monocytes from CRC patients in comparison to healthy individuals. In fact, this study is the first to demonstrate specific genetic changes in the highly versatile monocyte fraction, mediated by tumour-derived soluble factors. Moreover, we convincingly demonstrate with an in vitro model system that the alterations in gene expression are induced by tumour-released soluble factors, which adds to the value of our biologybound approach in comparison to mere high-throughput screenings. Our comparison of the reported gene signature in monocytes and PBMCs strongly supports our hypothesis that monocytes, more than any other immune cell in circulation, are highly plastic and responsive to microstimuli in the blood. Since the induced expression changes are higher in vitro, it is tempting to speculate that these are dependent on the concentration of cytokines and signals, which remain to be identified.

Interestingly, our analysis indicated that the induced gene signature stays robust over progression of the disease, which is consistent with our *in vitro* findings and not entirely surprising given recent evidence for the molecular similarity between the primary tumour and its metastases³⁶.

The diagnostic gene signature established here proved to be robust independent of the technique (genomewide expression microarray vs. qPCR) and has been validated independently (Supplementary Note 4). Its utilization for diagnosis of CRC

most likely depends on the development of a one-step assay with capture of monocytes from whole blood and gene expression analysis in a multiplex qPCR assay with absolute quantification, avoiding extensive preanalytical processing steps. However, the analytical reliability of this assay needs to be thoroughly established, most likely requiring centralization of the analysis during the first phase of distribution.

The finding that the specific gene signature is reversible if the stimulating cues are withdrawn, was not only demonstrated *in vitro*, but also in a pilot analysis *in vivo* in samples of patients after curative treatment. Although not completely unexpected in view of the plastic nature of the monocyte-macrophage lineage, this analysis opens avenues for treatment monitoring and companion diagnostics and will be assessed in detail in a prospective study during patient follow-up.

If supported by further prospective validation studies, this gene signature may outperform other published non-invasive test for CRC diagnosis⁶ (including single surface markers in monocytes^{37 38}) or score similar to the most recent evaluation of fecal tumour DNA²⁴. Moreover, we are the first to demonstrate that a potential diagnostic biomarker obtained in patients at the time of primary diagnosis might also be suitable for disease follow-up and thus assessment of treatment response, owing to its high plasticity.

We acknowledge the limited conclusions that can be drawn from our case-control study. Despite the confirmation in independent samples, we cannot fully exclude possible confounders that can only be unveiled by a blinded, prospective sample collection in screening individuals. These include, but are not limited to, the bias of selecting patients that underwent colonoscopy for a clinical indication; the differences

https://mc.manuscriptcentral.com/gut

in age, nutrition status, diet, and potentially lifestyle between patients and healthy volunteers; the unblinded sample collection and processing. It is therefore of paramount importance that a prospective validation study initiated by our group includes screening individuals prospectively with blinded sample processing. In addition, strategies to minimize false negatives and false positives (with potential morbidity resulting from colonoscopy and treatment) will need to be developed. This can be achieved by calculating a risk ratio on the basis of the individual expression profile, which could replace the current binary output (cancer vs. healthy) and thus define groups at risk that need to be subjected to colonoscopy as the gold standard. An informed choice on the thresholds would, at least in first instance, emphasize a high sensitivity at the expense of specificity. The resulting morbidity has to be correlated to the morbidity of screening colonoscopy.

Our study raises important questions, which will need to be addressed in further studies. First, the biological mechanisms and pivotal regulatory pathways in directing the fate of the monocyte gene signature are still unexplored. Of note, only a few genes appear to be commonly upregulated in CRC in comparison to gastric cancer and pancreatic cancer. While this demonstrates specificity for CRC, it also means further studies will be required to identify gene signatures specific to other tumours and possibly benign pathological conditions. Second, we will need to assess if the gene signature is already imprinted in pre-neoplastic lesions (*i.e.*, polyps) and determine the transformation steps at which the specific upregulation occurs. Third, as monocyte plasticity is the starting hypothesis of this study, we will need to assess if treatment regimes (e.g., steroids, chemotherapy, irradiation, postoperative stress conditions) affect the behaviour of the gene expression profile or interfere with its diagnostic capabilities. Fourth, we are currently investigating in a prospective setting

<text><text>

Gut

CONCLUSIONS

Taken together, these data provide unprecedented evidence that tumour-educated monocytes exhibit a distinct and plastic gene signature, which may not only be ' up op en u suitable for diagnosis of CRC, but potentially allows to monitor for success of therapy or for relapse. As monocytes can be obtained in a non-invasive way, these findings offer exciting new opportunities for both improving CRC diagnosis and enriching the armamentarium of therapeutic strategies, provided that the data obtained here can be replicated in an independent broad screening setting.

ACKNOWLEDGEMENTS

We would like to express our gratitude to all patients and healthy volunteers contributing to our study. The authors are indebted to Joke Allemeersch and Christos Sotiriou for critical advice. We thank DNAlytics (Belgium) for critical independent statistical review of the raw data, Brian Wong for critical review of the manuscript, and Martin Pejcinovski, Jens Serneels, Yannick Jönsson, Isabelle Terrasson, and Naïma Kheddoumi for technical assistance.

Competing interests:

Mazzone has submitted a world-wide patent pending for diagnostic use of gene expression profiles in monocytes. All other authors declare no conflict of interest.

Funding/Support:

Hamm was funded by the Deutsche Forschungsgemeinschaft (DFG), Prenen by the Leuven University Hospitals Clinical Research Foundation, Rothé by Actions de Recherche Concertée (ARC). This work was supported by grants from the European Research Council (OxyMo to Mazzone), the Fournier-Majoie Foundation (FFMI), FWO (G.0.793.11.N.10), Belgian Foundation Against Cancer (2010-198) and Italian Association for Cancer Research (AIRC 12214).

Role of the funding sources:

The funders had no role in the design and conduct of the study; collection, management, analysis, and interpretation of the data; preparation, review, or approval of the manuscript; and decision to submit the manuscript for publication.

FIGURE LEGENDS

Figure 1: Flowchart of patient inclusion and sample analysis

Inclusion criteria for patients were sporadic histologically confirmed adenocarcinoma of the colon and/or rectum for cohort I-III and VI, patients in remission from CRC for a treatment-free interval of minimum 3 months for cohort V, histologically confirmed adenocarcinoma of the stomach or gastroesophageal junction or of the pancreas, or histologically confirmed gastritis for cohort IV.

Figure 2: Development and validation of a gene signature in circulating monocytes for diagnosis of CRC

a, **b**, Differentially expressed genes between all CRC patients (P,PM) and healthy volunteers (HV). The MA plot (**a**) shows the fold change versus the average expression intensity, while the Volcano plot (**b**) shows fold change in relation to the p values. Green, significantly downregulated genes; red, significantly upregulated genes; corrected p<0.05. **c**, Final gene signature for diagnosis of CRC, comprised of 23 genes, validated in a multicentric test set of patients. Expression levels are displayed as expression relative to the HV mean; boxes, first to third quartile; Whiskers, range; dots, values outside 1.5-times the interquartile distance; horizontal line, median; +, mean. *, p<0.05; **, p<0.01; ***, p<0.001. **d**, ROC analysis for P,PM versus HV in multicentric cohort II. **e**, ROC analysis for P,PM versus HV (negative to screening colonoscopy) in cohort III. See Supplementary Methods for classification approaches.

Figure 3: Tumour-released soluble factors induce the specific upregulation of the gene signature

a-d, Stimulating freshly isolated, naïve monocytes with medium containing soluble factors demonstrates that the genetic fingerprint in monocytes used for the diagnostic gene signature is specifically induced by the transformed colon epithelium (HCT) in comparison to a benign cell line (CCD), as demonstrated by expression analysis comparing selective marker genes in stimulated monocytes to mock control (**a**). Genetic alterations are independent of hypoxic cues (**b**). The gene signature is specific to CRC in comparison to monocytes stimulated by a gastric cancer cell line (MKN) (**c**). The gene signature is reverted after withdrawal of the stimulus *i.e.*, the conditioned medium (**d**). n=6 (biological replicates from 6 different healthy donors); bars, mean with SEM; *, p<0.05; **, p<0.01; ***, p<0.001; ****, p<0.0001; #, p<0.05 towards mock control, assessed by ANOVA with Bonferroni correction. All experiments were repeated at least twice.

Figure 4: The diagnostic gene signature is specific for CRC of all stages and reverts upon curative treatment

a, Expression of the gene signature in patients with cancer of the gastro-esophageal junction (PG), demonstrating no upregulation and thus specificity of the diagnostic signature for CRC. See Figure 2 for details on graphic elements. **b**, Gene signature in patients after curative treatment (patients in remission, PR), in which the expression levels revert to those of healthy volunteers in comparison to CRC patients. **c**, **d**, ROC analyses corresponding to Figure 4a. **e**, Four patients with isolated metastatic recurrence at the time of analysis (black dots) in a 2D-projection of the multi-gene expression levels. The gene signature of metachronously metastasized patients clusters with those patients with primary tumours (red), distinct from healthy individuals (blue). *, p<0.05; **, p<0.01; ***, p<0.001

Figure 5: Specificity of the gene signature to monocytes in comparison to **PBMCs**

a, Expression study assessing the gene signature in PBMCs in comparison to monocytes (PBMs). While the entire signature is confirmed in PBMs in this independent sample set, it is impossible to detect robust genetic alterations in s. S. pis in PBN. ng the previously "", p<0.001 PBMCs, demonstrating specificity to PBMs. See Figure 2 for details on graphic elements. **b**, Corresponding ROC analysis in PBMCs. **c**, ROC analysis of P,PM versus HV in monocytes, confirming the previously established classification performance. *, p<0.05; **, p<0.01; ***, p<0.001

Gut

REFERENCES

- 1. Siegel R, Naishadham D, Jemal A. Cancer statistics, 2012. CA: a cancer journal for clinicians 2012;**62**(1):10-29.
- 2. Weitz J, Koch M, Debus J, et al. Colorectal cancer. Lancet 2005;**365**(9454):153-65.
- Lieberman DA. Clinical practice. Screening for colorectal cancer. The New England journal of medicine 2009;361(12):1179-87.
- Stoop EM, de Haan MC, de Wijkerslooth TR, et al. Participation and yield of colonoscopy versus non-cathartic CT colonography in population-based screening for colorectal cancer: a randomised controlled trial. The lancet oncology 2012;13(1):55-64.
- Quintero E, Castells A, Bujanda L, et al. Colonoscopy versus fecal immunochemical testing in colorectal-cancer screening. The New England journal of medicine 2012;366(8):697-706.
- Pawa N, Arulampalam T, Norton JD. Screening for colorectal cancer: established and emerging modalities. Nature reviews Gastroenterology & hepatology 2011;8(12):711-22.
- Murdoch C, Muthana M, Coffelt SB, et al. The role of myeloid cells in the promotion of tumour angiogenesis. Nat Rev Cancer 2008;8(8):618-31.
- Shi C, Pamer EG. Monocyte recruitment during infection and inflammation. Nature reviews Immunology 2011;11(11):762-74.
- Sandel MH, Dadabayev AR, Menon AG, et al. Prognostic value of tumor-infiltrating dendritic cells in colorectal cancer: role of maturation status and intratumoral localization. Clin Cancer Res 2005;11(7):2576-82.

- 10. Sica A, Mantovani A. Macrophage plasticity and polarization: in vivo veritas. The Journal of clinical investigation 2012;**122**(3):787-95.
- 11. Wynn TA, Chawla A, Pollard JW. Macrophage biology in development, homeostasis and disease. Nature 2013;**496**(7446):445-55.
- Irvine KM, Gallego P, An X, et al. Peripheral blood monocyte gene expression profile clinically stratifies patients with recent-onset type 1 diabetes. Diabetes 2012;61(5):1281-90.
- Zawada AM, Rogacev KS, Schirmer SH, et al. Monocyte heterogeneity in human cardiovascular disease. Immunobiology 2012;217(12):1273-84.
- 14. Ma H, Hong M, Duan J, et al. Altered cytokine gene expression in peripheral blood monocytes across the menstrual cycle in primary dysmenorrhea: a case-control study. PloS one 2013;8(2):e55200.
- 15. Honda T, Inagawa H, Yamamoto I. Differential expression of mRNA in human monocytes following interaction with human colon cancer cells. Anticancer research 2011;**31**(7):2493-7.
- Fletcher RH. Carcinoembryonic antigen. Annals of internal medicine 1986;104(1):66-73.
- 17. Schwarzenbach H, Hoon DS, Pantel K. Cell-free nucleic acids as biomarkers in cancer patients. Nature reviews Cancer 2011;**11**(6):426-37.
- Huang Z, Huang D, Ni S, et al. Plasma microRNAs are promising novel biomarkers for early detection of colorectal cancer. International journal of cancer Journal international du cancer 2010;**127**(1):118-26.
- 19. Church TR, Wandell M, Lofton-Day C, et al. Prospective evaluation of methylated SEPT9 in plasma for detection of asymptomatic colorectal cancer. Gut 2013.

- Gut
- 20. Xu Y, Xu Q, Yang L, et al. Gene expression analysis of peripheral blood cells reveals toll-like receptor pathway deregulation in colorectal cancer. PloS one 2013;8(5):e62870.
- 21. Han M, Liew CT, Zhang HW, et al. Novel blood-based, five-gene biomarker set for the detection of colorectal cancer. Clinical cancer research : an official journal of the American Association for Cancer Research 2008;**14**(2):455-60.
- Marshall KW, Mohr S, Khettabi FE, et al. A blood-based biomarker panel for stratifying current risk for colorectal cancer. International journal of cancer Journal international du cancer 2010;**126**(5):1177-86.
- 23. Nichita C, Ciarloni L, Monnier-Benoit S, et al. A novel gene expression signature in peripheral blood mononuclear cells for early detection of colorectal cancer. Alimentary pharmacology & therapeutics 2014;**39**(5):507-17.
- Imperiale TF, Ransohoff DF, Itzkowitz SH, et al. Multitarget stool DNA testing for colorectal-cancer screening. The New England journal of medicine 2014;**370**(14):1287-97.
- Nyugen J, Agrawal S, Gollapudi S, et al. Impaired functions of peripheral blood monocyte subpopulations in aged humans. Journal of clinical immunology 2010;**30**(6):806-13.
- 26. Smyth GK. Linear models and empirical bayes methods for assessing differential expression in microarray experiments. Statistical applications in genetics and molecular biology 2004;**3**:Article3.
- 27. 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(1):289-300.

28. Sample size for microarray experiments. Secondary Sample size for microarray experiments. <u>http://bioinformatics.mdanderson.org/MicroarraySampleSize/</u>.

29. Burges CJC. A Tutorial on Support Vector Machines for Pattern Recognition. Data Min Knowl Discov 1998;**2**(2):121-67.

30. Breiman L. Random Forests. Mach Learn 2001;45(1):5-32.

- Dietterich TG. Ensemble Methods in Machine Learning. Proceedings of the First International Workshop on Multiple Classifier Systems: Springer-Verlag, 2000:1-15.
- 32. Wessels LF, Reinders MJ, Hart AA, et al. A protocol for building and evaluating predictors of disease state based on microarray data. Bioinformatics 2005;21(19):3755-62.
- Edgar R, Domrachev M, Lash AE. Gene Expression Omnibus: NCBI gene expression and hybridization array data repository. Nucleic acids research 2002;30(1):207-10.
- 34. Piehler A, Grimholt R, Ovstebo R, et al. Gene expression results in lipopolysaccharide-stimulated monocytes depend significantly on the choice of reference genes. BMC Immunology 2010;**11**(1):21.
- 35. Guo C, Liu S, Wang J, et al. ACTB in cancer. Clinica chimica acta; international journal of clinical chemistry 2013;**417**:39-44.
- 36. Jones S, Chen WD, Parmigiani G, et al. Comparative lesion sequencing provides insights into tumor evolution. Proceedings of the National Academy of Sciences of the United States of America 2008;**105**(11):4283-8.
- Goede V, Coutelle O, Shimabukuro-Vornhagen A, et al. Analysis of Tie2expressing monocytes (TEM) in patients with colorectal cancer. Cancer investigation 2012;30(3):225-30.
<text><text><text><text> 38. Schauer D, Starlinger P, Reiter C, et al. Intermediate monocytes but not TIE2-

Gut

Tumour-Educated Circulating Monocytes are Powerful Candidate Biomarkers for Diagnosis and Disease Follow-up of Colorectal Cancer

Alexander Hamm^{1,2#}, Hans Prenen^{3#}, Wouter Van Delm^{4#}, Mario Di Matteo^{1,2}, Mathias Wenes^{1,2}, Estelle Delamarre^{1,2}, Thomas Schmidt⁵, Jürgen Weitz^{5,6}, Roberta Sarmiento⁷, Angelo Dezi⁷, Giampietro Gasparini⁷, Françoise Rothé⁸, Robin Schmitz⁵, André D'Hoore⁹, Hannes Iserentant¹⁰, Alain Hendlisz⁸ & Massimiliano Mazzone^{1,2}

¹Lab of Molecular Oncology and Angiogenesis, Vesalius Research Center, VIB, Leuven, Belgium ²Lab of Molecular Oncology and Angiogenesis, Vesalius Research Center, Department of Oncology, KU Leuven, Leuven, Belgium

³Digestive Oncology, University Hospitals Leuven and Department of Oncology, KU Leuven, Leuven, Belgium

⁴Nucleomics Core, VIB, Leuven, Belgium

⁵Department of General, Visceral, and Transplantation Surgery, University of Heidelberg, Heidelberg, Germany

⁶Department of Visceral, Thoracic, and Vascular Surgery, University Hospital Carl Gustav Carus,

Technical University Dresden, Dresden, Germany

⁷Department of Oncology, San Filippo Neri, Rome, Italy

⁸Medical Oncology Clinic, Institut Jules Bordet, Brussels, Belgium

⁹Department of Abdominal Surgery, University Hospitals Leuven, KU Leuven, Leuven, Belgium

¹⁰VIB, Zwijnaarde, Belgium

[#]contributed equally to this study

Correspondence:

Massimiliano Mazzone, <u>massimiliano.mazzone@vib-kuleuven.be</u>, Tel: +32-16-373213, Fax +32-16-372585 VIB Vesalius Research Center, KU Leuven, Herestraat 49, Bus 912, 3000 Leuven, Belgium

Hans Prenen, hans.prenen@uzleuven.be, Tel: +32-16-340238

Word Count: 4127

Key Words: Monocytes, colorectal cancer, screening, inflammation

https://mc.manuscriptcentral.com/gut

LIST OF ABBREVIATIONS

AUC	area under the curve
BER	balanced error rate
CEA	carcino-embryonic antigen
CRC	colorectal cancer
ENS	ensemble method
FIT	fecal immunochemical test
FOBT	fecal occult blood test
MACS	magnet-associated cell sorting
MCCV	Monte Carlo cross validation
NSAID	non-steroid anti-inflammatory drugs
PBM	peripheral blood monocytes
PBMC	peripheral blood mononuclear cells
<u>qPCR</u>	quantitative RT-PCR
RF	random forest
ROC	receiver operating characteristics
RT-PCR	reverse-transcription polymerase chain reaction
Se	sensitivity
SGMV	single gene majority vote
Sp	specificity
SVM	support vector machine
UICC	Union internationale contre le cancer
Labels of pa	atient groups:
HV	healthy volunteer
<u>P</u>	non-metastatic CRC patient

P. PM non-metastatic and metastatic CRC patients

- PC pancreatic cancer patient
- PG gastric cancer patient
- PGT gastritis patient

ΡM metastatic CRC patient patient in remission from CRC PR

 Gut

ABSTRACT

Objective: Cancer immunology is a growing field of research <u>whose aim is</u> to develop innovative therapies and diagnostic tests. <u>Starting from the hypothesis that immune</u> <u>cells promptly respond to harmful stimuli, we utilized</u> peripheral blood monocytes (PBM) <u>in order</u> to characterize a distinct gene expression profile and to evaluate its potential as a candidate diagnostic biomarker <u>in colorectal cancer (CRC) patients, a</u> <u>still unmet clinical need.</u>

Design: We performed a case-control study including <u>360</u> PBM samples from four European oncological centres <u>and defined</u> a gene expression profile specific to CRC. The robustness of the genetic profile and disease specificity, were assessed in an independent setting.

Results: This screen returned 43 putative diagnostic markers, which we refined and validated in the confirmative multicentric analysis to 23 genes with outstanding Se=100.0% [100.0%:100.0%]. diagnostic accuracy (AUC=0.99 [0.99:1.00]. Sp=92.9% [78.6%;100.0%] in multiple-gene ROC analysis). The diagnostic accuracy was robustly maintained in prospectively collected independent samples (AUC=0.95 [0.85;1.00], Se=92.6% [81.5%;100.0%], Sp=92.3% [76.9%;100.0%]. This monocyte signature was expressed at early disease onset, remained robust over the course of disease progression, and was specific for the monocytic fraction of mononuclear cells. The gene modulation was induced specifically by soluble factors derived from transformed colon epithelium in comparison to normal colon or other cancer histotypes. Moreover, expression changes were plastic and reversible, as they were abrogated upon withdrawal of these tumour-released factors. Consistently, the modified set of genes reverted to normal expression upon curative treatment and was specific for CRC.

 <text><text><text><text>

SUMMARY BOX

What is already known on this subject?

- Early diagnosis of colorectal cancer is crucial for curative surgical treatment, highlighting the need for efficient screening tools.
- Colorectal cancer screening is a rapidly evolving field, as several strategies for supplementing the invasive colonoscopic screening are explored.
- Circulating cells of the immune system in the blood stream are easily accessible, yet understudied with regard to their precise role in tumour immunology.
- Tumour-associated macrophages deriving from circulating monocytes can display diverse phenotypes and affect tumour growth and metastasis by different means, depending on the cellular context.

What are the new findings?

- Monocytes are plastic cells that are modified by early occurrence of colorectal cancer, resulting in a highly specific genetic fingerprint, which is independent of tumour stage.
- The changes in monocyte expression profiles are reversible, highly specific to the tissue type and cancer histotype, and induced in response to soluble factors released by the cancer cells in the primary or metastatic site.
- The specific genetic fingerprint in circulating monocytes can be harnessed for diagnosis and disease follow-up of colorectal cancer.

How might it impact on clinical practice in the foreseeable future?

> If the initiated prospective validation study supports our sound results, our

INTRODUCTION

Colorectal cancer (CRC) is the second leading cause of cancer-related deaths in the US¹. Its incidence and the difficulty in early-diagnosis make CRC a primary focus in the oncology community². Early CRC is symptomless, and, consequently, is frequently diagnosed when already advanced. Metastatic disease (found in 30 to 40% of CRC patients) is associated with a poor 5-year survival rate of less than 10%. In contrast, up to 80% of patients can be cured by early tumour resection, rendering timely diagnosis a crucial factor for proper disease management². Nevertheless, endoscopic screening as well as stool tests (fecal immunochemical test, FIT, or fecal occult blood test, FOBT⁵) are not widely accepted by the target population, while the socioeconomical burden of these procedures is high². Thus, there is urgent need to identify specific, non-invasive biomarkers for early CRC diagnosis and treatment monitoring to avoid disease progression to advanced stages that are difficult to cure⁶. Peripheral blood is one of the least invasive sample sources that can be intensively screened for CRC biomarkers. Within the blood stream, peripheral blood monocytes (PBM) represent a reservoir of inflammatory cells that contribute to disease progression by different means^{7 8}. These cells are recognized to be plastic and versatile cells, which can change their phenotype in response to microenvironmental stimuli, yielding either tumouricidal or pro-tumourigenic features depending on the stromal context or tumour type¹⁰¹¹. Interestingly, recent studies have suggested distinct expression profiles in circulating monocytes in several pathological conditions such as diabetes¹², atherosclerosis¹³, and dysmenorrhea¹⁴, though none have convincingly demonstrated a specific regulation of monocyte heterogeneity by malignantly transformed cells apart from descriptive studies in vitro on monocytic cell lines¹⁵.

Several novel accessible diagnostic tools share the major opportunity to make frequent screening more appealing to a greater number of patients, as a less invasive method is likely to increase compliance and allow for decreased screening intervals (recently comprehensively reviewed⁶). While conventional blood-based tumour markers (particularly carcino-embryonic antigen, CEA¹⁶) have been established as supplemental markers in treatment monitoring, they have failed to yield high diagnostic accuracy as primary screening tools. In addition to the established FIT or FOBT⁵, other potential diagnostic markers include serumassociated biomarkers (e.g. circulating tumour DNA¹⁷, micro-RNA¹⁸, methylation markers like SEPT9¹⁹), genetic marker sets in white blood cells²⁰⁻²³, and, most recently, fecal tumour DNA²⁴. However, all of these approaches display limited sensitivity and specificity⁶. In this study, we therefore assess the sensitivity and specificity of a novel gene signature in circulating monocytes for the diagnosis of CRC in comparison to healthy individuals or to other cancer types, and assess its robustness in prospectively obtained samples.

https://mc.manuscriptcentral.com/gut

PATIENTS AND METHODS

Patients

We collected a total of <u>360</u> samples between January <u>13</u>, 2010 and <u>January <u>26</u>, <u>2015</u>, comprised of the following cohorts: cohort I (genome-wide screening in <u>27</u> patients with non-metastatic stage I, stage II, or stage III CRC (P), <u>28</u> patients with metastatic stage IV CRC (PM), and <u>38</u> healthy volunteers (HV) (without history or evidence of acute or chronic disease)), cohort II (multicentric validation in <u>73</u> patients and <u>61</u> healthy volunteers <u>from four different oncological centres</u>), cohort III (robustness assessment in <u>27</u> patients and <u>13</u> asymptomatic healthy individuals with colonoscopy-confirmed absence of disease), cohort IV (<u>15 patients with gastric cancer (PG)</u>, <u>16 patients with pancreatic cancer (PC)</u>, <u>10 patients with gastritis (PGT)</u>, <u>all treatment-naïve</u>, <u>and <u>13 HV</u></u>), cohort V (<u>15 curatively treated patients</u>), and cohort VI (comparative expression analysis in PBM and PBMC in <u>17 patients and 7 healthy volunteers</u>). See Figure 1 for allocation of collected samples to analyses. <u>All participants gave written informed consent</u>, <u>and the study was approved by the respective institutional review boards</u>. <u>Details on inclusion and exclusion criteria</u>, <u>participating centers and ethical approval can be found in Supplementary Methods</u>.</u>

Identification of a gene signature

Genome-wide expression analysis was performed on the Illumina platform (Illumina) on RNA obtained from peripheral blood monocytes (PBM), isolated by a two-step procedure with density gradient centrifugation and positive selection for CD14 using the MACS system (Miltenyi). Details are reported in Supplementary Methods. Differential expression was assessed with the limma package of R²⁶. Putative candidate genes were confirmed on a random subset of cohort I and validated by

Gut

quantitative RT-PCR (qPCR) on the 7500Fast System (Applied Biosystems) using intron-spanning PrimeTime qPCR Assays (Integrated DNA Technologies) listed in Supplementary Table 1 as described in Supplementary Methods. For statistical analysis, we followed a three-step top-down approach to construct a gene signature for CRC, with details explained in Supplementary Methods.

Multicentric validation study

For validation of a diagnostic test, we used cohort II to train and validate a multi-gene classifier. Splits in training and test sets for validation were performed by stratified random sampling for centre of origin and class label as detailed out in Supplementary Methods. Samples with missing values for more than 25% of the genes were excluded from the analysis. We ruled out an effect of the class labeling on the percentage of missing values with Fisher's exact test (Supplementary Table 2).

The training dataset was used to build three types of classifiers: a support vector machine (SVM)²⁹ with linear kernel, a single-gene majority vote (SGMV) classifier, and a random forest classifier (RF³⁰). Subsequently, we applied an ensemble method³¹ that votes according to the majority of the three independent classifiers. Performance was validated both with ranking (AUC) and classification (balanced error rate, BER, Se, Sp) scores with 95% confidence intervals ([lower boundary; upper boundary]). We explicitly opted for relatively simple computational models in order to limit chances of over-fitting the training data and to maximize interpretability of the models' internal decision-making process. Model flexibility was further controlled through a Monte-Carlo cross-validation scheme (MCCV)³², before final estimation of the model parameters. Validation of the predictive models was done on

Gut

the test set of cohort II, which were not included during development of the models. Details on all classification methods are specified in Supplementary Methods.

In order to avoid biased conclusions, the analysis of the 23 genes was complemented with a study by an independent team (DNAlytics, Belgium) that adopted a slightly modified analysis protocol (see Supplementary Methods). All complementary analyses were performed in R with scripts designed by DNAlytics, fully independently from other analyses described in this paper.

In vitro model system

To study the effects of tumour-released soluble factors on gene expression in monocytes, we established an *in vitro* model system, where monocytes from healthy , son. donors were challenged with tumour-released soluble factors and changes in gene expression profile were analyzed by qPCR. See Supplementary Methods for details.

https://mc.manuscriptcentral.com/gut

Gut

RESULTS

Establishment of putative biomarkers by genome-wide expression analysis

<text><text><text> To obtain a set of putative biomarkers that might facilitate early diagnosis of CRC, we have performed a genome-wide expression analysis on PBMs from 55 untreated patients newly diagnosed with CRC and 38 healthy volunteers (cohort I). All relevant clinicopathological information on patient cohorts can be found in Table 1.

Cohort			I										V	VI	
	P,PM	HV	P,PM HV							P,PM	HV	Р	P,PM HV		
			LEU ^a	HD⁵	SFN ^c	IJB ^d	LEU ^a	HD⁵	SFN℃	IJB₫					
Number of samples	55	38	39	19	10	5	20	12	14	15	27	13	15	17	7
Age															
median	67	55	66	69	72	59	49	55	47	49	66	62	69	78	42
range	44-87	42-79	47-78	42-76	50-85	52-82	42-69	46-75	40-63	42-62	44-90	43-74	45-81	62-89	42-57
Gender															
male	22	15	24 🤇	11	5	1	15	7	11	2	14	8	8	11	5
female	33	23	15	8	5	4	5	5	3	13	13	5	7	6	2
metastatic	28	1	16	3	2	2	1	/	/	/	16	/	0	6	/
non-metastatic	27	/	23	16	8	3	/	/	/	/	11	/	15	11	/
UICC stage															
1	3	/	7	2	1	1	1	/	/	/	2	/	4	2	/
2	12	1	8	8	2	0	1	/	/	1	3	/	7	7	/
3	12	1	8	6	5	2	1	/	/	1	6	/	4	2	/
4	28	1	16	3	2	2		1	/	1	16	/	0	6	/
Tumour localization															
Caecum	5	1	3	3	0	0	1		/	/	2	/	1	2	/
Ascendens	11	1	4	3	1	0	1	Ι	1	/	6	/	4	4	/
Transversum	0	1	4	3	0	0	1			1	2	/	1	0	/
Descendens	4	1	2	1	3	3	1	1		▲ 1	0	/	0	1	/
Sigmoid	28	1	15	3	2	0	1	1	1	1	10	1	5	6	/
Rectum	6	1	8	5	3	2	1	1	1		7	1	4	2	/
Double	1	1	3	1	1	0	1	/	/		0	/	0	2	/

^aLeuven, ^bHeidelberg, ^cRome, ^dBrussels. See Supplementary Methods for the detailed description of contributing centres

The purity of the monocyte fraction was >90%, as assessed by FACS analysis in the pilot phase (Supplementary Figure 1a) and verified by hemocytometric analysis for each individual sample (Supplementary Figure 1b). Both absolute and relative monocyte counts were not different between patients and healthy volunteers (Supplementary Figure 1c). We therefore investigated differentially expressed genes by genome-wide expression analysis using the Illumina HumanHT-12 v4 Expression BeadChip Kit. The data discussed in this publication have been deposited in NCBI's Gene Expression Omnibus³³ and are accessible through GEO Series accession number GSE47756

(http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?token=hvmpvoswuqaeybc&acc=GSE 47756). In first instance, we compared the average expression values of all CRC patients (P,PM), comprised of non-metastatic (P) and metastatic (PM) patients, to that of healthy volunteers (HV). The resulting gene signature of (P,PM) versus HV consisted of 36 upregulated and 4 downregulated probes (Figure 2a, b, Table 2). In second instance, we were interested if the gene signature in patients with synchronous metastases *i.e.*, at the time of diagnosis (PM, n=28) was different from that in non-metastatic patients (P, n=27). Interestingly, the number of up- and downregulated genes was comparable in both P and PM (in comparison to HV) (Table 2 and Supplementary Figure 2a, b), while there were no genes found to be differentially expressed between the two patient groups (Supplementary Figure 2a, b), indicating that the gene signature induced at early onset stays robust over disease progression. Indeed, when post-hoc assessing those samples from patients with early stages (Tis and T1), they clustered with the rest of the patient samples (data not shown). A power analysis revealed that, for the given number of genes, samples and observed variation, chances were very low $(<10^{-10})$ that truly differentially expressed genes with

Gut

fold changes larger than 1.5 had been missed. Therefore, adding more samples would probably have changed little to the panel of candidate genes that our screen returned.

Confirmation of the gene signature in independently processed samples

To validate the genetic signature, we performed quantitative RT-PCR (qPCR) analysis on a random subset of PBM from 8 samples of each of the three groups (P, PM, and HV), normalizing to reference gene *B2M*, which was selected after an extensive screening procedure (Supplementary Note 1). To avoid bias in the confirmation procedure, we freshly extracted RNA from independently stored samples for confirmative expression analysis. In analyzing 43 putative marker genes with probes listed in Supplementary Table 1, 23 genes showed differential expression between (P, PM) and HV (Supplementary Figure 3b, Table 2, and Supplementary Table 4). Thus, we were able to confirm a subset of the previously established gene signature, independent of the RNA extraction and the platform used for expression analysis. Information on the annotated biological function of the genes of the diagnostic signature can be found in Supplementary Table 5 and Supplementary Note 2.

Confirmation of the gene signature in a multicentric validation set

For a rigorous validation of the gene signature, we collected an independent multicentric validation set (cohort II) from a total of 4 different European oncological centres with stratified training and test sets as described in Supplementary Methods. Using the panel of 23 genes confirmed previously, we found consistently differential expression between all patients and the healthy volunteers (Figure 2c and

Supplementary Figure 4). In line with the findings from the screening phase, there were no detectable differences in expression levels between P and PM (Supplementary Figure 5), while either patient group alone compared to HV was differentially expressed (data not shown).

In ROC analysis for single genes, we found that some, but not all of the genes that displayed significantly differential expression were able to discriminate patient samples from healthy individual samples with acceptable AUCs (Supplementary Figure 6 and data not shown). We therefore hypothesized that a marker panel consisting of multiple genes might yield better results in discriminating sample identity. To address this question, we decided to test three different classification algorithms on this data set, namely a support vector machine (SVM)²⁹ with linear kernel, a single-gene majority vote (SGMV) classifier, a random forest classifier (RF³⁰), and a combined classification by an ensemble method³¹, using the outcome of the three classification algorithms for a final diagnostic decision. <u>To limit overestimation of the performance by the particular training and test set, we performed a MCCV as a conservative estimate with 1,000 cross-validations. Performance of all classification algorithms in cohorts II – VI, including the conservative estimate of the MCCV in cohort II, is given in detail in Table 2.</u>

TABLE 2: PERFORMANCE SCORES OF MULTIGENE CLASSIFIER

	<u>SGMV</u>	<u>SVM</u>	<u>RF</u>	<u>ENS</u>			
Cohort II (Validation)							
AUC [95% CI]	<u>0.99 [0.99;1.00]</u>	1.00 [1.00;1.00]	0.99 [0.97;1.00]	0.99 [0.99;1.00]			
BER [%]	3.6	3.3	3.6	3.6			
Sensitivity [95% CI]	100 [100;100]	93.3 [80.0;100]	100 [100;100]	100 [100;100]			
Specificity [95% CI]	92.9 [78.6;100]	100 [100;100]	92.9 [78.6;100]	92.9 [78.6;100]			
Cohort II (MCCV)							
AUC [95% CI]	<u>0.94 [0.86;1.00]</u>	0.92 [0.83;0.99]	0.93 [0.83;1.00]	0.86 [0.72;0.99]			
BER	13.3	20.0	13.3	13.3			
Sensitivity [95% CI]	80.0 [60.0;100]	66.7 [20.0;93.3]	86.7 [60.0;100]	80.0 [60.0;100]			
Specificity [95% CI]	93.3 [66.7;100]	93.3 [80.0;100]	93.3 [73.3;100]	93.3 [80.0;100]			
Cohort III							
AUC [95% CI]	0.96 [0.89;0.99]	0.91 [0.80;0.99]	0.93 [0.79;1.00]	0.95 [0.85;1.00]			
BER	7.7	15.0	7.6	7.6			
Sensitivity [95% CI]	100 [100;100]	77.8 [59.3;92.6]	92.6 [81.5;100]	92.6 [81.5;100]			
Specificity [95% CI]	84.6 [61.5;100]	92.1 [76.9;100]	92.3 [76.9;100]	92.3 [76.9;100]			
Cohort IV (gastric cancer)							
Sensitivity [95% CI]	33.3 [13.3;60.0]	26.7 [6.7;46.7]	20.0 [0.0;40.0]	20.0 [0.0;40.0]			
Cohort IV(pancreatic cance	<u>er)</u>						
Sensitivity [95% CI]	0.0 [0.0;0.0]	<u>0.0 [0.0;0.0]</u>	<u>0.0 [0.0;0.0]</u>	0.0 [0.0;0.0]			
Cohort IV (gastritis)							
Sensitivity [95% CI]	<u>10.0 [0.0;30.0]</u>	<u>10.0 [0.0;30.0]</u>	<u>10.0 [0.0;30.0]</u>	<u>10.0 [0.0;30.0]</u>			
Cohort V (PR)							
Sensitivity [95% CI]	<u>50.0 [20.0;80.0]</u>	10.0 [0.0;30.0]	<u>20.0 [0.0;50.0]</u>	20.0 [0.0;50.0]			
Cohort VI (PBMC)							
<u>AUC [95% CI]</u>	<u>0.51 [0.19;0.80]</u>	0.44 [0.13;0.74]	0.64 [0.31;0.94]	<u>0.44 [0.19;0.66]</u>			
BER	<u>59.3</u>	49.3	52.1	52.1			
Sensitivity [95% CI]	<u>10.0 [0.0;30.0]</u>	30.0 [10.0;60.0]	<u>10.0 [0.0;30.0]</u>	<u>10.0 [0.0;30.0]</u>			
Specificity [95% CI]	71.4 [28.6;100]	71.4 [42.5;100]	85.7 [57.1;100]	85.7 [57.1;100]			
Cohort VI (PBM)							
<u>AUC [95% CI]</u>	<u>1.00 [1.00;1.00]</u>	<u>0.79 [0.54;1.00]</u>	<u>1.00 [1.00;1.00]</u>	<u>1.00 [1.00;1.00]</u>			
BER	<u>0.0</u>	<u>30.8</u>	<u>0.0</u>	<u>0.0</u>			
Sensitivity [95% CI]	<u>100 [100;100]</u>	<u>38.5 [15.4;61.5]</u>	<u>100 [100;100]</u>	<u>100 [100;100]</u>			
Specificity [95% CI]	<u>100 [100;100]</u>	<u>100 [100;100]</u>	<u>100 [100;100]</u>	100 [100;100]			

Listed are the performance scores of all multi-gene classifiers (SGMV, SVM, RF) and their combined ensemble method (ENS) of all different cohorts – please see methods for details. ^aSensitivity for labeling a gastric cancer sample as CRC

^bSensitivity for labeling a curatively treated patient in full remission as CRC

Strikingly, we achieved a remarkably high AUC of 0.99 [0.99;1.00] with a BER of 3.6% (Figure 2d and Table 2), translating into a sensitivity of 100.0% [100.0%;100.0%] and a specificity of 92.9% [78.6%;100.0%] (Table 2). <u>Neither of the classification algorithms was capable of separating P from PM or detect differences dependent on tumour localization (Supplementary Note 3).</u>

In order to assess whether the diagnostic gene signature is actually suitable for diagnosis of CRC in a screening setting, we have initiated a prospective sample collection in both patients and healthy individuals who are subjected to colonoscopy. In a pilot analysis in 27 patients (newly diagnosed with CRC by screening colonoscopy) and 13 healthy individuals negative to screening colonoscopy (cohort III), we found an AUC of 0.95 [0.85;1.00] with a BER of 7.6%, yielding a sensitivity of 92.6% [81.5%;100.0%] and a specificity of 92.3% [76.9%;100.0%] (Figure 2e and Table 2). The complementary data analysis (independently performed by DNAlytics, Belgium) on the same panel of 23 genes led to the matching conclusions in terms of performance. The first experiment consisted in cross-validating a model on Cohort II (BER: 8.4% [3.4%;13.4%]; AUC: 0.93 [0.88;0.98]). A second experiment consisted in learning the same type of model on Cohort II and having it make predictions on Cohort III (BER: 13.2%; AUC: 0.92).

Soluble factors released by colorectal cancer cells induce <u>an early, tumour</u> <u>type-specific and reversible genetic fingerprint in monocytes</u>

We hypothesized that tumour-released soluble factors are the key players in inducing the genetic signature in circulating monocytes. Thus, we established an *in vitro* model system where we cultured freshly isolated human monocytes from healthy donors in different conditions. In order to assess alterations in gene expression, we

first analyzed which of the 23 genes comprising the gene signature was up- or downregulated in culture after 72 hours without any additional stimulus and excluded these from the further *in vitro* studies (Supplementary Figure 7). Out of the remaining gene signature, the majority (7/9) was specifically upregulated when culturing naïve monocytes in medium conditioned by the CRC cell line HCT116, while expression levels were not affected by mock medium (Figure 3a). Moreover, in line with the coherent induction of the specific signature independent of the stage of the disease, the induction *in vitro* was independent of hypoxic cues, as HCT116-conditioned medium in hypoxia did not induce any different expression levels than medium obtained in normoxia (Figure 3b). Likewise, the changes in expression levels of all these genes occurred already 18 hours after stimulating monocytes with the conditioned medium, consistent with the fact that already early stages are detectable by the diagnostic signature.

To rule out an off-target effect of conditioned medium *i.e.*, unspecific cues from cell metabolites, apoptotic bodies, pH, etc., we assessed the expression levels of the genes upregulated by HCT116-conditioned medium in comparison to a benign colon epithelium cell line, CCD 841 CoN (CCD), which did not induce alterations in gene expression levels different from the Mock control (Figure 3a).

Prompted by this finding, we investigated if the induction of the genetic signature was a general effect of malignant transformation or might be specific to the histotype of cancer. To address this question, we conditioned medium with a gastric cancer cell line, MKN-45 (MKN), to compare CRC to another frequent gastrointestinal solid neoplasm. Remarkably, when comparing the expression levels in naïve monocytes upon stimulation with the different conditioned media, we found that MKN-45

conditioned medium did not induce the same upregulation of the genes of interest as HCT116 conditioned medium (Figure 3c).

As immune cells are highly versatile and plastic cells mirroring the microenvironment, where they are embedded, we reasoned that the genetic signature induced by CRC in monocytes might be dependent on the continuous presence of the stimulating agents and thus be reversible upon inversion of the conditions. We therefore incubated naïve monocytes first with HCT116-conditioned medium for 18 hours and then refreshed the medium with plain culture medium, thus withdrawing the tumour-released soluble factors. Strikingly, the previously elevated expression levels of a set of marker genes were almost entirely reverted to the original (and to the mock control) expression levels 72 hours after withdrawing the tumour-cell conditioned medium (Figure 3d), whereas they remained constantly overexpressed when the conditioned medium was maintained (data not shown).

The monocyte signature is specific for CRC and might serve as a candidate biomarker of disease follow-up

Based on the *in vitro* results showing that the genetic signature is specific to CRC, we sought to confirm these findings *in vivo*. We therefore assessed the diagnostic signature in patients with i. cancer of the stomach and gastro-esophageal junction (PG, n=15) and ii. pancreatic ductal adenocarcinoma (PC, n=16), two other frequent cancers of the gastrointestinal tract¹. In addition, we analysed iii. patients with gastritis (PGT, n=10) in order to compare the gene signature in CRC to a benign inflammatory condition of the gastrointestinal tract (cohort IV). In line with the *in vitro* results we saw that the vast majority of all genes were not significantly different between either of the patient groups and healthy volunteers, indicating the specificity

of this monocyte imprinting by colorectal cancer cells (Figure 4a and data not shown). <u>Moreover, the classifier established to diagnose CRC could not separate</u> patients with gastric cancer (AUC 0.63 [0.48;0.77]), pancreatic cancer (AUC 0.41 [0.27;0.50]), or gastritis (AUC 0.52 [0.35;0.68]) from healthy individuals (Figure 4c, d and Table 2).

The finding that the genetic signature is reverted upon withdrawal of the stimulating agents prompted us to investigate in a pilot phase the behaviour of the entire diagnostic signature in patients upon curative treatment *i.e.*, patients with surgically removed tumours without any evidence of residual disease. To this end, we isolated monocytes from 15 patients of stages I to III treated with curative intent (with or without adjuvant treatment) and presenting at follow-up without detectable residual disease (PR) (cohort V). Here, we found that virtually all of the previously upregulated genes were reverted to expression levels comparable to those of healthy volunteers (Figure 4b). Consequently, when applying the previously established classifier, we found that it was able to distinguish accurately between patients in remission and patients with tumour, while it could not detect differences between patients in remission and healthy volunteers (Table 2).

Finally, as the plasticity of the signature offers the perspective to use the gene signature for follow-up of treated patients, we became interested if the same signature could be used to diagnose relapse (frequently as metachronous metastases rather than local recurrence²). Although our dataset was not powered to address this question with sufficient significance, we post-hoc identified four patients from cohorts I and II included at presentation with metachronous metastases. All four clustered clearly in the group of patients, separately from the healthy volunteers (Figure 4e), suggesting that the signature might be used to detect disease relapse in

line with the previous results that show coherent expression over disease progression.

The gene signature is specific to monocytes in comparison to all peripheral blood mononuclear cells (PBMCs)

To rigorously assess if the genetic fingerprint identified in monocytes was specific to this cell type or an epiphenomenon of genetic shifts in the entire population of PBMCs, we isolated both monocytes and full PBMC fractions from 17 patients and 7 healthy volunteers for a comparative analysis (cohort VI). Interestingly, we found that while in the monocyte population, the diagnostic marker set of 23 genes was upregulated in all patients (both P and PM) in accordance with our previous results (Figure 5a), there were no significant differences in the expression levels of the analyzed genes in the full PBMC compartment when comparing patients to healthy volunteers (Figure 5a). Consistently, applying the previously established classifier with the defined cut-off values, it was impossible to separate the patient group from the healthy volunteer group in PBMC (Figure 5c and Table 2), while the classifier confirmed its accuracy in PBM (Figure 5c and Table 2). Thus, the differential regulation of the gene signature in PBM used for CRC diagnosis is specific to the monocytic lineage, reinforcing our initial working hypothesis that these cells are specifically affected by tumour-secreted factors.

DISCUSSION

The dismal prognosis of CRC can be effectively attenuated by an early and accurate diagnosis, which is however hampered by low compliance rates to the available screening strategies²⁶. With this study, we present a hypothesis-driven approach to screen for specific biomarkers for diagnosis of CRC, which exploits the canonical knowledge on tumour-stroma interactions¹⁰. By using genome-wide expression analysis, we show that a distinct gene signature is detectable in circulating monocytes from CRC patients in comparison to healthy individuals. In fact, this study is the first to demonstrate specific genetic changes in the highly versatile monocyte fraction, mediated by tumour-derived soluble factors. Moreover, we convincingly demonstrate with an in vitro model system that the alterations in gene expression are induced by tumour-released soluble factors, which adds to the value of our biologybound approach in comparison to mere high-throughput screenings. Our comparison of the reported gene signature in monocytes and PBMCs strongly supports our hypothesis that monocytes, more than any other immune cell in circulation, are highly plastic and responsive to microstimuli in the blood. Since the induced expression changes are higher in vitro, it is tempting to speculate that these are dependent on the concentration of cytokines and signals, which remain to be identified.

Interestingly, our analysis indicated that the induced gene signature stays robust over progression of the disease, which is consistent with our *in vitro* findings and not entirely surprising given recent evidence for the molecular similarity between the primary tumour and its metastases³⁶.

<u>The diagnostic gene signature established here proved to be robust independent of</u> <u>the technique (genomewide expression microarray vs. qPCR) and has been</u> validated independently (Supplementary Note 4). Its utilization for diagnosis of CRC

most likely depends on the development of a one-step assay with capture of monocytes from whole blood and gene expression analysis in a multiplex qPCR assay with absolute quantification, avoiding extensive preanalytical processing steps. However, the analytical reliability of this assay needs to be thoroughly established, most likely requiring centralization of the analysis during the first phase of distribution.

The finding that the specific gene signature is reversible if the stimulating cues are withdrawn, was not only demonstrated *in vitro*, but also in a pilot analysis *in vivo* in samples of patients after curative treatment. Although not completely unexpected in view of the plastic nature of the monocyte-macrophage lineage, this analysis opens avenues for treatment monitoring and companion diagnostics and will be assessed in detail in a prospective study during patient follow-up.

If supported by further prospective validation studies, <u>this gene signature may</u> <u>outperform other published non-invasive test for CRC diagnosis⁶ (including single surface markers in monocytes^{37 38})</u> or score similar to the most recent evaluation of fecal tumour DNA²⁴. Moreover, we are the first to demonstrate that a potential diagnostic biomarker obtained in patients at the time of primary diagnosis might also be suitable for disease follow-up <u>and thus assessment of treatment response</u>, owing to its high plasticity.

We acknowledge the limited conclusions that can be drawn from our case-control study. Despite the confirmation in independent samples, we cannot fully exclude possible confounders that can only be unveiled by a blinded, prospective sample collection in screening individuals. These include, but are not limited to, the bias of selecting patients that underwent colonoscopy for a clinical indication; the differences

in age, nutrition status, diet, and potentially lifestyle between patients and healthy volunteers; the unblinded sample collection and processing. It is therefore of paramount importance that a prospective validation study initiated by our group includes screening individuals prospectively with blinded sample processing. In addition, strategies to minimize false negatives and false positives (with potential morbidity resulting from colonoscopy and treatment) will need to be developed. This can be achieved by calculating a risk ratio on the basis of the individual expression profile, which could replace the current binary output (cancer vs. healthy) and thus define groups at risk that need to be subjected to colonoscopy as the gold standard. An informed choice on the thresholds would, at least in first instance, emphasize a high sensitivity at the expense of specificity. The resulting morbidity has to be correlated to the morbidity of screening colonoscopy.

Our study raises important questions, which will need to be addressed in further studies. First, the biological mechanisms and pivotal regulatory pathways in directing the fate of the monocyte gene signature are still unexplored. Of note, only a few genes appear to be commonly upregulated in CRC in comparison to gastric cancer and pancreatic cancer. While this demonstrates specificity for CRC, it also means further studies will be required to identify gene signatures specific to other tumours and possibly benign pathological conditions. Second, we will need to assess if the gene signature is already imprinted in pre-neoplastic lesions (*i.e.*, polyps) and determine the transformation steps at which the specific upregulation occurs. Third, as monocyte plasticity is the starting hypothesis of this study, we will need to assess if treatment regimes (e.g., steroids, chemotherapy, irradiation, postoperative stress conditions) affect the behaviour of the gene expression profile or interfere with its diagnostic capabilities. Fourth, we are currently investigating in a prospective setting

<text>

Gut

CONCLUSIONS

Taken together, these data provide unprecedented evidence that tumour-educated monocytes exhibit a distinct and plastic gene signature, which may not only be <text> suitable for diagnosis of CRC, but potentially allows to monitor for success of therapy or for relapse. As monocytes can be obtained in a non-invasive way, these findings offer exciting new opportunities for both improving CRC diagnosis and enriching the armamentarium of therapeutic strategies, provided that the data obtained here can be replicated in an independent broad screening setting.

ACKNOWLEDGEMENTS

We would like to express our gratitude to all patients and healthy volunteers contributing to our study. The authors are indebted to Joke Allemeersch and Christos Sotiriou for critical advice. We thank DNAlytics (Belgium) for critical independent statistical review of the raw data, <u>Brian Wong for critical review of the manuscript</u>, and Martin Pejcinovski, Jens Serneels, Yannick Jönsson, Isabelle Terrasson, and Naïma Kheddoumi for technical assistance.

Competing interests:

Mazzone has submitted a world-wide patent pending for diagnostic use of gene expression profiles in monocytes. All other authors declare no conflict of interest.

Funding/Support:

Hamm was funded by the Deutsche Forschungsgemeinschaft (DFG), Prenen by the Leuven University Hospitals Clinical Research Foundation, Rothé by Actions de Recherche Concertée (ARC). This work was supported by grants from the European Research Council (OxyMo to Mazzone), the Fournier-Majoie Foundation (FFMI), FWO (G.0.793.11.N.10), Belgian Foundation Against Cancer (2010-198) and Italian Association for Cancer Research (AIRC 12214).

Role of the funding sources:

The funders had no role in the design and conduct of the study; collection, management, analysis, and interpretation of the data; preparation, review, or approval of the manuscript; and decision to submit the manuscript for publication.

FIGURE LEGENDS

Figure 1: Flowchart of patient inclusion and sample analysis

Inclusion criteria for patients were sporadic histologically confirmed adenocarcinoma of the colon and/or rectum for cohort I-III and VI, patients in remission from CRC for a treatment-free interval of minimum 3 months for cohort V, histologically confirmed adenocarcinoma of the stomach or gastroesophageal junction or of the pancreas, or histologically confirmed gastritis for cohort IV.

Figure 2: Development and validation of a gene signature in circulating monocytes for diagnosis of CRC

a, **b**, Differentially expressed genes between all CRC patients (P,PM) and healthy volunteers (HV). The MA plot (**a**) shows the fold change versus the average expression intensity, while the Volcano plot (**b**) shows fold change in relation to the p values. Green, significantly downregulated genes; red, significantly upregulated genes; corrected p<0.05. **c**, Final gene signature for diagnosis of CRC, comprised of 23 genes, validated in a multicentric test set of patients. Expression levels are displayed as expression relative to the HV mean; boxes, first to third quartile; Whiskers, range; dots, values outside 1.5-times the interquartile distance; horizontal line, median; +, mean. *, p<0.05; **, p<0.01; ***, p<0.001. **d**, ROC analysis for P,PM versus HV in multicentric cohort II. **e**, ROC analysis for P,PM versus HV (negative to screening colonoscopy) in cohort III. See Supplementary Methods for classification approaches.

Figure 3: Tumour-released soluble factors induce the specific upregulation of the gene signature

a-d, Stimulating freshly isolated, naïve monocytes with medium containing soluble factors demonstrates that the genetic fingerprint in monocytes used for the diagnostic gene signature is specifically induced by the transformed colon epithelium (HCT) in comparison to a benign cell line (CCD), as demonstrated by expression analysis comparing selective marker genes in stimulated monocytes to mock control (**a**). Genetic alterations are independent of hypoxic cues (**b**). The gene signature is specific to CRC in comparison to monocytes stimulated by a gastric cancer cell line (MKN) (**c**). The gene signature is reverted after withdrawal of the stimulus *i.e.*, the conditioned medium (**d**). n=6 (biological replicates from 6 different healthy donors); bars, mean with SEM; *, p<0.05; **, p<0.01; ***, p<0.001; ****, p<0.0001; #, p<0.05 towards mock control, assessed by ANOVA with Bonferroni correction. All experiments were repeated at least twice.

Figure 4: The diagnostic gene signature is specific for CRC of all stages and reverts upon curative treatment

a, Expression of the gene signature in patients with cancer of the gastro-esophageal junction (PG), demonstrating no upregulation and thus specificity of the diagnostic signature for CRC. See Figure 2 for details on graphic elements. **b**, Gene signature in patients after curative treatment (patients in remission, PR), in which the expression levels revert to those of healthy volunteers in comparison to CRC patients. **c**, **d**, ROC analyses corresponding to Figure 4a. **e**, Four patients with isolated metastatic recurrence at the time of analysis (black dots) in a 2D-projection of the multi-gene expression levels. The gene signature of metachronously metastasized patients clusters with those patients with primary tumours (red), distinct from healthy individuals (blue). *, p<0.05; **, p<0.01; ***, p<0.001

Figure 5: Specificity of the gene signature to monocytes in comparison to **PBMCs**

a, Expression study assessing the gene signature in PBMCs in comparison to monocytes (PBMs). While the entire signature is confirmed in PBMs in this independent sample set, it is impossible to detect robust genetic alterations in s. S. pis in PBN. ng the previously "", p<0.001 PBMCs, demonstrating specificity to PBMs. See Figure 2 for details on graphic elements. **b**, Corresponding ROC analysis in PBMCs. **c**, ROC analysis of P,PM versus HV in monocytes, confirming the previously established classification performance. *, p<0.05; **, p<0.01; ***, p<0.001

Gut

REFERENCES

- 1. Siegel R, Naishadham D, Jemal A. Cancer statistics, 2012. CA: a cancer journal for clinicians 2012;**62**(1):10-29.
- 2. Weitz J, Koch M, Debus J, et al. Colorectal cancer. Lancet 2005;**365**(9454):153-65.
- Lieberman DA. Clinical practice. Screening for colorectal cancer. The New England journal of medicine 2009;361(12):1179-87.
- Stoop EM, de Haan MC, de Wijkerslooth TR, et al. Participation and yield of colonoscopy versus non-cathartic CT colonography in population-based screening for colorectal cancer: a randomised controlled trial. The lancet oncology 2012;13(1):55-64.
- Quintero E, Castells A, Bujanda L, et al. Colonoscopy versus fecal immunochemical testing in colorectal-cancer screening. The New England journal of medicine 2012;366(8):697-706.
- Pawa N, Arulampalam T, Norton JD. Screening for colorectal cancer: established and emerging modalities. Nature reviews Gastroenterology & hepatology 2011;8(12):711-22.
- Murdoch C, Muthana M, Coffelt SB, et al. The role of myeloid cells in the promotion of tumour angiogenesis. Nat Rev Cancer 2008;8(8):618-31.
- Shi C, Pamer EG. Monocyte recruitment during infection and inflammation. Nature reviews Immunology 2011;11(11):762-74.
- Sandel MH, Dadabayev AR, Menon AG, et al. Prognostic value of tumor-infiltrating dendritic cells in colorectal cancer: role of maturation status and intratumoral localization. Clin Cancer Res 2005;11(7):2576-82.

- 10. Sica A, Mantovani A. Macrophage plasticity and polarization: in vivo veritas. The Journal of clinical investigation 2012;**122**(3):787-95.
- 11. Wynn TA, Chawla A, Pollard JW. Macrophage biology in development, homeostasis and disease. Nature 2013;**496**(7446):445-55.
- Irvine KM, Gallego P, An X, et al. Peripheral blood monocyte gene expression profile clinically stratifies patients with recent-onset type 1 diabetes. Diabetes 2012;61(5):1281-90.
- Zawada AM, Rogacev KS, Schirmer SH, et al. Monocyte heterogeneity in human cardiovascular disease. Immunobiology 2012;217(12):1273-84.
- 14. Ma H, Hong M, Duan J, et al. Altered cytokine gene expression in peripheral blood monocytes across the menstrual cycle in primary dysmenorrhea: a case-control study. PloS one 2013;8(2):e55200.
- Honda T, Inagawa H, Yamamoto I. Differential expression of mRNA in human monocytes following interaction with human colon cancer cells. Anticancer research 2011;31(7):2493-7.
- Fletcher RH. Carcinoembryonic antigen. Annals of internal medicine 1986;104(1):66-73.
- 17. Schwarzenbach H, Hoon DS, Pantel K. Cell-free nucleic acids as biomarkers in cancer patients. Nature reviews Cancer 2011;**11**(6):426-37.
- Huang Z, Huang D, Ni S, et al. Plasma microRNAs are promising novel biomarkers for early detection of colorectal cancer. International journal of cancer Journal international du cancer 2010;**127**(1):118-26.
- 19. Church TR, Wandell M, Lofton-Day C, et al. Prospective evaluation of methylated SEPT9 in plasma for detection of asymptomatic colorectal cancer. Gut 2013.

- Gut
- Xu Y, Xu Q, Yang L, et al. Gene expression analysis of peripheral blood cells reveals toll-like receptor pathway deregulation in colorectal cancer. PloS one 2013;8(5):e62870.
- 21. Han M, Liew CT, Zhang HW, et al. Novel blood-based, five-gene biomarker set for the detection of colorectal cancer. Clinical cancer research : an official journal of the American Association for Cancer Research 2008;**14**(2):455-60.
- Marshall KW, Mohr S, Khettabi FE, et al. A blood-based biomarker panel for stratifying current risk for colorectal cancer. International journal of cancer Journal international du cancer 2010;**126**(5):1177-86.
- 23. Nichita C, Ciarloni L, Monnier-Benoit S, et al. A novel gene expression signature in peripheral blood mononuclear cells for early detection of colorectal cancer. Alimentary pharmacology & therapeutics 2014;**39**(5):507-17.
- Imperiale TF, Ransohoff DF, Itzkowitz SH, et al. Multitarget stool DNA testing for colorectal-cancer screening. The New England journal of medicine 2014;370(14):1287-97.
- Nyugen J, Agrawal S, Gollapudi S, et al. Impaired functions of peripheral blood monocyte subpopulations in aged humans. Journal of clinical immunology 2010;**30**(6):806-13.
- 26. Smyth GK. Linear models and empirical bayes methods for assessing differential expression in microarray experiments. Statistical applications in genetics and molecular biology 2004;**3**:Article3.
- 27. 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(1):289-300.

28. Sample size for microarray experiments. Secondary Sample size for microarray experiments. <u>http://bioinformatics.mdanderson.org/MicroarraySampleSize/</u>.

29. Burges CJC. A Tutorial on Support Vector Machines for Pattern Recognition. Data Min Knowl Discov 1998;**2**(2):121-67.

30. Breiman L. Random Forests. Mach Learn 2001;45(1):5-32.

- 31. Dietterich TG. Ensemble Methods in Machine Learning. Proceedings of the First International Workshop on Multiple Classifier Systems: Springer-Verlag, 2000:1-15.
- 32. Wessels LF, Reinders MJ, Hart AA, et al. A protocol for building and evaluating predictors of disease state based on microarray data. Bioinformatics 2005;21(19):3755-62.
- Edgar R, Domrachev M, Lash AE. Gene Expression Omnibus: NCBI gene expression and hybridization array data repository. Nucleic acids research 2002;30(1):207-10.
- 34. Piehler A, Grimholt R, Ovstebo R, et al. Gene expression results in lipopolysaccharide-stimulated monocytes depend significantly on the choice of reference genes. BMC Immunology 2010;**11**(1):21.
- 35. Guo C, Liu S, Wang J, et al. ACTB in cancer. Clinica chimica acta; international journal of clinical chemistry 2013;**417**:39-44.
- 36. Jones S, Chen WD, Parmigiani G, et al. Comparative lesion sequencing provides insights into tumor evolution. Proceedings of the National Academy of Sciences of the United States of America 2008;**105**(11):4283-8.
- 37. Goede V, Coutelle O, Shimabukuro-Vornhagen A, et al. Analysis of Tie2expressing monocytes (TEM) in patients with colorectal cancer. Cancer investigation 2012;**30**(3):225-30.
<text><text><text><text><text> 38. Schauer D, Starlinger P, Reiter C, et al. Intermediate monocytes but not TIE2-

Gut



193x153mm (300 x 300 DPI)



154x226mm (300 x 300 DPI)







169x164mm (300 x 300 DPI)



164x264mm (300 x 300 DPI)





164x179mm (300 x 300 DPI)

0,1

Gut

Supplementary Material

Tumour-Educated Circulating Monocytes are Powerful Candidate Biomarkers for Diagnosis and Disease Follow-up of Colorectal Cancer

Alexander Hamm, Hans Prenen, Wouter Van Delm, Mario Di Matteo, Mathias Wenes, Estelle Delamarre, Thomas Schmidt, Jürgen Weitz, Roberta Sarmiento, Angelo Dezi, Giampietro Gasparini, Françoise Rothé, Robin Schmitz, André D'Hoore, Hannes Iserentant, Alain Hendlisz & Massimiliano Mazzone

https://mc.manuscriptcentral.com/gut

CONTEN	ITS
Supplementary Methods	Page 3
Supplementary Notes	Page 14
Supplementary Figures	Page 17
Supplementary Tables	Page 28
Supplementary References	Page 35
https://mc.manuscripto	central.com/gut

Page 80 of 149

Gut

SUPPLEMENTARY METHODS

Patients

The composition of patient cohorts is given in detail in the main manuscript. Inclusion criteria for patients were sporadic histologically confirmed adenocarcinoma of the colon and/or rectum for cohort I-III and VI, patients in remission from CRC for a treatment-free interval of minimum 3 months for cohort V, histologically confirmed adenocarcinoma of the stomach or gastroesophageal junction or of the pancreas, or histologically confirmed gastritis for cohort IV. All patient samples were prospectively collected after histological diagnosis upon screening colonoscopy (reference standard defined by international clinical guidelines¹), prior to any treatment, at clinically indicated regular appointments separate of medical interventions (such as colonoscopy, surgical preparations etc.). All newly diagnosed patients presenting to the responsible clinicians were consecutively included when they met criteria and gave written informed consent. Healthy volunteers were included when there was no evidence or record of acute or chronic disease, with identical exclusion criteria as the patients. A subset of healthy individuals (within cohort III) was included upon screening colonoscopy without any pathological findings. Exclusion criteria were age of less than 40 years (to exclude cancers suspicious of genetic syndromes and restrict possible age-related variations in the monocyte phenotype reported previously²), history of oncological, chronic inflammatory, and autoimmune diseases within 10 years prior to this study, clinical or laboratory evidence of acute infection, anti-inflammatory and/or immunosuppressive medication within 90 days of blood sampling with the exception of occasional NSAID, commencement of medical or surgical anti-cancer treatment, medication with sedatives or opioid-based analgesics within 72 hours prior to blood sampling, clinical or microbiological evidence of altered

Gut

gut flora. Samples were excluded from further analysis when final histology of the surgical specimen did not confirm adenocarcinoma of the large intestine (assessed by board-certified pathologists within clinical routine procedures).

The following four oncological centres contributed samples to this study: Digestive Oncology, University Hospitals Leuven and Department of Oncology, KU Leuven, Leuven, Belgium; Department of General, Visceral, and Transplantation Surgery, University of Heidelberg, Heidelberg, Germany; Department of Oncology, San Filippo Neri, Rome, Italy; Medical Oncology Clinic, Institut Jules Bordet, Brussels, Belgium. The responsible scientists in each centre (1-2 per centre) were trained in the protocol for isolation of PBM to ensure uniformity of the procedure. All participants gave written informed consent, and the study was approved by the respective institutional review boards (Leuven: B322201215873, Brussels: CE1950, Heidelberg: 323/2004, Rome: 319/51). No adverse events from blood collection or colonoscopy were recorded in included participants.

Isolation of PBM

20ml of EDTA-anticoagulated peripheral venous blood was collected following clinical routine procedure, stored at 4°C and processed within 2 hours of blood collection. For further isolation, blood was diluted 1:2 with DPBS (free of Ca2+ and Mg2+) and layered carefully on Lymphoprep (Axis-Shield) in two separate tubes. All blood collection and isolation steps were performed identical for samples of all origin. Density gradient centrifugation was performed at 1,200g for 20 minutes at low acceleration and no brake. Samples with macroscopically visible hemolysis were excluded from further analysis. The PBMC interphase was collected carefully and washed twice for 12 minutes at 250g and 175g with PBS. Hemocytrometric analysis

Gut

was performed to ensure purity of PBMCs, and the pellet was pooled for further processing and washed once for 10 minutes at 300g. Cells were then incubated with CD14 magnetically-conjugated beads (BD) for 15 minutes at 4°C, washed 10 minutes at 300g and positively separated with the MACS system (Miltenvi) following the manufacturer's instructions. The CD14+ fraction was flushed out and washed once 10 minutes at 300g. Purity was assessed by FACS analysis for CD14 in the pilot phase and by hemocytometric analysis (CellDyn 3700, Abbott) in every further sample. Only samples with purity of >90% and viability >95% (assessed by Trypan Blue staining) were retained for further analysis. Cell pellets were lysed in Buffer RLT (Qiagen) at 10⁶ monocytes in 350µl of Buffer RLT and stored at -80°C. For each respective expression study, all samples were extracted simultaneously with the RNeasy Mini Kit (Qiagen) following the manufacturer's instructions. Quality control was performed by checking RNA quality on the Nanodrop system, and RNA integrity was checked for microarray samples on the Agilent Bio-Analyzer. Only samples with an extinction fraction 260/280 > 1.8 and 260/230 > 1.5, and an RNA integrity index of >6 were retained for further analysis.

Genome-wide expression analysis

For genome-wide expression analysis, RNA was amplified and biotinylated using Illumina TotalPrep RNA Amplification Kit (Ambion) following the manufacturer's instructions to obtain biotinylated cRNA, which was hybridized to Illumina HumanHT-12 v4 Expression BeadChips (Illumina) with the Illumina Whole-Genome Gene Expression Direct Hybridization Assay (Illumina) following the manufacturer's instructions. The Illumina HumanHT-12 v4 Expression BeadChip Kit contains 47,323 probes and 887 controls. After scanning, background-corrected expression values

and detection scores were extracted with GenomeStudio GX (version 1.5.4). For each array, we used the summarized expression level (AVG_Signal), standard error of the bead replicates (BEAD_STERR), number of beads used (AVG_NBEADS) and a detection score, which estimates the probability of a gene being detected above the background. Resulting expression data was analyzed with R, using the lumi package³. A variance stabilizing transformation⁴ was applied, followed by quantile normalization to compensate for batch effects of the individual bead chips. For each probe, the number of present calls over all samples was determined (the threshold on the detection was p<0.01), and probes absent in all samples were omitted in the analysis. This omitted subset consisted of 18,396 probes. Hence, analysis was performed for 28,927 probes. Differential expression was assessed with the limma package of \mathbb{R}^5 .

Quantitative RT-PCR (qPCR)

For qPCR analyses, 400ng of RNA was reverse transcribed with SuperScript III First Strand Kit (Invitrogen) following the manufacturer's instructions, and qPCR was performed in duplicates on a 7500Fast System (Applied Biosystems) using intronspanning PrimeTime qPCR Assays (Integrated DNA Technologies) listed in Supplementary Table 2. Wherever possible, qPCR assays were selected that covered the exon in which the Illumina Expression BeadChip probe was located. Raw data was analyzed with SDS v1.4 (Applied Biosystems), and expression was normalized within samples with the $\Delta\Delta$ CT method to reference gene *B2M*. Data was expressed relative to the average expression of that gene in the healthy volunteers in the dataset. Data points where duplicates differed by more than 1 CT were discarded. Inter-run validity was verified by both processing and running previously

Gut

analyzed samples as internal controls and ensuring correct clustering within their respective groups. Where necessary for normalization purposes, stored and validated healthy volunteer samples were re-profiled along with samples from cohorts IV and V.

Identification of a gene signature

For each pair-wise comparison between HV, P and PM, we evaluated all probes with a moderated t-test, as implemented in the limma-package⁵ of R. P-values were adjusted for multiple testing with Benjamini-Hochberg to control the false discovery rate⁶. A probe was selected as being differentially expressed between two groups when the adjusted p-value was smaller than 0.05 and the fold change exceeded 1.5 times up- or down-regulation ($\log_2 > 0.58$ or < -0.58, respectively). For the comparison between PM/P and HV, differential expression of the selected genes was further validated with qPCR in 8 randomly selected individuals from each of the groups in cohort I. The panel of 35 candidate genes derived from the 40 Illumina probes differentially expressed in cohort I was augmented by 8 genes which marginally missed the applied cutoff criteria and had been identified in unpublished in vitro and in vivo screens during the pilot phase. Minimal sample size for further cohorts was chosen to be 15 after conducting a statistical power analysis with the data from cohort I to estimate the expected variation in gene expression. Sample size was chosen to achieve a statistical power of 0.9 with an ordinary t-test when fold changes of 1.5 are considered and 5% false positives are accepted. Power calculations were done with the online tool from the Department of Bioinformatics and Computational Biology of MD Anderson Cancer Center⁷. Differential expression was considered to be confirmed by qPCR when the p-value after a two-tailed

unpaired t-test was smaller than 0.1 and/or the associated area under the ROC curve (AUC) was larger than 0.7. as calculated with Prism (GraphPad, Inc.). We chose deliberately for loose cut-offs on p-value and AUC for the confirmation, since less distinctly differentially expressed genes could in theory still add value to a (later developed) multiple-gene classification strategy.

Multicentric validation study

Overview. The diagnostic test consists of a gene panel assay in combination with software for decision support. The software implements an algorithm that takes the data from the assay as input and outputs a binary decision: whether the profiled sample comes from a CRC patient or not. The algorithm is an ensemble method (ENS)⁸ that consults 3 subroutines, then counts the number of votes in favor of CRC and finally proposes the decision that is supported by at least 2 subroutines. The 3 subroutines form a heterogeneous set of alternative classification algorithms: an easily interpretable ensemble stump classifier (SGMV – single gene majority vote), a linear support vector machine (SVM) and a more complex random forest (RF). The parameters of the 3 subroutines were fitted in parallel to a subset of samples from the multi-centric cohort II. This training subset was constructed via stratified random sampling. Performance of the algorithm was assessed through a Monte Carlo cross-validation (MCCV) procedure on the training data and further validated on the samples from cohort II that were excluded during training.

Stratified random sampling. We identified combinations of the four oncology centres and two sample classes (i.e. HV or CRC) as 8 strata. From each stratum, we sampled 2 times as much training samples as validation samples. The actual number of samples per stratum was chosen so that i. there was no evidence of dependence

Gut

of class labeling on centre in either validation or training dataset, ii. the final datasets were balanced (i.e. as much HV as CRC). Dependence between class labeling and centre of origin was excluded by testing with a Fisher's exact test (p > 0.93). The random split was performed prior to fitting parameters and retained for all further analyses to obtain realistic measures of classification performance. Since our subroutines required complete data, we imputed missing values after assembling the training and validation datasets for each dataset separately using nearest neighbor averaging, as implemented in the impute-package in \mathbb{R}^9 .

Subroutines. The SGMV compares the expression value of each input gene first to a gene-specific cut-off and then assigns a defined class to an unknown sample depending on whether the cut-offs are exceeded for at least half of the genes (i.e. majority vote). The SGMV parameters hence consist of gene-specific cut-offs. The gene-specific cut-offs are fitted by taking that value that corresponds to the point closest to the top-left corner of the gene-associated ROC curve, using the pROCpackage in R¹⁰. The SVM with linear kernel is similar to linear discriminant analysis, taking as input the expression values of a set of genes and comparing a linear combination of the input values to a threshold in order to assign a defined class to an unknown sample, thereby giving higher weight to more informative genes. The SVM parameters hence consist of gene-specific weights and one threshold. We fitted the parameters with the kernlab-package in R¹¹. The RF pushes the expression values of a set of genes through a multitude of decision trees (each looking at a random subset of genes and built from a random subset of samples from the training data), notes down for each class the proportion of supporting individual trees and finally assigns the class with highest support. The RF parameters hence consist of individual decision trees. We fitted the parameters with the randomForest-package in R¹².

Avoiding over-fitting. Fitting the parameters of the SVM and RF subroutines was conditioned on hyper-parameters that influence the flexibility of the subroutines to fit the training data. Too flexible procedures lead to over-fitting of training samples at the cost of bad performance on unseen samples. Flexibility was therefore constrained by selecting hyper-parameters from a range of options with Monte Carlo cross-validation (MCCV), prior to final determination of the common parameters. We divided the training dataset during 100 cycles in 2/3 and 1/3, trained the SVM/RF each time on the largest part with a given hyper-parameter, tested the SVM/RF each time on the smallest part and finally averaged the AUC and BER of all cycles for a particular hyper-parameter value. We chose the hyper-parameter with best average AUC, or in case of multiple options, the one with best average BER. Note that this MCCV procedure to select hyper-parameters was also run as an inner loop within the outer MCCV loop when algorithm performance was assessed (see above)¹³.

Performance metrics. The classifiers were validated on the qPCR test dataset, constructed from healthy volunteers and patients of multi-centric cohort II who were not included during development of the models (see above). To verify the similarity of the test set to the training set, a Spearman-correlation between all assays was performed, ensuring that test assays did not cluster separately from training assays. A separate clustering would have been an indication that the training dataset was not representative for the test samples. Two types of performance were finally reported: ranking performance and classification performance. Ranking performance is the capability of an algorithm to give a higher score to an individual from class CRC than to an individual from class HV. We measured ranking performance by the area under the ROC curve (AUC). For all 4 routines (SGMV, SVM, RF and ENS), we provided the AUC as well as the lower bound and upper bound of its 95% confidence interval,

Gut

as computed after 2,000 bootstraps with the pROC-package in R¹⁰. Classification performance measures the capability of an algorithm to assign an individual to the correct class. We reported for all routines the balanced error rate (BER), sensitivity (Se) and specificity (Sp). For Se and Sp, we also computed the lower bound and upper bound of the 95% confidence interval after 2,000 bootstraps.

Complementary data analysis

A complementary data analysis by an independent team (DNAlytics, Belgium) on the same 23-marker signature led to the same conclusions in terms of performances. Another (per-marker) normalization procedure has been proposed. This normalization is applied on the log-transformed gene expression (i.e. Δ CT values) and consists in computing, on the training set (for example Cohort II, both HV and CRC), the mean and standard deviation of each marker. When a prediction has to be made on a new, potentially isolated sample, each marker measurement of this new sample is normalized by subtracting the corresponding mean, and by dividing by the corresponding standard deviation. A modified procedure has also been proposed for the imputation of missing values, making it dependent on the reference cohort only. This avoids the need for a new reference HV batch as prediction has to be made on a new (set of) sample(s).

The first experiment consisted in cross-validating a model on Cohort II (BER: 8.4% [3.4%;13.4%]; AUC: 0.93 [0.88;0.98]). A second experiment consisted in learning the same type of model on Cohort II and having it make predictions on Cohort III (BER: 13.2%; AUC: 0.92). All analyses were performed in R with scripts designed by DNAlytics, fully independent from other analyses described in this paper.

In vitro model system

To study the effects of tumour-released soluble factors on gene expression in monocytes, we established an in vitro model system. Medium conditioned with cell-released soluble factors was obtained by seeding the following cell lines at 40% confluence at 37°C at 21% O₂, 5% CO₂ in a moist atmosphere in their respective medium and ultra-filtering the conditioned medium 72 hours later: HCT116 (new from ATCC, CCL-247) in RPMI (10% FBS, 1% Glutamine, 1% PenStrep), grown in normoxia or hypoxia (1% O2), CCD 841 CoN (new from ATCC CRL-1790) in EMEM (10% FBS, 1% Glutamine, 1% PenStrep), MKN-45 (a kind gift from Frans van Roy, UGent, Belgium) in RPMI (10% FBS, 1% Glutamine, 1% PenStrep, 1% Na-Pyruvate). Each medium was also incubated separately without cells to obtain the respective mock controls. Absence of Mycoplasma species was verified with MycoAlert Mycoplasma Detection Kit (Lonza).

Monocytes from healthy volunteers (n=6) were isolated as described above and were seeded at 200,000 cells / well in a tissue-culture treated 24-well plate (Costar) in IMDM (10% autologous serum, 1% Glutamine), supplemented 1:5 with conditioned medium. Cells were lysed in Buffer RLT (Qiagen) after 18 hours. For experiments on reversion of the gene signature after withdrawing the stimulus, monocytes were washed with PBS after 18 hours of culture in conditioned medium, and medium was refreshed with plain IMDM (10% autologous serum, 1% Glutamine). After 72 hours, cells were then lysed in Buffer RLT. All experiments were performed in technical quadruplicates and repeated at least twice.

All RNA was extracted simultaneously with the RNeasy MicroKit (Qiagen) following the manufacturer's instructions, and RNA quality was verified with the Nanodrop system as described above.

https://mc.manuscriptcentral.com/gut

Gut

<page-header><text><text><text>

 Gut

SUPPLEMENTARY NOTES

Supplementary Note 1

To select a robust reference gene, we checked in the available microarray data for stably expressed genes that met all of the following criteria: *i*. p>0.5 for any pair-wise comparison of groups, *ii*. lowest coefficient of variation among all samples, *iii*. good annotation of the gene, *iv*. consistent high expression levels. After further screening of available literature on potential reference genes ("housekeeping genes"), we selected in a pilot phase the following genes from the stably expressed genes for analysis: *ACTB*, *B2M*, *HPRT*, *PGK1*, *RPS14*, and *RPS27*. We found most stable expression for *B2M*, which in addition showed a lower coefficient of variation than *ACTB*, recently suggested to be a less-than-ideal housekeeping gene depending on the cellular context¹⁴ ¹⁵. To rule out any inconsistency in the use of the reference gene, we opted to use *B2M* and compared the qPCR expression data of cohort II to normalization against *ACTB*, which yielded similar results (Supplementary Figure 3a and data not shown).

Supplementary Note 2

We assessed the annotated biological function of the 23 genes comprising the final diagnostic signature, as well as their putative role in monocyte function and/or phenotype. An overview can be found in Supplementary Table 5. A pathway analysis by Ingenuity Pathway Analysis (www.ingenuity.com) revealed that top pathways and functions included acute phase response signalling, free radical scavenging, immune cell trafficking, inflammatory disease, and cell death and survival. Taking those 7 genes upregulated in the in vitro model system, their annotated function suggests that immune signals may be the underlying mechanism in driving their expression

Gut

shift. However, we could not identify key regulators of known pathways, probably due to the limited information on reciprocal effects of PBM and tumour cells¹⁶. Though of high interest with regards to the biological function, functional biological knowledge is dispensable to exploit the full potential of the gene signature as a diagnostic tool in analogy to other important clinical tests, which are devoid of a biological understanding (e.g., prostate specific antigen, PSA, and pro-calcitonin, PCT).

Supplementary Note 3

In accordance with our initial screening results, we found no differences in expression patterns of P versus PM (data not shown). Moreover, as cumulating evidence is suggesting subcategories of CRC according to its location¹⁷, we investigated if the gene signature was capable of separating left versus right CRC or colon versus rectal cancer, respectively. In line with the homogeneous clustering of samples, we found no differences by location (AUC of 0.45 [0.20-0.73] for left versus right CRC and AUC of 0.47 [0.28-0.70] for colon versus rectal cancer).

Supplementary Note 4

We sought to confirm our findings from the screening in independent samples by independent techniques to rule out bias by the chosen technique and maximize chances of extrapolation to other clinical centres. Our first step was a random reprocessing of collected samples and assessment by qPCR, which led to an initial refinement of the gene signature, while some genes in this subset of samples performed well even as single markers. By assessing Spearman correlation values between expression data in the Illumina platform (used for screening) and the qPCR technique (used for confirmation), we could rule out discrepancies in expression

<text>



Supplementary Figure 1: Isolation of PBM and monocyte counts

a, Quality control of PBM isolation procedure in the pilot phase: FACS staining as histogram for CD14 (FITC). Comparison of the CD14⁺ flow-through (left) and the CD14⁺ purified monocytes (right). b, Representative hemocytometric assessment of PBM purity, which was performed for each individual sample. c, Monocyte counts in whole blood were not different between (P,PM) and HV, neither relative (left), nor absolute (right).

Gut

а

μ

p value (-log₁₀) 10

-2.0

-1.0

Fold change (log₂)

b

Fold change (log₂)

P vs HV

۰.

Average intensity (log₂)

Differential gene expression

P vs HV

2.0

1.0



PM vs HV





P vs PM

Fold change (log₂)

Supplementary Figure 2: Differentially expressed genes in PBM

1.0

a

0.5 0.5

₽

value (-bg₁₀)

a, b, Differentially expressed genes in groupwise comparison of P, PM, and HV. The MA plots (a) show the fold change versus the average expression intensity, while the Volcano plots (b) show fold change in relation to the p values. Green, significantly downregulated genes; red, significantly upregulated genes; corrected p<0.05.

Fold change (\log_2)

-0.5

r he i n rela.



Supplementary Figure 3: Technical validation (subset of cohort I)

Ē 🗖

÷

a, Comparative dot plot of raw CT values in qPCR for *ACTB* and *B2M*, revealing that the distribution is similar for both genes, and box-and-whiskers plot comparing normalization against both reference genes. **b**, Expression levels of all 43 putative candidates identified by genome-wide screening and assessed by qPCR. Expression levels are displayed as expression relative to the HV mean; boxes, first to third quartile; Whiskers, range; dots, values outside 1.5-times the interquartile distance; horizontal line, median; +, mean; *, p<0.1; **, p<0.01; ***, p<0.01.

CD68

40-

20- _____

0-

8-

0-

2-

٥Ţ

₽÷

🖳 😐

TAF15

RLPL2

i l

IFR2

ē 🖻

GPER



Supplementary Figure 4: Gene expression levels of non-confirmed candidates in the multicentric validation (cohort II)

Expression levels are displayed as expression relative to the HV mean; boxes, first to third quartile; Whiskers, range; dots, values outside 1.5-times the interquartile distance; horizontal line, median; +, mean; *, p<0.05; **, p<0.01; ***, p<0.001.

C.

https://mc.manuscriptcentral.com/gut

Gut



Supplementary Figure 5: The gene signature stays robust over disease progression (cohort II)

Multicentric validation of the finding that the gene signature cannot discriminate between P and PM. Expression levels are displayed as expression relative to the HV mean; boxes, first to third quartile; Whiskers, range; dots, values outside 1.5-times the interquartile distance; horizontal line, median; +, mean; *, p<0.05; **, p<0.01; ***, p<0.001.

iminate . .,rtile; Whiske, .,<0.01; ***, p<0.0c

Page 99 of 149

1 2 3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38 39

40

41

42

43

44

45

46

47

48

49

50



AUC = 0.7885

AUC = 0.6619

0.4 0.6 0.8 1.0

1 - Specificity

0.6 0.8 1.0

0.4

1 - Specificity

SOCS3

0.2

AUC = 0.8101 0.01 0.0 1.0

BAX

0.2

AUC = 0.8356

AUC = 0.6658

AUC = 0.8305

AUC = 0.7314

0.4 0.6 0.8 1.0

1 - Specificity

S100P

1.0

0.4 0.6 0.8

1 - Specificity

HP

0.2

0.4 0.6 0.8

1 - Specificity

HBA1

1.0

1.0

0.4 0.6 0.8

1 - Specificity

DDIT4

0.6 0.8 0.4 0.2 1 - Specificity

AUC = 0.7046

AUC = 0.7597

0.4 0.6 0.8 1.0

1 - Specificity

1.0

0.4 0.6 0.8

1 - Specificity

ТКТ

0.0

1.0

0.8

₹ 0.€ Sensil 0.4

0.2

0.0

0.2

Supplementary Figure 6: Single gene ROC analysis ROC analyses for each individual in cohort II. AUC, area under the curve.

0.0

1.0-

0.8

≹ 0.6-

ын С 0.4-

0.2

0.0

0.2

AUC = 0.6187

AUC = 0.7104

0.4 0.6 0.8 1.0

1 - Specificity

1.0

0.4 0.6 0.8

1 - Specificity

SLPI

0.0

1.01

0.8

₹ 0.6

349S 0.4

0.2

0.0

0.2

0.2

Gut



Supplementary Figure 7: Identification of putative markers in the in vitro model

Shown are the expression levels of the 16 genes not selected out of the gene signature, which show alterated expression levels in culture without any stimulus. Expression levels are shown as mean with SEM at 18 hours and 72 hours later (90 hours). *, p<0.05; **, p<0.01; ****, p<0.001; ****, p<0.0001; n.e., not expressed *in vitro*. Vem goodf; n.e., ... Page 101 of 149



Gut

Scatter plots of Cohort I displaying correlation between Illumina microarray (x axis) and qPCR data (y axis). Spearman correlation values and p values are noted in the figures.





Page 103 of 149







Supplementary Figure 8 – continued

6

SUPPLEMENTARY TABLES

Supplementary Table 1: IDT PrimeTime qPCR Assays

Gene Name	Assay ID
ACP5	Hs.PT.47.311649.g
ACTB	Hs.PT.47.227970.g
ADM	Hs.PT.47.59577.a
ALDH1A1	Hs.PT.47.4497955
APP	Hs.PT.47.3063778
ARPC1B	Hs.PT.47.18828860
B2M	Hs.PT.47.18818394
BAX	Hs.PT.47.18828862
CCR1	Hs.PT.47.18828864
CD68	Hs.PT.47.18828865
CTSZ	Hs.PT.47.18828866
CXCR4	Hs.PT.47.512220
DDIT4	Hs.PT.47.18828867
DNAJC7	Hs.PT.47.18828868
ENSA	Hs.PT.47.18828869
FCER1A	Hs.PT.47.18828870
FKBP5	Hs.PT.47.18828871
GPER	Hs.PT.47.18828872
HBA1 / HBA2	Hs.PT.47.18828873
HBB	Hs.PT.47.18828874
HLA-DQA1	Hs.PT.47.18828891
HLA-DRB4	Hs.PT.47.18828875
HMOX1	Hs.PT.47.18828876
HNRNPK	Hs.PT.47.18828877
HP	Hs.PT.47.18828878
HPRT1	Hs.PT.47.1231226
IER2	Hs.PT.47.18828880
IL1R2	Hs.PT.47.18828881
LAPTM4A	Hs.PT.47.18828882
LOC100008589	Hs.PT.47.18828883
LOC100130707	Hs.PT.47.18828884
LOC100132394	Hs.PT.47.18828885
LOC100170939	Hs.PT.47.18828886
LOC644063	Hs.PT.47.18828888
LOC653888	Hs.PT.47.18828889
LOC723972	Hs.PT.47.18828890
PGK1	Hs.PT.47.18828893
RILPL2	Hs.PT.47.18828894
RPS14	Hs.PT.47.18828895
RPS27	Hs.PT.47.18828896
S100P	Hs.PT.47.18828897

https://mc.manuscriptcentral.com/gut

29

2	
3	
4	
5	
6	
7	
0	
0	
9	
10	
11	
12	
12	
13	
14	
15	
16	
17	
10	
IÖ	
19	
20	
21	
22	
22	
23	
24	
25	
26	
27	
21	
28	
29	
30	
31	
22	
32	
33	
34	
35	
36	
27	
31	
38	
39	
40	
<u>⊿1</u>	
40	
42	
43	
44	
45	
16	
40	
41	
48	
49	
50	
51	
51	
52	
53	
54	
55	
56	
50	
57	
58	

59 60

SEPT5 Hs.PT.47.2501884 SLC39A1 Hs.PT.47.18828898 SLPI Hs.PT.47.18828900 TAF15 Hs.PT.47.18828901 TKT Hs.PT.47.18828902 TNF Hs.PT.47.18828903	Gene Name	Assay ID	
SLC39A1 Hs.PT.47.18828898 SLPI Hs.PT.47.18828900 TAF15 Hs.PT.47.18828901 TKT Hs.PT.47.18828902 TNF Hs.PT.47.14765639.g TNPO1 Hs.PT.47.18828903	SEPT5	Hs.PT.47.2501884	
SLPI Hs.PT.47.18828899 SOCS3 Hs.PT.47.18828900 TAF15 Hs.PT.47.18828902 TKT Hs.PT.47.18828902 TNF Hs.PT.47.14765639.g TNPO1 Hs.PT.47.18828903	SLC39A1	Hs.PT.47.18828898	
SOCS3 Hs.PT.47.18828900 TAF15 Hs.PT.47.18828901 TKT Hs.PT.47.18828902 TNF Hs.PT.47.14765639.g TNPO1 Hs.PT.47.18828903	SLPI	Hs.PT.47.18828899	
TAF15 Hs.PT.47.18828901 TKT Hs.PT.47.18828902 TNF Hs.PT.47.14765639.g TNPO1 Hs.PT.47.18828903	SOCS3	Hs.PT.47.18828900	
TKT Hs.PT.47.18828902 TNF Hs.PT.47.14765639.g TNPO1 Hs.PT.47.18828903	TAF15	Hs.PT.47.18828901	
TNF Hs.PT.47.14765639.g TNPO1 Hs.PT.47.18828903	ТКТ	Hs.PT.47.18828902	
<i>TNPO1</i> Hs.PT.47.18828903	TNE	Hs.PT.47.14765639.g	
		Hs PT 47 18828903	
	INPU1	HS.P1.47.18828903	

-	q
ACP5	0.6760
ADM	1
ALDH1A1	0.2020
APP	1
ARPC1B	1
BAX	0.7569
CCR1	1
CD68	1
CTSZ	1
CXCR4	0.2020
DDIT4	
DNAJC7	1
ENSA	0.3600
FCER1A	1
FKBP5	0.3600
GPER	0.2401
HBA1	0.2411
HBB	1
HLA-DQ1	0.1095
HLA-DRB4	1
HMOX1	1
HNRNPK	0.6160
HP	0.6160
IER2	0.8236
IL1R2	0.1773
LAPTM4A	0.2651
LOC100008589	1
LOC100170939	1
LOC643888	1
LOC644063	0.0147
LOC723972	0.0552
RLPL2	0.4941
RN28S1	1
S100P	1
SDHC	0.6160
SEPT5	1
SLC39A1	1
SLPI	0.8103
SOCS3	1
TAF15	0.3600
ткт	0.1162
TNF	0.5485
TNPO1	0.6160
Gut

Supplementary Table 3: Overview of development of a validated gene signature from putative candidates

			Genomewi	de Screen	ning				Confirm	ation and Valio	lation
	P,PM vs HV ^a		Ρv	/s HV			PM vs HV			P,PM vs. HV	
	Ratio	р		Ratio	р		Ratio	р		Ratio	р
ADM	2,00	<0,0001	ADM	1,75	0,0059	ADM	2,27	<0,0001	ACP5 [⊳]	1,61	<0,0001
ALDH1A1	0,66	0,0002	CTSZ	1,76	0,0103	ALDH1A1	0,56	<0,0001	ADM	2,16	<0,0001
ARPC1B	1,55	0,0209	DDIT4	1,78	0,0226	AQP9	1,62	<0,0001	ALDH1A1	0,88	<0,0001
BAX	1,50	0,0001	DNAJC7	1,59	0,0005	BAX	1,52	0,0008	APP	1,61	<0,0001
CTSZ	1,79	0,0007	FCER1A	0,61	0,0296	CTSZ	1,81	0,0056	ARPC1B	0,98	0,6497
DDIT4	1,71	0,0063	HBA1	3,51	0,0008	DDIT4	1,65	0,0477	BAX	1,76	<0,0001
DNAJC7	1,59	<0,0001	HBA2	4,31	0,0004	DNAJC7	1,60	0,0004	CCR1	0,90	0,3981
FCER1A	0,52	0,0002	HBB	3,95	0,0004	DYSF	1,52	0,0002	CD68	1,76	<0,0001
FKBP5	1,61	0,0001	HMOX1	1,55	0,0017	FCER1A	0,45	0,0001	CTSZ	1,96	<0,0001
GPER	1,58	0,0006	HNRNPK	1,60	0,0497	FCGR1A	1,52	0,0003	CXCR4	2,24	<0,0001
HBA1	2,33	0,0078	HS.143909	1,56	<0,0001	FKBP5	1,85	<0,0001	DDIT4	1,47	0,0025
HBA2	2,69	0,0051	HS.581828	1,52	<0,0001	GPER	1,78	0,0001	DNAJC7	1,07	0,1045
HBB	2,39	0,0099	HS.61208	1,65	<0,0001	HLA-DRB6	0,42	0,0102	ENSA	0,89	0,1122
HMOX1	1,54	0,0001	IER3	1,50	0,0009	HMOX1	1,53	0,0020	FCER1A	0,97	0,7541
HNRNPK	1,58	0,0125	LOC100008589	1,68	0,0131	HP	1,75	0,0080	FKBP5	2,45	<0,0001
HP	1,54	0,0131	LOC100128274	0,66	0,0195	HS.61208	1,56	<0,0001	GPER	5,29	<0,0001
HS.143909	1,51	<0,0001	LOC100130707	1,51	0,0232	LOC100170	939 1,65	0,0001	HBA1	15,07	0,0165
HS.61208	1,60	<0,0001	LOC100132394	1,79	0,0095	LOC100190	986 1,53	0,0001	HBB	11,96	0,0281
IL1R2	1,50	0,0482	LOC100132727	0,66	0,0282	LOC153561	1,73	0,0001	HLA-DQA1	1,01	0,8425
LOC100008	3589 1,55	0,0079	LOC100134364	1,57	0,0057	LOC441087	1,54	0,0177	HLA-DRB4	0,77	0,3931
LOC100129	9685 1,71	0,0356	LOC153561	1,50	0,0049	RNF146	1,50	0,0002	НМОХ1	0,95	0,7338
LOC100132	2394 1,65	0,0045	LOC649143	1,90	0,0133	S100P	1,75	0,0007	HNRNPK	0,92	0,2280
LOC100134	4364 1,53	0,0009	LOC723972	1,51	0,0001	SEPT5	1,59	0,0347	HP	1,92	<0,0001
LOC100170	939 1,54	<0,0001	LOC728755	0,64	0,0210	SLC39A1	1,54	0,0025	IER2	0,97	0,8782
LOC153561	1,61	<0,0001	SLC39A1	1,50	0,0058	SOCS3	1,73	0,0014	IL1R2	0,86	0,4209

			Genom	newide Screen	ing				Confirmation	and Valic	lation
P,PN	// vs H V ^a			P vs HV			PM vs HV		P,PM	vs. HV	
	Ratio	р		Ratio	р		Ratio	р		Ratio	р
LOC649143	1,56	0,0356	TAF15	2,06	0,0001	TAF15	1,73	0,0002	LAPTM4A	1,59	<0,000
LOC653156	1,73	0,0443	ΤΚΤ	1,58	0,0034	ΤΚΤ	1,55	0,0048	LOC100008589	0,99	0,931
LOC653737	1,86	0,0472	ZNF223	0,66	0,0478	TNPO1	1,55	0,0001	LOC100170939	1,09	0,000
LOC728755	0,66	0,0066				UPP1	1,58	<0,0001	LOC643888	1,05	0,261
S100P	1,53	0,0020	S'A			ZBTB16	1,52	0,0252	LOC644063	1,54	<0,000
SEPT5	1,57	0,0094							LOC723972	1,03	0,190
SLC39A1	1,52	0,0003							RLPL2	0,95	0,361
SOCS3	1,51	0,0043							RN28S1	1,03	0,300
TAF15	1,76	<0,0001							S100P	3,35	<0,000
ТКТ	1,56	0,0003							SDHC	1,04	0,301
									SEPT5	3,47	0,002
									SLC39A1	1,02	0,382
									SLPI	15,76	0,009
									SOCS3	1,60	0,015
									TAF15	0,84	0,015
									ткт	1,79	<0,000
									TNF	0,75	0,020
									TNPO1	1.02	0,316

^a Listed are the gene symbols to which probes correspond. Note that the identified 40 probes correspond to 35 genes, as several probes may exist for one gene. See Supplementary methods for details on gene numbers. ^b Genes confirmed by qPCR are shown in bold print (23 genes).

expression* p AUC ACP5 81,336 1,73 0.0081c 0.79 ALDM 73,107 4,23 0.0941 0.95 ALDH1A1 47,649 0,99 0.9624 0.51 APP 176,332 1,32 0.1576 0.73 ARPC1B 1873,873 1,15 0.3336 0.57 BAX 11,474 1,29 0.0978 0.67 CCR1 338,797 1,01 0.9292 0.52 CD68 1821,912 1,27 0.0640 0.76 CTSZ 1580,418 1,28 0.0637 0.76 CXCR4 508,754 1,52 0.0065 0.84 DD174 38,963 4,34 0.0010 0.96 NAJC7 105,494 0.92 0.531 0.62 ENSA 250,287 1,21 0.2580 0.67 FCER1A 410,118 0.54 0.0768 0.73 FKBP5 39,131 2.08 </th <th></th> <th>Mean</th> <th>Fold ratio^b</th> <th></th> <th>_</th>		Mean	Fold ratio ^b		_
ACP5 81,336 1,73 0.0084° 0.79 ADM 73,107 4,23 0.0941 0.95 ALDH1A1 47,649 0,99 0.9624 0.51 APP 176,332 1,32 0.1576 0.73 ARPC1B 1873,873 1,15 0.3336 0.57 BAX 11,474 1,29 0.0978 0.67 CCR1 338,797 1,01 0.9292 0.52 CD68 1821,912 1,27 0.0640 0.76 CTSZ 1580,418 1,28 0.0637 0.76 CXCR4 508,754 1,52 0.0065 0.84 DDI74 36,963 4,34 0.0010 0.96 DNA/C7 105,494 0.92 0.5431 0.62 ENSA 250,287 1,21 0.2580 0.67 FCER1A 410,118 0.54 0.0768 0.73 FKBP5 39,131 2.08 0.0013 0.89 GPER 1,875 7.59 0.0138 0.93 HBA1 5243		expression ^a		р	AUC
ADM 73,107 4,23 0.0941 0.95 ALDH1A1 47,649 0,99 0.9624 0.51 APP 176,332 1,32 0.1576 0.73 ARPC1B 1873,873 1,15 0.3336 0.57 BAX 11,474 1,29 0.0978 0.67 CCR1 338,797 1,01 0.9292 0.52 CD68 1821,912 1,27 0.0640 0.76 CXCR4 508,754 1,52 0.0065 0.84 DDI74 36,963 4,34 0.0010 0.96 DNAJC7 105,494 0,92 0.5431 0.62 ENSA 250,287 1,21 0.2580 0.67 FKBP5 39,131 2.08 0.0013 0.89 GPER 1,875 7.59 0.0138 0.93 HBA 40,188 31,10 0.0773 0.85 HLA-DQ1 1748,345 0,53 0.0918 0.77 HADN1 5243,339 41,40 0.0861 0.86 HBB 440,	ACP5	81,336	1,73	0.0081 [°]	0.79
ALDH1A1 47,649 0,99 0.9624 0.51 APP 176,332 1,32 0.1576 0.73 ARPC1B 1873,873 1,15 0.3336 0.57 BAX 11,474 1,29 0.0978 0.67 CCR1 338,797 1,01 0.9292 0.52 CD68 1821,912 1,27 0.0640 0.76 CTSZ 1580,418 1,28 0.0637 0.76 CXCR4 508,754 1,52 0.0065 0.84 DDI74 36,963 4,34 0.010 0.96 DNAJC7 105,494 0.92 0.5431 0.62 ENSA 250,287 1,21 0.2580 0.67 FCER1A 410,118 0.54 0.0768 0.73 FKBP5 39,131 2,08 0.0013 0.89 GPER 1,875 7,59 0.0138 0.93 HBA1 5243,339 41,40 0.0661 0.53 HBA1 129,478 1,50 0.0218 0.77 HLA-DR4	ADM	73,107	4,23	0.0941	0.95
APP 176,332 1,32 0.1576 0.73 ARPC1B 1873,873 1,15 0.3336 0.57 BAX 11,474 1,29 0.0978 0.67 CCR1 338,797 1,01 0.9292 0.52 CD88 1821,912 1,27 0.0640 0.76 CTSZ 1580,418 1,28 0.0637 0.76 CXCR4 508,754 1,52 0.0665 0.84 DDIT4 36,963 4,34 0.0010 0.96 DNAJC7 105,494 0,92 0.5431 0.62 ENSA 250,287 1,21 0.2580 0.67 FCER1A 410,118 0.54 0.0768 0.73 FKBP5 39,131 2,08 0.0013 0.89 GPER 1,875 7,59 0.0138 0.93 HBA1 5243,339 41,40 0.0861 0.86 HACD1 1748,345 0.53 0.0918 0.77 HLA-	ALDH1A1	47,649	0,99	0.9624	0.51
ARPC1B 1873,873 1,15 0.3336 0.57 BAX 11,474 1,29 0.0978 0.67 CCR1 338,797 1,01 0.9292 0.52 CD68 1821,912 1,27 0.0640 0.76 CTSZ 1580,418 1,28 0.0637 0.76 CXCR4 508,754 1,52 0.0065 0.84 DDI74 36,963 4,34 0.0010 0.96 DNAJC7 105,494 0,92 0.5431 0.62 ENSA 250,287 1,21 0.2580 0.67 FCER1A 410,118 0.54 0.0768 0.73 FKBP5 39,131 2,08 0.0013 0.89 GPER 1,875 7,59 0.0138 0.93 HBA1 5243,339 41,40 0.0861 0.86 HBB 440,188 31,10 0.0773 0.85 HLA-DQ1 1748,345 0,53 0.0918 0.77 H	APP	176,332	1,32	0.1576	0.73
BAX 11,474 1,29 0.0978 0.67 CCR1 338,797 1,01 0.9292 0.52 CD68 1821,912 1,27 0.0640 0.76 CTSZ 1580,418 1,28 0.0637 0.76 CXCR4 508,754 1,52 0.0065 0.84 DDIT4 36,963 4,34 0.0010 0.96 DNAJC7 105,494 0.92 0.5431 0.62 ENSA 250,287 1,21 0.2580 0.67 FCER1A 410,118 0.54 0.0768 0.73 FKBP5 39,131 2,08 0.0013 0.89 GPER 1,875 7,59 0.0138 0.93 HBA1 5243,339 41,40 0.0861 0.86 HBB 440,188 31,10 0.0773 0.85 HLA-DQ1 1748,345 0,53 0.0918 0.77 HNNNFK 1648,472 1,05 0.7356 0.52 HP 129,478 1,50 0.0218 0.76 LAPTM4A <t< td=""><td>ARPC1B</td><td>1873,873</td><td>1,15</td><td>0.3336</td><td>0.57</td></t<>	ARPC1B	1873,873	1,15	0.3336	0.57
CCR1 338,797 1,01 0.9292 0.52 CD68 1821,912 1,27 0.0640 0.76 CTSZ 1580,418 1,28 0.0637 0.76 CXCR4 508,754 1,52 0.0065 0.84 DDIT4 36,963 4,34 0.0010 0.96 DNAJC7 105,494 0.92 0.5431 0.62 ENSA 250,287 1,21 0.2580 0.67 FCER1A 410,118 0.54 0.0768 0.73 FKBP5 39,131 2.08 0.0013 0.89 GPER 1,875 7,59 0.0138 0.93 HBA1 5243,339 41,40 0.0861 0.86 HBB 440,188 31,10 0.0773 0.85 HLA-DQ1 1748,345 0,53 0.0918 0.77 HNAX1 405,989 1,30 0.0729 0.70 HNRNPK 1648,472 1,05 0.256 0.63 <td< td=""><td>BAX</td><td>11,474</td><td>1,29</td><td>0.0978</td><td>0.67</td></td<>	BAX	11,474	1,29	0.0978	0.67
CD68 1821,912 1,27 0.0640 0.76 CTSZ 1580,418 1,28 0.0637 0.76 CXCR4 508,754 1,52 0.0065 0.84 DDIT4 36,963 4,34 0.0010 0.96 DNAJC7 105,494 0.92 0.5431 0.62 ENSA 250,287 1.21 0.2580 0.67 FCER1A 410,118 0.54 0.0768 0.73 FKBP5 39,131 2.08 0.0013 0.89 GPER 1.875 7,59 0.0138 0.93 HBA1 5243,339 41,40 0.0861 0.86 HBB 440,188 31,10 0.0773 0.85 HLA-DQ1 1748,345 0,53 0.0918 0.77 HANNYL 1648,472 1.05 0.7356 0.52 HP 129,478 1,50 0.0218 0.76 LAPTIM4A 388,391 1,35 0.0206 0.78 <	CCR1	338,797	1,01	0.9292	0.52
CTSZ 1580,418 1,28 0.0637 0.76 CXCR4 508,754 1,52 0.0065 0.84 DDI74 36,963 4,34 0.0010 0.96 DNAJC7 105,494 0,92 0.5431 0.62 ENSA 250,287 1,21 0.2580 0.67 FCER1A 410,118 0,54 0.0768 0.73 FKBP5 39,131 2,08 0.0013 0.89 GPER 1,875 7,59 0.0138 0.93 HBA1 5243,339 41,40 0.0861 0.86 HBB 440,188 31,10 0.0773 0.85 HLA-DQ1 1748,345 0,53 0.0918 0.77 HLA-DRB4 1135,072 0,71 0.6316 0.52 HP 129,478 1,50 0.0218 0.63 LAPTMAA 38,391 1,35 0.0226 0.78 LOC100008589 18500877,250 1,12 0.5782 0.63 <t< td=""><td>CD68</td><td>1821,912</td><td>1,27</td><td>0.0640</td><td>0.76</td></t<>	CD68	1821,912	1,27	0.0640	0.76
CXCR4 508,754 1,52 0.0065 0.84 DDIT4 36,963 4,34 0.0010 0.96 DNAJC7 105,494 0,92 0.5431 0.62 ENSA 250,287 1,21 0.2580 0.67 FCER1A 410,118 0,54 0.0768 0.73 FKBP5 39,131 2,08 0.0013 0.89 GPER 1,875 7,59 0.0138 0.93 HBA1 5243,339 41,40 0.0861 0.86 HBB 440,188 31,10 0.0773 0.85 HLA-DQ1 1748,345 0,53 0.0918 0.77 HLA-DRB4 1135,072 0,71 0.6316 0.53 HMOX1 405,989 1,30 0.0729 0.70 HNRNPK 1648,472 1,05 0.7356 0.52 HP 129,478 1,50 0.0218 0.63 LAPTM4A 338,391 1,35 0.0206 0.78	CTSZ	1580,418	1,28	0.0637	0.76
DDIT4 36,963 4,34 0.0010 0.96 DNAJC7 105,494 0,92 0.5431 0.62 ENSA 250,287 1,21 0.2580 0.67 FCER1A 410,118 0,54 0.0768 0.73 FKBP5 39,131 2,08 0.0013 0.89 GPER 1,875 7,59 0.0138 0.93 HBA1 5243,339 41,40 0.0861 0.86 HBB 440,188 31,10 0.0773 0.85 HLA-DQ1 1748,345 0,53 0.0918 0.77 HLA-DRH4 1135,072 0,71 0.6316 0.52 HP 129,478 1,50 0.0218 0.76 IRR2 0,657 0,65 0.2556 0.63 IL1R2 4,794 5,85 0.028 0.87 LOC10008589 18500877,250 1,12 0.5782 0.63 LOC100170939 252,768 0,96 0.8506 0.56 <tr< td=""><td>CXCR4</td><td>508,754</td><td>1,52</td><td>0.0065</td><td>0.84</td></tr<>	CXCR4	508,754	1,52	0.0065	0.84
DNAJC7 105,494 0,92 0.5431 0.62 ENSA 250,287 1,21 0.2580 0.67 FCER1A 410,118 0,54 0.0768 0.73 FKBP5 39,131 2,08 0.0013 0.89 GPER 1,875 7,59 0.0138 0.93 HBA1 5243,339 41,40 0.0861 0.86 HBB 440,188 31,10 0.0773 0.85 HLA-DQ1 1748,345 0,53 0.0918 0.77 HLA-DRB4 1135,072 0,71 0.6316 0.53 HMOX1 405,989 1,30 0.0729 0.70 HNRNPK 1648,472 1,05 0.7356 0.52 HP 129,478 1,50 0.0218 0.67 IER2 0,657 0,655 0.2556 0.63 LOC100008589 18500877,250 1,12 0.5782 0.63 LOC644063 1504,972 1,07 0.0415 0.71	DDIT4	36,963	4,34	0.0010	0.96
ENSA 250,287 1,21 0.2580 0.67 FCER1A 410,118 0,54 0.0768 0.73 FKBP5 39,131 2,08 0.0013 0.89 GPER 1,875 7,59 0.0138 0.93 HBA1 5243,339 41,40 0.0861 0.86 HBB 440,188 31,10 0.0773 0.85 HLA-DQ1 1748,345 0,53 0.0918 0.77 HLA-DRB4 1135,072 0,71 0.6316 0.53 HMOX1 405,989 1,30 0.0729 0.70 HNRNPK 1648,472 1,05 0.7356 0.52 HP 129,478 1,50 0.0218 0.76 IER2 0,657 0,65 0.2556 0.63 IL1R2 4,794 5,85 0.0288 0.87 LOC10008589 18500877,250 1,12 0.5782 0.63 LOC644063 1504,972 1,07 0.0415 0.71 <t< td=""><td>DNAJC7</td><td>105,494</td><td>0,92</td><td>0.5431</td><td>0.62</td></t<>	DNAJC7	105,494	0,92	0.5431	0.62
FCER1A 410,118 0,54 0.0768 0.73 FKBP5 39,131 2,08 0.0013 0.89 GPER 1,875 7,59 0.0138 0.93 HBA1 5243,339 41,40 0.0861 0.86 HBB 440,188 31,10 0.0773 0.85 HLA-DQ1 1748,345 0,53 0.0918 0.77 HLA-DRB4 1135,072 0,71 0.6316 0.53 HMOX1 405,989 1,30 0.0729 0.70 HNRNPK 1648,472 1,05 0.7356 0.52 HP 129,478 1,50 0.0218 0.76 IER2 0,657 0,65 0.2556 0.63 IL1R2 4,794 5,85 0.0288 0.87 LAPTM4A 338,391 1,35 0.2066 0.78 LOC10008589 18500877,250 1,12 0.5782 0.63 LOC100170939 252,768 0,96 0.8506 0.56 LOC644063 1504,972 1,07 0.0415 0.71	ENSA	250.287	1.21	0.2580	0.67
FKBP5 $39,131$ $2,08$ 0.0013 0.89 GPER $1,875$ $7,59$ 0.0138 0.93 HBA1 $5243,339$ $41,40$ 0.0861 0.86 HBB $440,188$ $31,10$ 0.0773 0.85 HLA-DQ1 $1748,345$ $0,53$ 0.0918 0.77 HLA-DRB4 $1135,072$ $0,71$ 0.6316 0.53 HMOX1 $405,989$ $1,30$ 0.0729 0.70 HNRNPK $1648,472$ $1,05$ 0.7356 0.52 HP $129,478$ $1,50$ 0.0218 0.76 IER2 $0,657$ $0,65$ 0.2556 0.63 IL1R2 $4,794$ $5,85$ 0.0268 0.87 LAPTM4A $338,391$ $1,35$ 0.2066 0.78 LOC10008589 $18500877,250$ $1,12$ 0.5782 0.63 LOC644063 $1504,972$ $1,07$ 0.0415 0.71 LOC723972 $568,957$ $1,00$ 0.9645 0.56 LOC644063 $1504,972$ $1,07$ 0.903 0.91 SDHC $313,373$ $1,02$ 0.9145 0.53 SEPT5 0.298 $1,22$ 0.6497 0.50 SLC39A1 1666 $5,39$ 0.2477 0.71 SOCS3 $128,894$ $3,36$ 0.0081 0.91 TAF15 $327,194$ 0.83 0.3084 0.65 TKT $986,111$ $1,48$ 0.0061 0.82 TNPO1 $140,447$ $1,00$ 0.9722 0.52 <	FCER1A	410,118	0.54	0.0768	0.73
GPER 1,875 7,59 0.0138 0.93 HBA1 5243,339 41,40 0.0861 0.86 HBB 440,188 31,10 0.0773 0.85 HLA-DQ1 1748,345 0,53 0.0918 0.77 HLA-DRB4 1135,072 0,71 0.6316 0.53 HMOX1 405,989 1,30 0.0729 0.70 HNRNPK 1648,472 1,05 0.7356 0.52 HP 129,478 1,50 0.0218 0.76 IER2 0,657 0,65 0.2556 0.63 IL1R2 4,794 5,85 0.0288 0.87 LAPTM4A 338,391 1,35 0.0206 0.78 LOC100008589 18500877,250 1,12 0.5782 0.63 LOC40008589 18500877,250 1,12 0.5782 0.63 LOC644063 1504,972 1,07 0.0415 0.71 LOC723972 568,957 1,00 0.9645 0.56 RN28S1 14924567,167 0,92 0.5908 0.53 </td <td>FKBP5</td> <td>39 131</td> <td>2.08</td> <td>0 0013</td> <td>0.89</td>	FKBP5	39 131	2.08	0 0013	0.89
HEA1 5243,339 41,40 0.0861 0.86 HBB 440,188 31,10 0.0773 0.85 HLA-DQ1 1748,345 0,53 0.0918 0.77 HLA-DRB4 1135,072 0,71 0.6316 0.53 HMOX1 405,989 1,30 0.0729 0.70 HNRNPK 1648,472 1,05 0.7356 0.52 HP 129,478 1,50 0.0218 0.76 IER2 0,657 0,65 0.2556 0.63 IL1R2 4,794 5,85 0.0288 0.87 LAPTM4A 338,391 1,35 0.2066 0.78 LOC10008589 18500877,250 1,12 0.5782 0.63 LOC100170939 252,768 0,96 0.8506 0.56 LOC644063 1504,972 1,07 0.0415 0.71 LOC723972 568,957 1,00 0.9645 0.56 RLPL2 186,476 1,02 0.9078 0.53 STIOP 8,494 2,90 0.0003 0.91 <t< td=""><td>GPER</td><td>1 875</td><td>7 59</td><td>0.00138</td><td>0.93</td></t<>	GPER	1 875	7 59	0.00138	0.93
HBR 440,188 31,10 0.0773 0.85 HBB 440,188 31,10 0.0773 0.85 HLA-DQ1 1748,345 0,53 0.0918 0.77 HLA-DRB4 1135,072 0,71 0.6316 0.53 HMOX1 405,989 1,30 0.0729 0.70 HNRNPK 1648,472 1,05 0.7356 0.52 HP 129,478 1,50 0.0218 0.76 IER2 0,657 0,65 0.2556 0.63 ILR2 4,794 5,85 0.0206 0.78 LOC100008589 18500877,250 1,12 0.5782 0.63 LOC100170939 252,768 0,96 0.8506 0.56 LOC643888 308,104 1,46 0.6143 0.55 LOC644063 1504,972 1,07 0.0415 0.71 LOC723972 568,957 1,00 0.9645 0.56 RLPL2 186,476 1,02 0.9078 0.53 STIOP 8,494 2,90 0.0003 0.91 <t< td=""><td></td><td>52/3 330</td><td>11.40</td><td>0.0861</td><td>0.86</td></t<>		52/3 330	11.40	0.0861	0.86
HLD 140,100 01,10 0.0113 0.0316 HLA-DQ1 1748,345 0,53 0.0918 0.77 HLA-DRB4 1135,072 0,71 0.6316 0.53 HMOX1 405,989 1,30 0.0729 0.70 HNRNPK 1648,472 1,05 0.7356 0.52 HP 129,478 1,50 0.0218 0.76 IER2 0,657 0,65 0.2556 0.63 IL1R2 4,794 5,85 0.0206 0.78 LOC100008589 18500877,250 1,12 0.5782 0.63 LOC100008589 18500877,250 1,12 0.5782 0.63 LOC100008589 18500877,250 1,07 0.0415 0.71 LOC643888 308,104 1,46 0.6143 0.55 LOC643888 308,104 1,46 0.6143 0.55 LOC643888 308,104 1,46 0.6143 0.55 LOC643888 308,104 1,46 0.6143 0.53 RN28S1 14924567,167 0,92 0.5908	HBB	<i>14</i> 0 188	31 10	0.0001	0.85
HLA-DQ1 1740,343 0,33 0.0316 0.77 HLA-DRB4 1135,072 0,71 0.6316 0.53 HMOX1 405,989 1,30 0.0729 0.70 HNRNPK 1648,472 1,05 0.7356 0.52 HP 129,478 1,50 0.0218 0.76 IER2 0,657 0,65 0.2556 0.63 IL1R2 4,794 5,85 0.0288 0.87 LAPTM4A 338,391 1,35 0.0206 0.78 LOC100008589 18500877,250 1,12 0.5782 0.63 LOC100170939 252,768 0,96 0.8506 0.56 LOC644063 1504,972 1,07 0.0415 0.71 LOC723972 568,957 1,00 0.9645 0.56 RLPL2 186,476 1,02 0.9078 0.53 S100P 8,494 2,90 0.0003 0.91 SDHC 313,373 1,02 0.9145 0.53 SEPT5 0,298 1,22 0.6497 0.50		1749 345	0.53	0.0775	0.03
HLA-DRB4 1135,072 0,71 0.0316 0.53 HMOX1 405,989 1,30 0.0729 0.70 HNRNPK 1648,472 1,05 0.7356 0.52 HP 129,478 1,50 0.0218 0.76 IER2 0,657 0,65 0.2556 0.63 IL1R2 4,794 5,85 0.0288 0.87 LAPTM4A 338,391 1,35 0.0206 0.78 LOC100008589 18500877,250 1,12 0.5782 0.63 LOC100170939 252,768 0,96 0.8506 0.56 LOC644388 308,104 1,46 0.6143 0.55 LOC644063 1504,972 1,07 0.0415 0.71 LOC723972 568,957 1,00 0.9645 0.56 RLPL2 186,476 1,02 0.9078 0.53 S100P 8,494 2,90 0.0003 0.91 SDHC 313,373 1,02 0.9145 0.53 SEPT5 0,298 1,22 0.6497 0.50 <tr< td=""><td></td><td>1125 072</td><td>0,33</td><td>0.0910</td><td>0.77</td></tr<>		1125 072	0,33	0.0910	0.77
HMRNPK 1648,472 1,50 0.0725 0.70 HP 129,478 1,50 0.0218 0.76 IER2 0,657 0,65 0.2556 0.63 IL1R2 4,794 5,85 0.0288 0.87 LAPTM4A 338,391 1,35 0.0206 0.78 LOC100008589 18500877,250 1,12 0.5782 0.63 LOC100170939 252,768 0,96 0.8506 0.56 LOC643888 308,104 1,46 0.6143 0.55 LOC644063 1504,972 1,07 0.0415 0.71 LOC723972 568,957 1,00 0.9645 0.56 RLPL2 186,476 1,02 0.9078 0.53 S100P 8,494 2,90 0.0003 0.91 SDHC 313,373 1,02 0.9145 0.53 SEPT5 0,298 1,22 0.6497 0.50 SLC39A1 166,812 1,12 0.4069 0.62 SLPI 1,656 5,39 0.2477 0.71		1155,072	0,71	0.0310	0.53
HNRNPK $1048,4/2$ $1,05$ 0.7356 0.52 HP $129,478$ $1,50$ 0.0218 0.76 IER2 $0,657$ $0,65$ 0.2556 0.63 IL1R2 $4,794$ $5,85$ 0.0288 0.87 LAPTM4A $338,391$ $1,35$ 0.0206 0.78 LOC100008589 $18500877,250$ $1,12$ 0.5782 0.63 LOC100170939 $252,768$ $0,96$ 0.8506 0.56 LOC643888 $308,104$ $1,46$ 0.6143 0.55 LOC644063 $1504,972$ $1,07$ 0.0415 0.71 LOC723972 $568,957$ $1,00$ 0.9645 0.56 RLPL2 $186,476$ $1,02$ 0.9078 0.53 S100P $8,494$ $2,90$ 0.0003 0.91 SDHC $313,373$ $1,02$ 0.9145 0.53 SEPT5 $0,298$ $1,22$ 0.6497 0.50 SLC39A1 $166,812$ $1,12$ 0.4069 0.62 SLPI $1,656$ $5,39$ 0.2477 0.71 SOCS3 $128,894$ $3,36$ 0.0081 0.91 TAF15 $327,194$ $0,83$ 0.3084 0.65 TKT $986,111$ $1,48$ 0.0061 0.82 TNF $28,056$ $0,94$ 0.7325 0.53 TNPO1 $140,447$ $1,00$ 0.9722 0.52		405,969	1,30	0.0729	0.70
HP 129,478 1,50 0.0218 0.76 IER2 0,657 0,65 0.2556 0.63 IL1R2 4,794 5,85 0.0288 0.87 LAPTM4A 338,391 1,35 0.0206 0.78 LOC100008589 18500877,250 1,12 0.5782 0.63 LOC100170939 252,768 0,96 0.8506 0.56 LOC644063 1504,972 1,07 0.0415 0.71 LOC723972 568,957 1,00 0.9645 0.56 RLPL2 186,476 1,02 0.9078 0.53 RN28S1 14924567,167 0,92 0.5908 0.58 S100P 8,494 2,90 0.0003 0.91 SDHC 313,373 1,02 0.9145 0.53 SEPT5 0,298 1,22 0.6497 0.50 SLC39A1 166,812 1,12 0.4069 0.62 SLPI 1,656 5,39 0.2477 0.71 SOCS3 128,894 3,36 0.0081 0.91		1648,472	1,05	0.7356	0.52
IER2 0,657 0,65 0.2556 0.63 IL1R2 4,794 5,85 0.0288 0.87 LAPTM4A 338,391 1,35 0.0206 0.78 LOC100008589 18500877,250 1,12 0.5782 0.63 LOC100170939 252,768 0,96 0.8506 0.56 LOC643888 308,104 1,46 0.6143 0.55 LOC644063 1504,972 1,07 0.0415 0.71 LOC723972 568,957 1,00 0.9645 0.56 RLPL2 186,476 1,02 0.9078 0.53 RN28S1 14924567,167 0,92 0.5908 0.58 S100P 8,494 2,90 0.0003 0.91 SDHC 313,373 1,02 0.9145 0.53 SEPT5 0,298 1,22 0.6497 0.50 SLC39A1 166,812 1,12 0.4069 0.62 SLPI 1,656 5,39 0.2477 0.71 SOCS3 128,894 3,36 0.0081 0.91 <	HP	129,478	1,50	0.0218	0.76
L1R24,7945,850.02880.87LAPTM4A338,3911,350.02060.78LOC10000858918500877,2501,120.57820.63LOC100170939252,7680,960.85060.56LOC643888308,1041,460.61430.55LOC6440631504,9721,070.04150.71LOC723972568,9571,000.96450.56RLPL2186,4761,020.90780.53RN28S114924567,1670,920.59080.58S100P8,4942,900.00030.91SDHC313,3731,020.91450.53SEPT50,2981,220.64970.50SLC39A1166,8121,120.40690.62SLPI1,6565,390.24770.71SOCS3128,8943,360.00810.91TAF15327,1940,830.30840.65TKT986,1111,480.00610.82TNF28,0560,940.73250.53TNPO1140,4471,000.97220.52	IER2	0,657	0,65	0.2556	0.63
LAPTM4A 338,391 1,35 0.0206 0.78 LOC100008589 18500877,250 1,12 0.5782 0.63 LOC100170939 252,768 0,96 0.8506 0.56 LOC643888 308,104 1,46 0.6143 0.55 LOC644063 1504,972 1,07 0.0415 0.71 LOC723972 568,957 1,00 0.9645 0.56 RLPL2 186,476 1,02 0.9078 0.53 RN28S1 14924567,167 0,92 0.5908 0.58 S100P 8,494 2,90 0.0003 0.91 SDHC 313,373 1,02 0.9145 0.53 SEPT5 0,298 1,22 0.6497 0.50 SLC39A1 166,812 1,12 0.4069 0.62 SLPI 1,656 5,39 0.2477 0.71 SOCS3 128,894 3,36 0.0081 0.91 TAF15 327,194 0,83 0.3084 0.65 TKT 986,111 1,48 0.0061 0.82	IL1R2	4,794	5,85	0.0288	0.87
LOC10000858918500877,2501,120.57820.63LOC100170939252,7680,960.85060.56LOC643888308,1041,460.61430.55LOC6440631504,9721,070.04150.71LOC723972568,9571,000.96450.56RLPL2186,4761,020.90780.53RN28S114924567,1670,920.59080.58S100P8,4942,900.00030.91SDHC313,3731,020.91450.53SEPT50,2981,220.64970.50SLC39A1166,8121,120.40690.62SLPI1,6565,390.24770.71SOCS3128,8943,360.00810.91TAF15327,1940,830.30840.65TKT986,1111,480.00610.82TNF28,0560,940.73250.53TNPO1140,4471,000.97220.52	LAPTM4A	338,391	1,35	0.0206	0.78
LOC100170939252,7680,960.85060.56LOC643888308,1041,460.61430.55LOC6440631504,9721,070.04150.71LOC723972568,9571,000.96450.56RLPL2186,4761,020.90780.53RN28S114924567,1670,920.59080.58S100P8,4942,900.00030.91SDHC313,3731,020.91450.53SEPT50,2981,220.64970.50SLC39A1166,8121,120.40690.62SLPI1,6565,390.24770.71SOCS3128,8943,360.00810.91TAF15327,1940,830.30840.65TKT986,1111,480.00610.82TNF28,0560,940.73250.53TNPO1140,4471,000.97220.52	LOC100008589	18500877,250	1,12	0.5782	0.63
LOC643888 $308,104$ $1,46$ 0.6143 0.55 LOC644063 $1504,972$ $1,07$ 0.0415 0.71 LOC723972 $568,957$ $1,00$ 0.9645 0.56 RLPL2 $186,476$ $1,02$ 0.9078 0.53 RN28S1 $14924567,167$ $0,92$ 0.5908 0.58 S100P $8,494$ $2,90$ 0.0003 0.91 SDHC $313,373$ $1,02$ 0.9145 0.53 SEPT5 $0,298$ $1,22$ 0.6497 0.50 SLC39A1 $166,812$ $1,12$ 0.4069 0.62 SLPI $1,656$ $5,39$ 0.2477 0.71 SOCS3 $128,894$ $3,36$ 0.0081 0.91 TAF15 $327,194$ $0,83$ 0.3084 0.65 TKT $986,111$ $1,48$ 0.0061 0.82 TNF $28,056$ $0,94$ 0.7325 0.53 TNFO1 $140,447$ $1,00$ 0.9722 0.52	LOC100170939	252,768	0,96	0.8506	0.56
LOC6440631504,9721,070.04150.71LOC723972568,9571,000.96450.56 $RLPL2$ 186,4761,020.90780.53 $RN28S1$ 14924567,1670,920.59080.58 $S100P$ 8,4942,900.00030.91 $SDHC$ 313,3731,020.91450.53 $SEPT5$ 0,2981,220.64970.50 $SLC39A1$ 166,8121,120.40690.62 $SLPI$ 1,6565,390.24770.71 $SOCS3$ 128,8943,360.00810.91 $TAF15$ 327,1940,830.30840.65 TKT 986,1111,480.00610.82 TNF 28,0560,940.73250.53 $TNPO1$ 140,4471,000.97220.52	LOC643888	308,104	1,46	0.6143	0.55
LOC723972 $568,957$ $1,00$ 0.9645 0.56 RLPL2 $186,476$ $1,02$ 0.9078 0.53 RN28S1 $14924567,167$ $0,92$ 0.5908 0.58 S100P $8,494$ $2,90$ 0.0003 0.91 SDHC $313,373$ $1,02$ 0.9145 0.53 SEPT5 $0,298$ $1,22$ 0.6497 0.50 SLC39A1 $166,812$ $1,12$ 0.4069 0.62 SLPI $1,656$ $5,39$ 0.2477 0.71 SOCS3 $128,894$ $3,36$ 0.0081 0.91 TAF15 $327,194$ $0,83$ 0.3084 0.65 TKT $986,111$ $1,48$ 0.0061 0.82 TNF $28,056$ $0,94$ 0.7325 0.53 TNPO1 $140,447$ $1,00$ 0.9722 0.52	LOC644063	1504,972	1,07	0.0415	0.71
RLPL2 186,476 1,02 0.9078 0.53 RN28S1 14924567,167 0,92 0.5908 0.58 S100P 8,494 2,90 0.0003 0.91 SDHC 313,373 1,02 0.9145 0.53 SEPT5 0,298 1,22 0.6497 0.50 SLC39A1 166,812 1,12 0.4069 0.62 SLPI 1,656 5,39 0.2477 0.71 SOCS3 128,894 3,36 0.0081 0.91 TAF15 327,194 0,83 0.3084 0.65 TKT 986,111 1,48 0.0061 0.82 TNF 28,056 0,94 0.7325 0.53 TNPO1 140,447 1,00 0.9722 0.52	LOC723972	568,957	1,00	0.9645	0.56
RN28S1 14924567,167 0,92 0.5908 0.58 S100P 8,494 2,90 0.0003 0.91 SDHC 313,373 1,02 0.9145 0.53 SEPT5 0,298 1,22 0.6497 0.50 SLC39A1 166,812 1,12 0.4069 0.62 SLPI 1,656 5,39 0.2477 0.71 SOCS3 128,894 3,36 0.0081 0.91 TAF15 327,194 0,83 0.3084 0.65 TKT 986,111 1,48 0.0061 0.82 TNF 28,056 0,94 0.7325 0.53 TNPO1 140,447 1,00 0.9722 0.52	RLPL2	186,476	1,02	0.9078	0.53
S100P 8,494 2,90 0.0003 0.91 SDHC 313,373 1,02 0.9145 0.53 SEPT5 0,298 1,22 0.6497 0.50 SLC39A1 166,812 1,12 0.4069 0.62 SLPI 1,656 5,39 0.2477 0.71 SOCS3 128,894 3,36 0.0081 0.91 TAF15 327,194 0,83 0.3084 0.65 TKT 986,111 1,48 0.0061 0.82 TNF 28,056 0,94 0.7325 0.53 TNPO1 140,447 1,00 0.9722 0.52	RN28S1	14924567,167	0,92	0.5908	0.58
SDHC 313,373 1,02 0.9145 0.53 SEPT5 0,298 1,22 0.6497 0.50 SLC39A1 166,812 1,12 0.4069 0.62 SLPI 1,656 5,39 0.2477 0.71 SOCS3 128,894 3,36 0.0081 0.91 TAF15 327,194 0,83 0.3084 0.65 TKT 986,111 1,48 0.0061 0.82 TNF 28,056 0,94 0.7325 0.53 TNPO1 140,447 1,00 0.9722 0.52	S100P	8,494	2,90	0.0003	0.91
SEPT5 0,298 1,22 0.6497 0.50 SLC39A1 166,812 1,12 0.4069 0.62 SLPI 1,656 5,39 0.2477 0.71 SOCS3 128,894 3,36 0.0081 0.91 TAF15 327,194 0,83 0.3084 0.65 TKT 986,111 1,48 0.0061 0.82 TNF 28,056 0,94 0.7325 0.53 TNPO1 140,447 1,00 0.9722 0.52	SDHC	313,373	1,02	0.9145	0.53
SLC39A1166,8121,120.40690.62SLPI1,6565,390.24770.71SOCS3128,8943,360.00810.91TAF15327,1940,830.30840.65TKT986,1111,480.00610.82TNF28,0560,940.73250.53TNPO1140,4471,000.97220.52	SEPT5	0,298	1,22	0.6497	0.50
SLPI 1,656 5,39 0.2477 0.71 SOCS3 128,894 3,36 0.0081 0.91 TAF15 327,194 0,83 0.3084 0.65 TKT 986,111 1,48 0.0061 0.82 TNF 28,056 0,94 0.7325 0.53 TNPO1 140,447 1,00 0.9722 0.52	SLC39A1	166,812	1,12	0.4069	0.62
SOCS3 128,894 3,36 0.0081 0.91 TAF15 327,194 0,83 0.3084 0.65 TKT 986,111 1,48 0.0061 0.82 TNF 28,056 0,94 0.7325 0.53 TNPO1 140,447 1,00 0.9722 0.52	SLPI	1,656	5,39	0.2477	0.71
TAF15 327,194 0,83 0.3084 0.65 TKT 986,111 1,48 0.0061 0.82 TNF 28,056 0,94 0.7325 0.53 TNPO1 140,447 1,00 0.9722 0.52	SOCS3	128,894	3,36	0.0081	0.91
TKT 986,111 1,48 0.0061 0.82 TNF 28,056 0,94 0.7325 0.53 TNPO1 140,447 1,00 0.9722 0.52	TAF15	327,194	0,83	0.3084	0.65
TNF 28,056 0,94 0.7325 0.53 TNPO1 140,447 1,00 0.9722 0.52	ткт	986,111	1,48	0.0061	0.82
TNPO1 140,447 1,00 0.9722 0.52 ^a Moon expression of gaps of interset (10,000 contics of B214)	TNF	28.056	0.94	0.7325	0.53
^a Maan avaraasian of gang of interact / 10,000 apping of D0/4	TNPO1	140.447	1.00	0.9722	0.52
	aMoon overcosier	of appoint inter-	ot / 10,000 occ	ion of POM	3.02

^bFold ratio of patients compared to healthy volunteers ^cBold print indicates where cutoff criteria (p<0.1, AUC>0.7) are met. See main manuscript and Supplementary methods for more detailed information to the standard stand

ouppieme			
Gene	Full Name	Biological Function	Potential Function in Monocytes
ACP5	acid phosphatase 5,	iron containing glycoprotein	negative regulation of inflammatory
	tartrate resistant	involved in adhesion and	response in interleukin pathways
		migration	
ADM	adrenomedullin	vasodilation, regulation of	antimicrobial activity, wound healing
		hormone secretion,	
		promotion of angiogenesis	
APP	amyloid beta (A4)	protein basis of amyloid	antimicrobial activity, mitotic activity
	precursor protein	plaques in Alzheimer disease	
BAX	BCL2-associated X	p53-mediated activator of	myeloid cell homeostasis
	protein	apoptosis	
CD68	CD68 molecule	integral membran	highly expressed on monocytes and
		glycoprotein of scavenger	macrophages, mediator of recruitment
		receptor family	and activation
CTSZ	cathepsin Z	lysosomal cystein	unknown
		proteinase, involved in	
		migration and adhesion	
CXCR4	chemokine (C-X-C motif)	CXC chemokine receptor	mediator of recruitment. chemotaxis.
	receptor 4	specific for stromal cell-	and activation
		derived factor-1	
DDIT4	DNA-damage-inducible	negative regulation of mTOR	defense response to microbial signals
	transcript 4	signalling upon cellular	
		stress	
FCFR1A	Ec fragment of IgE high	alpha subunit of IgE-	positive regulation of type-Limmune
I OLIVIII	affinity L recentor for:	mediated allergic response	response and macrophage
	alpha polypentide	inculated allergic response	differentiation
EKBD5	EK506 binding protein 5	member of immunophilin	recentor for EK506 and ranamycin
I NDI 5	r 1000 binding protein 5	protein family	mediating calcineurin inhibition
		immunorogulation	
CDED	C protoin coupled		pogativo regulator of loukoovto
GPER	G protein-coupled	non-genomic signaling of	
	estrogen receptor 1	estrogen stimulus	
пват	nemoglobin, alpha 1		UNKNOWN
	have all him hate	NDA	negitive regulation of situic sydds
пвв	nemoglobin, beta		positive regulation of hitric oxide
	we sign bists a sure stibility		synthesis,
HLA-DQA1	major histocompatibility	MHC class II receptor	antigen processing and presentation
	complex, class II, DQ	activity; peptide antigen	
	alpha 1	binding	
HMOX1	heme oxygenase	heme catabolism	regulation of phagocytosis and
	(decycling) 1		migration, chemokine synthesis, wound
			healing, and angiogenesis
HP	haptoglobin	preproprotein of haptoglobin	acute-phase defense response
		subunit	
IL1R2	interleukin 1 receptor,	cytokine receptor for IL-1	cytokine-mediated immune response
	type II		
LAPTM4A	lysosomal protein	unknown	unknown
	transmembrane 4 alpha		
LOC644063	heterogeneous nuclear	unknown	unknown
	ribonucleoprotein K		
	pseudogene 4		
S100P	S100 calcium binding	cell cycle progression and	unknown
	protein P	differentiation	
SLPI	secretory leukocyte	secreted inhibitor of serin	negative regulation of endopeptidase
	peptidase inhibitor	proteinases	activity
SOCS3	suppressor of cytokine	negative regulator of	modulator of immune response
	signaling 3	cvtokine signalling	particularly IFN-y mediated
ткт	transketolase	enzyme of pentose	metabolic modulator
		phosphate pathway	

Supplementary Table 5: Identity and Function of the gene signature members

SUPPLEMENTARY REFERENCES

- 1. Weitz J, Koch M, Debus J, et al. Colorectal cancer. Lancet 2005;**365**(9454):153-65.
- Nyugen J, Agrawal S, Gollapudi S, et al. Impaired functions of peripheral blood monocyte subpopulations in aged humans. Journal of clinical immunology 2010;**30**(6):806-13.
- Du P, Kibbe WA, Lin SM. Iumi: a pipeline for processing Illumina microarray. Bioinformatics 2008;24(13):1547-8.
- 4. Lin SM, Du P, Huber W, et al. Model-based variance-stabilizing transformation for Illumina microarray data. Nucleic acids research 2008;**36**(2):e11.
- Smyth GK. Linear models and empirical bayes methods for assessing differential expression in microarray experiments. Statistical applications in genetics and molecular biology 2004;3:Article3.
- 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(1):289-300.
- 7. Sample size for microarray experiments. Secondary Sample size for microarray experiments. <u>http://bioinformatics.mdanderson.org/MicroarraySampleSize/</u>.
- 8. Dietterich TG. Ensemble methods in machine learning. Lecture Notes in Computer Science 2000;**1857**:1-15.

9. Impute: Imputation for microarray data. [program]. 1.32.0 version, 2013.

- 10. Robin X, Turck N, Hainard A, et al. pROC: an open-source package for R and S+ to analyze and compare ROC curves. BMC bioinformatics 2011;**12**:77.
- Burges CJC. A Tutorial on Support Vector Machines for Pattern Recognition.
 Data Min Knowl Discov 1998;2(2):121-67.

- Liaw A, Wiener M. Classification and Regression by randomForest. R News: The Newsletter of the R Project 2002;2(3):18-22.
 - 13. Ambroise C, McLachlan GJ. Selection bias in gene extraction on the basis of microarray gene-expression data. Proceedings of the National Academy of Sciences of the United States of America 2002;**99**(10):6562-6.
 - Piehler A, Grimholt R, Ovstebo R, et al. Gene expression results in lipopolysaccharide-stimulated monocytes depend significantly on the choice of reference genes. BMC Immunology 2010;11(1):21.
 - 15. Guo C, Liu S, Wang J, et al. ACTB in cancer. Clinica chimica acta; international journal of clinical chemistry 2013;**417**:39-44.
- 16. Khatri P, Sirota M, Butte AJ. Ten Years of Pathway Analysis: Current Approaches and Outstanding Challenges. PLoS Comput Biol 2012;**8**(2):e1002375.
- 17. Jess P, Hansen IO, Gamborg M, et al. A nationwide Danish cohort study challenging the categorisation into right-sided and left-sided colon cancer.
 BMJ open 2013;3(5).

Supplementary Material

Tumour-Educated Circulating Monocytes are Powerful Candidate Biomarkers for Diagnosis and Disease Follow-up of Colorectal Cancer

Alexander Hamm, Hans Prenen, Wouter Van Delm, Mario Di Matteo, Mathias Wenes, Estelle Delamarre, Thomas Schmidt, Jürgen Weitz, Roberta Sarmiento, Angelo Dezi, Giampietro Gasparini, Françoise Rothé, Robin Schmitz, André D'Hoore, Hannes Iserentant, Alain Hendlisz & Massimiliano Mazzone

https://mc.manuscriptcentral.com/gut

- 3 4	CONTENT	S
	Supplementary Methods	Page 3
3	Supplementary Notes	Page 14
, 0 1	Supplementary Figures	Page 17
12 13	Supplementary Tables	Page 28
14 15 16 17 18 9 22 22 22 22 22 22 22 22 22 22 22 22 2	Supplementary References	Page 35
56 57 58 59 60		
	https://mc.manuscriptce	ntral com/quit

SUPPLEMENTARY METHODS

Patients

The composition of patient cohorts is given in detail in the main manuscript. Inclusion criteria for patients were sporadic histologically confirmed adenocarcinoma of the colon and/or rectum for cohort I-III and VI, patients in remission from CRC for a treatment-free interval of minimum 3 months for cohort V, histologically confirmed adenocarcinoma of the stomach or gastroesophageal junction or of the pancreas, or histologically confirmed gastritis for cohort IV. All patient samples were prospectively collected after histological diagnosis upon screening colonoscopy (reference standard defined by international clinical guidelines¹), prior to any treatment, at clinically indicated regular appointments separate of medical interventions (such as colonoscopy, surgical preparations etc.). All newly diagnosed patients presenting to the responsible clinicians were consecutively included when they met criteria and gave written informed consent. Healthy volunteers were included when there was no evidence or record of acute or chronic disease, with identical exclusion criteria as the patients. A subset of healthy individuals (within cohort III) was included upon screening colonoscopy without any pathological findings. Exclusion criteria were age of less than 40 years (to exclude cancers suspicious of genetic syndromes and restrict possible age-related variations in the monocyte phenotype reported previously²), history of oncological, chronic inflammatory, and autoimmune diseases within 10 years prior to this study, clinical or laboratory evidence of acute infection, anti-inflammatory and/or immunosuppressive medication within 90 days of blood sampling with the exception of occasional NSAID, commencement of medical or surgical anti-cancer treatment, medication with sedatives or opioid-based analgesics within 72 hours prior to blood sampling, clinical or microbiological evidence of altered

Gut

gut flora. Samples were excluded from further analysis when final histology of the surgical specimen did not confirm adenocarcinoma of the large intestine (assessed by board-certified pathologists within clinical routine procedures). The following four oncological centres contributed samples to this study: Digestive Oncology, University Hospitals Leuven and Department of Oncology, KU Leuven, Leuven, Belgium; Department of General, Visceral, and Transplantation Surgery, University of Heidelberg, Heidelberg, Germany; Department of Oncology, San Filippo Neri, Rome, Italy; Medical Oncology Clinic, Institut Jules Bordet, Brussels, Belgium. The responsible scientists in each centre (1-2 per centre) were trained in the protocol for isolation of PBM to ensure uniformity of the procedure. All participants gave written informed consent, and the study was approved by the respective institutional review boards (Leuven: B322201215873, Brussels: CE1950, Heidelberg: 323/2004, Rome: 319/51). No adverse events from blood collection or colonoscopy were recorded in included participants.

Isolation of PBM

20ml of EDTA-anticoagulated peripheral venous blood was collected following clinical routine procedure, stored at 4°C and processed within 2 hours of blood collection. For further isolation, blood was diluted 1:2 with DPBS (free of Ca2+ and Mg2+) and layered carefully on Lymphoprep (Axis-Shield) in two separate tubes. All blood collection and isolation steps were performed identical for samples of all origin. Density gradient centrifugation was performed at 1,200g for 20 minutes at low acceleration and no brake. <u>Samples with macroscopically visible hemolysis were excluded from further analysis.</u> The PBMC interphase was collected carefully and washed twice for 12 minutes at 250g and 175g with PBS. Hemocytrometric analysis

was performed to ensure purity of PBMCs, and the pellet was pooled for further processing and washed once for 10 minutes at 300g. Cells were then incubated with CD14 magnetically-conjugated beads (BD) for 15 minutes at 4°C, washed 10 minutes at 300g and positively separated with the MACS system (Miltenyi) following the manufacturer's instructions. The CD14+ fraction was flushed out and washed once 10 minutes at 300g. Purity was assessed by FACS analysis for CD14 in the pilot phase and by hemocytometric analysis (CellDyn 3700, Abbott) in every further sample. Only samples with purity of >90% and viability >95% (assessed by Trypan Blue staining) were retained for further analysis. Cell pellets were lysed in Buffer RLT (Qiagen) at 10⁶ monocytes in 350µl of Buffer RLT and stored at -80 ℃. For each respective expression study, all samples were extracted simultaneously with the RNeasy Mini Kit (Qiagen) following the manufacturer's instructions. Quality control was performed by checking RNA guality on the Nanodrop system, and RNA integrity was checked for microarray samples on the Agilent Bio-Analyzer. Only samples with an extinction fraction 260/280 > 1.8 and 260/230 > 1.5, and an RNA integrity index of >6 were retained for further analysis.

Genome-wide expression analysis

For genome-wide expression analysis, RNA was amplified and biotinylated using Illumina TotalPrep RNA Amplification Kit (Ambion) following the manufacturer's instructions to obtain biotinylated cRNA, which was hybridized to Illumina HumanHT-12 v4 Expression BeadChips (Illumina) with the Illumina Whole-Genome Gene Expression Direct Hybridization Assay (Illumina) following the manufacturer's instructions. The Illumina HumanHT-12 v4 Expression BeadChip Kit contains 47,323 probes and 887 controls. After scanning, background-corrected expression values

and detection scores were extracted with GenomeStudio GX (version 1.5.4). For each array, we used the summarized expression level (AVG_Signal), standard error of the bead replicates (BEAD_STERR), number of beads used (AVG_NBEADS) and a detection score, which estimates the probability of a gene being detected above the background. Resulting expression data was analyzed with R, using the lumi package³. A variance stabilizing transformation⁴ was applied, followed by quantile normalization to compensate for batch effects of the individual bead chips. For each probe, the number of present calls over all samples was determined (the threshold on the detection was p<0.01), and probes absent in all samples were omitted in the analysis. This omitted subset consisted of 18,396 probes. Hence, analysis was performed for 28,927 probes. Differential expression was assessed with the limma package of \mathbb{R}^5 .

Quantitative RT-PCR (qPCR)

For qPCR analyses, 400ng of RNA was reverse transcribed with SuperScript III First Strand Kit (Invitrogen) following the manufacturer's instructions, and qPCR was performed in duplicates on a 7500Fast System (Applied Biosystems) using intronspanning PrimeTime qPCR Assays (Integrated DNA Technologies) listed in Supplementary Table 2. Wherever possible, qPCR assays were selected that covered the exon in which the Illumina Expression BeadChip probe was located. Raw data was analyzed with SDS v1.4 (Applied Biosystems), and expression was normalized within samples with the $\Delta\Delta$ CT method to reference gene *B2M*. Data was expressed relative to the average expression of that gene in the healthy volunteers in the dataset. Data points where duplicates differed by more than 1 CT were discarded. Inter-run validity was verified by both processing and running previously

analyzed samples as internal controls and ensuring correct clustering within their respective groups. Where necessary for normalization purposes, stored and validated healthy volunteer samples were re-profiled along with samples from cohorts IV and V.

Identification of a gene signature

For each pair-wise comparison between HV, P and PM, we evaluated all probes with a moderated t-test, as implemented in the limma-package⁵ of R. P-values were adjusted for multiple testing with Benjamini-Hochberg to control the false discovery rate⁶. A probe was selected as being differentially expressed between two groups when the adjusted p-value was smaller than 0.05 and the fold change exceeded 1.5 times up- or down-regulation ($\log_2 > 0.58$ or < -0.58, respectively). For the comparison between PM/P and HV, differential expression of the selected genes was further validated with qPCR in 8 randomly selected individuals from each of the groups in cohort I. The panel of 35 candidate genes derived from the 40 Illumina probes differentially expressed in cohort I was augmented by 8 genes which marginally missed the applied cutoff criteria and had been identified in unpublished in vitro and in vivo screens during the pilot phase. Minimal sample size for further cohorts was chosen to be 15 after conducting a statistical power analysis with the data from cohort I to estimate the expected variation in gene expression. Sample size was chosen to achieve a statistical power of 0.9 with an ordinary t-test when fold changes of 1.5 are considered and 5% false positives are accepted. Power calculations were done with the online tool from the Department of Bioinformatics and Computational Biology of MD Anderson Cancer Center⁷. Differential expression was considered to be confirmed by qPCR when the p-value after a two-tailed

unpaired t-test was smaller than 0.1 and/or the associated area under the ROC curve (AUC) was larger than 0.7. as calculated with Prism (GraphPad, Inc.). We chose deliberately for loose cut-offs on p-value and AUC for the confirmation, since less distinctly differentially expressed genes could in theory still add value to a (later developed) multiple-gene classification strategy.

Multicentric validation study

Overview. The diagnostic test consists of a gene panel assay in combination with software for decision support. The software implements an algorithm that takes the data from the assay as input and outputs a binary decision: whether the profiled sample comes from a CRC patient or not. The algorithm is an ensemble method (ENS)⁸ that consults 3 subroutines, then counts the number of votes in favor of CRC and finally proposes the decision that is supported by at least 2 subroutines. The 3 subroutines form a heterogeneous set of alternative classification algorithms: an easily interpretable ensemble stump classifier (SGMV – single gene majority vote), a linear support vector machine (SVM) and a more complex random forest (RF). The parameters of the 3 subroutines were fitted in parallel to a subset of samples from the multi-centric cohort II. This training subset was constructed via stratified random sampling. Performance of the algorithm was assessed through a Monte Carlo cross-validation (MCCV) procedure on the training data and further validated on the samples from cohort II that were excluded during training.

Stratified random sampling. We identified combinations of the four oncology centres and two sample classes (i.e. HV or CRC) as 8 strata. From each stratum, we sampled 2 times as much training samples as validation samples. The actual number of samples per stratum was chosen so that i. there was no evidence of dependence

of class labeling on centre in either validation or training dataset, ii. the final datasets were balanced (i.e. as much HV as CRC). Dependence between class labeling and centre of origin was excluded by testing with a Fisher's exact test (p > 0.93). The random split was performed prior to fitting parameters and retained for all further analyses to obtain realistic measures of classification performance. Since our subroutines required complete data, we imputed missing values after assembling the training and validation datasets for each dataset separately using nearest neighbor averaging, as implemented in the impute-package in \mathbb{R}^9 .

Subroutines. The SGMV compares the expression value of each input gene first to a gene-specific cut-off and then assigns a defined class to an unknown sample depending on whether the cut-offs are exceeded for at least half of the genes (i.e. majority vote). The SGMV parameters hence consist of gene-specific cut-offs. The gene-specific cut-offs are fitted by taking that value that corresponds to the point closest to the top-left corner of the gene-associated ROC curve, using the pROCpackage in R¹⁰. The SVM with linear kernel is similar to linear discriminant analysis, taking as input the expression values of a set of genes and comparing a linear combination of the input values to a threshold in order to assign a defined class to an unknown sample, thereby giving higher weight to more informative genes. The SVM parameters hence consist of gene-specific weights and one threshold. We fitted the parameters with the kernlab-package in R¹¹. The RF pushes the expression values of a set of genes through a multitude of decision trees (each looking at a random subset of genes and built from a random subset of samples from the training data), notes down for each class the proportion of supporting individual trees and finally assigns the class with highest support. The RF parameters hence consist of individual decision trees. We fitted the parameters with the randomForest-package in R¹².

Avoiding over-fitting. Fitting the parameters of the SVM and RF subroutines was conditioned on hyper-parameters that influence the flexibility of the subroutines to fit the training data. Too flexible procedures lead to over-fitting of training samples at the cost of bad performance on unseen samples. Flexibility was therefore constrained by selecting hyper-parameters from a range of options with Monte Carlo cross-validation (MCCV), prior to final determination of the common parameters. We divided the training dataset during 100 cycles in 2/3 and 1/3, trained the SVM/RF each time on the largest part with a given hyper-parameter, tested the SVM/RF each time on the smallest part and finally averaged the AUC and BER of all cycles for a particular hyper-parameter value. We chose the hyper-parameter with best average AUC, or in case of multiple options, the one with best average BER. Note that this MCCV procedure to select hyper-parameters was also run as an inner loop within the outer MCCV loop when algorithm performance was assessed (see above)¹³.

Performance metrics. The classifiers were validated on the qPCR test dataset, constructed from healthy volunteers and patients of multi-centric cohort II who were not included during development of the models (see above). To verify the similarity of the test set to the training set, a Spearman-correlation between all assays was performed, ensuring that test assays did not cluster separately from training assays. A separate clustering would have been an indication that the training dataset was not representative for the test samples. Two types of performance were finally reported: ranking performance and classification performance. Ranking performance is the capability of an algorithm to give a higher score to an individual from class CRC than to an individual from class HV. We measured ranking performance by the area under the ROC curve (AUC). For all 4 routines (SGMV, SVM, RF and ENS), we provided the AUC as well as the lower bound and upper bound of its 95% confidence interval,

as computed after 2,000 bootstraps with the pROC-package in R¹⁰. Classification performance measures the capability of an algorithm to assign an individual to the correct class. We reported for all routines the balanced error rate (BER), sensitivity (Se) and specificity (Sp). For Se and Sp, we also computed the lower bound and upper bound of the 95% confidence interval after 2,000 bootstraps.

Complementary data analysis

A complementary data analysis by an independent team (DNAlytics, Belgium) on the same 23-marker signature led to the same conclusions in terms of performances. Another (per-marker) normalization procedure has been proposed. This normalization is applied on the log-transformed gene expression (i.e. Δ CT values) and consists in computing, on the training set (for example Cohort II, both HV and CRC), the mean and standard deviation of each marker. When a prediction has to be made on a new, potentially isolated sample, each marker measurement of this new sample is normalized by subtracting the corresponding mean, and by dividing by the corresponding standard deviation. A modified procedure has also been proposed for the imputation of missing values, making it dependent on the reference cohort only. This avoids the need for a new reference HV batch as prediction has to be made on a new (set of) sample(s).

The first experiment consisted in cross-validating a model on Cohort II (BER: 8.4% [3.4%;13.4%]; AUC: 0.93 [0.88;0.98]). A second experiment consisted in learning the same type of model on Cohort II and having it make predictions on Cohort III (BER: 13.2%; AUC: 0.92). All analyses were performed in R with scripts designed by DNAlytics, fully independent from other analyses described in this paper.

In vitro model system

To study the effects of tumour-released soluble factors on gene expression in monocytes, we established an in vitro model system. Medium conditioned with cell-released soluble factors was obtained by seeding the following cell lines at 40% confluence at 37 °C at 21% O₂, 5% CO₂ in a moist atmosphere in their respective medium and ultra-filtering the conditioned medium 72 hours later: HCT116 (new from ATCC, CCL-247) in RPMI (10% FBS, 1% Glutamine, 1% PenStrep), grown in normoxia or hypoxia (1% O2), CCD 841 CoN (new from ATCC CRL-1790) in EMEM (10% FBS, 1% Glutamine, 1% PenStrep), MKN-45 (a kind gift from Frans van Roy, UGent, Belgium) in RPMI (10% FBS, 1% Glutamine, 1% PenStrep, 1% Na-Pyruvate). Each medium was also incubated separately without cells to obtain the respective mock controls. Absence of Mycoplasma species was verified with MycoAlert Mycoplasma Detection Kit (Lonza).

Monocytes from healthy volunteers (n=6) were isolated as described above and were seeded at 200,000 cells / well in a tissue-culture treated 24-well plate (Costar) in IMDM (10% autologous serum, 1% Glutamine), supplemented 1:5 with conditioned medium. Cells were lysed in Buffer RLT (Qiagen) after 18 hours. For experiments on reversion of the gene signature after withdrawing the stimulus, monocytes were washed with PBS after 18 hours of culture in conditioned medium, and medium was refreshed with plain IMDM (10% autologous serum, 1% Glutamine). After 72 hours, cells were then lysed in Buffer RLT. All experiments were performed in technical quadruplicates and repeated at least twice.

All RNA was extracted simultaneously with the RNeasy MicroKit (Qiagen) following the manufacturer's instructions, and RNA quality was verified with the Nanodrop system as described above.

Gut

<page-header><text><text><text>

SUPPLEMENTARY NOTES

Supplementary Note 1

To select a robust reference gene, we checked in the available microarray data for stably expressed genes that met all of the following criteria: *i*. p>0.5 for any pair-wise comparison of groups, *ii*. lowest coefficient of variation among all samples, *iii*. good annotation of the gene, *iv*. consistent high expression levels. After further screening of available literature on potential reference genes ("housekeeping genes"), we selected in a pilot phase the following genes from the stably expressed genes for analysis: *ACTB*, *B2M*, *HPRT*, *PGK1*, *RPS14*, and *RPS27*. We found most stable expression for *B2M*, which in addition showed a lower coefficient of variation than *ACTB*, recently suggested to be a less-than-ideal housekeeping gene depending on the cellular context¹⁴ ¹⁵. To rule out any inconsistency in the use of the reference gene, we opted to use *B2M* and compared the qPCR expression data of cohort II to normalization against *ACTB*, which yielded similar results (Supplementary Figure 3a and data not shown).

Supplementary Note 2

We assessed the annotated biological function of the 23 genes comprising the final diagnostic signature, as well as their putative role in monocyte function and/or phenotype. An overview can be found in Supplementary Table 5. A pathway analysis by Ingenuity Pathway Analysis (www.ingenuity.com) revealed that top pathways and functions included acute phase response signalling, free radical scavenging, immune cell trafficking, inflammatory disease, and cell death and survival. Taking those 7 genes upregulated in the in vitro model system, their annotated function suggests that immune signals may be the underlying mechanism in driving their expression

shift. However, we could not identify key regulators of known pathways, probably due to the limited information on reciprocal effects of PBM and tumour cells¹⁶. Though of high interest with regards to the biological function, functional biological knowledge is dispensable to exploit the full potential of the gene signature as a diagnostic tool in analogy to other important clinical tests, which are devoid of a biological understanding (e.g., prostate specific antigen, PSA, and pro-calcitonin, PCT).

Supplementary Note 3

In accordance with our initial screening results, we found no differences in expression patterns of P versus PM (data not shown). Moreover, as cumulating evidence is suggesting subcategories of CRC according to its location¹⁷, we investigated if the gene signature was capable of separating left versus right CRC or colon versus rectal cancer, respectively. In line with the homogeneous clustering of samples, we found no differences by location (AUC of 0.45 [0.20-0.73] for left versus right CRC and AUC of 0.47 [0.28-0.70] for colon versus rectal cancer).

Supplementary Note 4

We sought to confirm our findings from the screening in independent samples by independent techniques to rule out bias by the chosen technique and maximize chances of extrapolation to other clinical centres. Our first step was a random reprocessing of collected samples and assessment by qPCR, which led to an initial refinement of the gene signature, while some genes in this subset of samples performed well even as single markers. By assessing Spearman correlation values between expression data in the Illumina platform (used for screening) and the qPCR technique (used for confirmation), we could rule out discrepancies in expression

between both analyses (Supplementary Figure 8). Consistently, a multicentric <text> validation trial revealed that the established gene signature retained the promising

Gut



Supplementary Figure 1: Isolation of PBM and monocyte counts

a, Quality control of PBM isolation procedure in the pilot phase: FACS staining as histogram for CD14 (FITC). Comparison of the CD14 flow-through (left) and the CD14⁺ purified monocytes (right). **b**, Representative hemocytometric assessment of PBM purity, which was performed for each individual sample. **c**, Monocyte counts in whole blood were not different between (P,PM) and HV, neither relative (left), nor absolute (right).



Supplementary Figure 2: Differentially expressed genes in PBM

a, b, Differentially expressed genes in groupwise comparison of P, PM, and HV. The MA plots (a) show the fold change versus the average expression intensity, while the Volcano plots (b) show fold change in relation to the p values. Green, significantly downregulated genes; red, significantly upregulated genes; corrected p<0.05.

A. The N. in relativ

https://mc.manuscriptcentral.com/gut



Supplementary Figure 3: Technical validation (subset of cohort I)

÷

HV HV

P, PM

8-

ا الله 🗐

TNPO1

Ē

TNF

Ē 🗖

a, Comparative dot plot of raw CT values in qPCR for *ACTB* and *B2M*, revealing that the distribution is similar for both genes, and box-and-whiskers plot comparing normalization against both reference genes. **b**, Expression levels of all 43 putative candidates identified by genome-wide screening and assessed by qPCR. Expression levels are displayed as expression relative to the HV mean; boxes, first to third quartile; Whiskers, range; dots, values outside 1.5-times the interquartile distance; horizontal line, median; +, mean; *, p<0.1; **, p<0.01; ***, p<0.01.

🖻 Ŧ

CD68

40-

20- _____

0-

0-

8-

0-

2-

<u>م</u>ل

₽÷

÷

🖳 😐

TAF15

RLPL2

i l

IFR2

ē 🖻

GPER

0-

i 🖻 🖶

ткт

÷



Supplementary Figure 4: Gene expression levels of non-confirmed candidates in the multicentric validation (cohort II)

Expression levels are displayed as expression relative to the HV mean; boxes, first to third quartile; Whiskers, range; dots, values outside 1.5-times the interquartile distance; horizontal line, median; +, mean; *, p<0.05; **, p<0.01; ***, p<0.001.



Supplementary Figure 5: The gene signature stays robust over disease progression (cohort II)

Multicentric validation of the finding that the gene signature cannot discriminate between P and PM. Expression levels are displayed as expression relative to the HV mean; boxes, first to third quartile; Whiskers, range; dots, values outside 1.5-times the interquartile distance; horizontal line, median; +, mean; *, p<0.05; **, p<0.01; ***, p<0.001.

iminate μ μtile; Whiskei, μ<0.01; ***, p<0.0μ

Page 135 of 149



ROC analyses for each individual in cohort II. AUC, area under the curve.

1.0

1.0

1.0

1.0

https://mc.manuscriptcentral.com/gut

Gut





Supplementary Figure 7: Identification of putative markers in the *in vitro* model

Shown are the expression levels of the 16 genes not selected out of the gene signature, which show alterated expression levels in culture without any stimulus. Expression levels are shown as mean with SEM at 18 hours and 72 hours later (90 hours). *, p<0.05; **, p<0.01; ***, p<0.001; ****, p<0.0001; n.e., not expressed *in vitro*.

https://mc.manuscriptcentral.com/gut

Page 137 of 149



Gut

Scatter plots of Cohort I displaying correlation between Illumina microarray (x axis) and qPCR data (y axis). Spearman correlation values and p values are noted in the figures.





Page 139 of 149







Supplementary Figure 8 – continued

SUPPLEMENTARY TABLES

Supplementary Table 1: IDT PrimeTime qPCR Assays

Gene Name	Assay ID
ACP5	Hs.PT.47.311649.g
ACTB	Hs.PT.47.227970.g
ADM	Hs.PT.47.59577.g
ALDH1A1	Hs.PT.47.4497955
APP	Hs.PT.47.3063778
ARPC1B	Hs.PT.47.18828860
B2M	Hs.PT.47.18818394
BAX	Hs.PT.47.18828862
CCR1	Hs.PT.47.18828864
CD68	Hs.PT.47.18828865
CTSZ	Hs.PT.47.18828866
CXCR4	Hs.PT.47.512220
DDIT4	Hs.PT.47.18828867
DNAJC7	Hs.PT.47.18828868
ENSA	Hs.PT.47.18828869
FCER1A	Hs.PT.47.18828870
FKBP5	Hs.PT.47.18828871
GPER	Hs.PT.47.18828872
HBA1 / HBA2	Hs.PT.47.18828873
HBB	Hs PT 47 18828874
HI A-DOA1	Hs PT 47 18828891
HI A-DRB4	Hs PT 47 18828875
HMOX1	Hs PT 47 18828876
HNRNPK	Hs PT 47 18828877
HP	Hs PT 47 18828878
HPRT1	Hs PT 47 1231226
	Ho DT 47 19929990
	Lo DT 47.10020000
	Lo DT 47 10020001
LAF 11014A	Lo DT 47.10020002
LOC100000509	ПS.F 1.47.10020003
	HS.P1.47.18828884
LOC100132394	HS.P1.47.18828885
LOC100170939	HS.P1.47.18828886
LOC644063	Hs.P1.47.18828888
LOC653888	Hs.PT.47.18828889
LOC723972	Hs.PT.47.18828890
PGK1	Hs.PT.47.18828893
RILPL2	Hs.PT.47.18828894
RPS14	Hs.PT.47.18828895
RPS27	Hs.PT.47.18828896
S100P	Hs.PT.47.18828897

2
3
4
5
c
0
7
8
à
10
10
11
12
10
13
14
15
16
17
17
18
19
20
20
21
22
23
24
24
25
26
27
20
28
29
30
31
51
32
33
34
25
35
36
37
38
200
39
40
41
42
40
43
44
45
16
40
47
48
<u>4</u> 0
50
50
51
52
52
55
54
55
56
57
57
58
59

60

	р
CP5	0.6760
ADM	1
ALDH1A1	0.2020
APP	1
ARPC1B	1
BAX	0.7569
CCR1	1
CD68	1
CTSZ	1
CXCR4	0.2020
DDIT4	
DNAJC7	1
ENSA	0.3600
-CER1A	1
FKBP5	0.3600
GPER	0.2401
HBA1	0.2411
HBB	1
HLA-DQ1	0.1095
HLA-DRB4	1
HMOX1	1
HNRNPK	0.6160
HP	0.6160
IER2	0.8236
IL1R2	0.1773
LAPTM4A	0.2651
LOC100008589	1
LOC100170939	1
LOC643888	1
LOC644063	0.0147
_OC723972	0.0552
RLPL2	0.4941
RN28S1	1
S100P	1
SDHC	0.6160
SEPT5	1
SLC39A1	1
SLPI	0.8103
SOCS3	1
TAF15	0.3600
ТКТ	0.1162
TNF	0.5485
TNPO1	0.6160
Gut

Supplementary Table 3: Overview of development of a validated gene signature from putative candidates

Genomewide Screening										ation and Valio	lation
	P,PM vs HV ^a		Ρv	′s HV			PM vs HV			P,PM vs. HV	
	Ratio	р		Ratio	р		Ratio	р		Ratio	р
ADM	2,00	<0,0001	ADM	1,75	0,0059	ADM	2,27	<0,0001	<i>ACP5</i> ^₀	1,61	<0,0001
ALDH1A1	0,66	0,0002	CTSZ	1,76	0,0103	ALDH1A1	0,56	<0,0001	ADM	2,16	<0,0001
ARPC1B	1,55	0,0209	DDIT4	1,78	0,0226	AQP9	1,62	<0,0001	ALDH1A1	0,88	<0,0001
BAX	1,50	0,0001	DNAJC7	1,59	0,0005	BAX	1,52	0,0008	APP	1,61	<0,0001
CTSZ	1,79	0,0007	FCER1A	0,61	0,0296	CTSZ	1,81	0,0056	ARPC1B	0,98	0,6497
DDIT4	1,71	0,0063	HBA1	3,51	0,0008	DDIT4	1,65	0,0477	BAX	1,76	<0,0001
DNAJC7	1,59	<0,0001	HBA2	4,31	0,0004	DNAJC7	1,60	0,0004	CCR1	0,90	0,3981
FCER1A	0,52	0,0002	HBB	3,95	0,0004	DYSF	1,52	0,0002	CD68	1,76	<0,0001
FKBP5	1,61	0,0001	HMOX1	1,55	0,0017	FCER1A	0,45	0,0001	CTSZ	1,96	<0,0001
GPER	1,58	0,0006	HNRNPK	1,60	0,0497	FCGR1A	1,52	0,0003	CXCR4	2,24	<0,0001
HBA1	2,33	0,0078	HS.143909	1,56	<0,0001	FKBP5	1,85	<0,0001	DDIT4	1,47	0,0025
HBA2	2,69	0,0051	HS.581828	1,52	<0,0001	GPER	1,78	0,0001	DNAJC7	1,07	0,1045
HBB	2,39	0,0099	HS.61208	1,65	<0,0001	HLA-DRB6	0,42	0,0102	ENSA	0,89	0,1122
HMOX1	1,54	0,0001	IER3	1,50	0,0009	HMOX1	1,53	0,0020	FCER1A	0,97	0,7541
HNRNPK	1,58	0,0125	LOC100008589	1,68	0,0131	HP	1,75	0,0080	FKBP5	2,45	<0,0001
HP	1,54	0,0131	LOC100128274	0,66	0,0195	HS.61208	1,56	<0,0001	GPER	5,29	<0,0001
HS.143909	1,51	<0,0001	LOC100130707	1,51	0,0232	LOC100170	939 1,65	0,0001	HBA1	15,07	0,0165
HS.61208	1,60	<0,0001	LOC100132394	1,79	0,0095	LOC100190	986 1,53	0,0001	HBB	11,96	0,0281
IL1R2	1,50	0,0482	LOC100132727	0,66	0,0282	LOC153561	1,73	0,0001	HLA-DQA1	1,01	0,8425
LOC10008	<i>1,55</i> 1 ,55	0,0079	LOC100134364	1,57	0,0057	LOC441087	1,54	0,0177	HLA-DRB4	0,77	0,3931
LOC100129	<i>1,71</i> 685	0,0356	LOC153561	1,50	0,0049	RNF146	1,50	0,0002	HMOX1	0,95	0,7338
LOC100132	2394 1,65	0,0045	LOC649143	1,90	0,0133	S100P	1,75	0,0007	HNRNPK	0,92	0,2280
LOC100134	364 1,53	0,0009	LOC723972	1,51	0,0001	SEPT5	1,59	0,0347	HP	1,92	<0,0001
LOC100170	939 1,54	<0,0001	LOC728755	0,64	0,0210	SLC39A1	1,54	0,0025	IER2	0,97	0,8782
LOC153561	1,61	<0,0001	SLC39A1	1,50	0,0058	SOCS3	1,73	0,0014	IL1R2	0,86	0,4209

			Genon	newide Screeni	ing				Confirmation	and Valid	lation
P,PM	l vs HV ^a		P vs HV			PM vs HV			P,PM vs. HV		
	Ratio	р		Ratio	р		Ratio	р		Ratio	р
LOC649143	1,56	0,0356	TAF15	2,06	0,0001	TAF15	1,73	0,0002	LAPTM4A	1,59	<0,00
LOC653156	1,73	0,0443	ΤΚΤ	1,58	0,0034	TKT	1,55	0,0048	LOC100008589	0,99	0,93
LOC653737	1,86	0,0472	ZNF223	0,66	0,0478	TNPO1	1,55	0,0001	LOC100170939	1,09	0,00
LOC728755	0,66	0,0066				UPP1	1,58	<0,0001	LOC643888	1,05	0,26
S100P	1,53	0,0020	S'A			ZBTB16	1,52	0,0252	LOC644063	1,54	<0,00
SEPT5	1,57	0,0094							LOC723972	1,03	0,19
SLC39A1	1,52	0,0003							RLPL2	0,95	0,36
SOCS3	1,51	0,0043							RN28S1	1,03	0,30
TAF15	1,76	<0,0001							S100P	3,35	<0,00
ТКТ	1,56	0,0003							SDHC	1,04	0,30
									SEPT5	3,47	0,00
									SLC39A1	1,02	0,38
									SLPI	15,76	0,00
									SOCS3	1,60	0,01
									TAF15	0,84	0,01
									ткт	1,79	<0,00
									TNF	0,75	0,02
									TNPO1	1.02	0.31

t the identified 40 probes correspond to 35 genes, as several probes may exist for one gene. See Supplementary methods for details on gene numbers.

Genes confirmed by qPCR are shown in bold print (23 genes).

Sup	op	lementary	/ Table 4	4: Co	onfirma	tion in	random	subset	of	cohort	
-----	----	-----------	-----------	-------	---------	---------	--------	--------	----	--------	--

	Mean	Fold ratio ^b			
	expression ^a		р	AUC	
ACP5	81,336	1,73	0.0081 ^c	0.79	
ADM	73,107	4,23	0.0941	0.95	
ALDH1A1	47,649	0,99	0.9624	0.51	
APP	176,332	1,32	0.1576	0.73	
ARPC1B	1873,873	1,15	0.3336	0.57	
BAX	11,474	1,29	0.0978	0.67	
CCR1	338,797	1,01	0.9292	0.52	
CD68	1821,912	1,27	0.0640	0.76	
CTSZ	1580,418	1,28	0.0637	0.76	
CXCR4	508,754	1,52	0.0065	0.84	
DDIT4	36,963	4,34	0.0010	0.96	
DNAJC7	105,494	0,92	0.5431	0.62	
NSA	250,287	1,21	0.2580	0.67	
CER1A	410,118	0,54	0.0768	0.73	
-KBP5	39,131	2.08	0.0013	0.89	
PER	1.875	7.59	0.0138	0.93	
IBA1	5243.339	41.40	0.0861	0.86	
IBB	440.188	31.10	0.0773	0.85	
LA-DQ1	1748.345	0.53	0.0918	0.77	
LA-DRB4	1135.072	0.71	0.6316	0.53	
IMOX1	405.989	1.30	0.0729	0.70	
NRNPK	1648.472	1.05	0.7356	0.52	
1P	129.478	1.50	0.0218	0.76	
- 	0.657	0.65	0.2556	0.63	
1R2	4 794	5,85	0.0288	0.87	
PTM4A	338,391	1.35	0.0206	0.78	
OC100008589	18500877 250	1.12	0.5782	0.63	
OC100170939	252,768	0.96	0.8506	0.56	
OC643888	308 104	1 46	0.6143	0.55	
OC644063	1504 972	1.07	0.0415	0.71	
C723972	568 957	1 00	0.9645	0.56	
1 PI 2	186 476	1 02	0.9078	0.53	
IN2851	14924567 167	0.92	0.5908	0.58	
3100P	8 494	2,90	0.0003	0.91	
SDHC	313 373	1 02	0.9145	0.53	
SEPT5	0 298	1 22	0.6497	0.50	
LC39A1	166 812	1 12	0.0407 0.4069	0.62	
	1 656	5 39	0.4003	0.02	
20053	128 80/	3.35	0.2477	0.71	
νουσσ ΓΔΕ15	120,034	0,00	0.0001	0.51	
	026 111	0,00	0.0004	0.00	
	200,111	0.04	0.0001	0.02	
	20,000	0,94	0.7323	0.55	
	140,447	1,00	0.9122	0.52	

^aMean expression of gene of interest / 10,000 copies of *B2M*

^bFold ratio of patients compared to healthy volunteers ^cBold print indicates where cutoff criteria (p<0.1, AUC>0.7) are met. See main manuscript and Supplementary methods for more detailed information technical.com/gut

Gut

		· · · · · · · · · · · · · · · · · · ·	
Gene	Full Name	Biological Function	Potential Function in Monocytes
ACP5	acid phosphatase 5.	iron containing glycoprotein	negative regulation of inflammatory
	tartrate resistant	involved in adhesion and	response in interleukin pathways
		<u>migration</u>	
<u>ADM</u>	<u>adrenomedullin</u>	vasodilation, regulation of	antimicrobial activity, wound healing
		hormone secretion,	

Supplementary Table 5: Identity and Function of the gene signature members

		hormone secretion,	
		promotion of angiogenesis	
APP	amyloid beta (A4)	protein basis of amyloid	antimicrobial activity, mitotic activity
	precursor protein	plaques in Alzheimer disease	
BAX	BCL2-associated X	p53-mediated activator of	mveloid cell homeostasis
	protein	apoptosis	
CD68	CD68 molecule	integral membran	highly expressed on monocytes and
0200		alycoprotein of scavenger	macrophages mediator of recruitment
		recentor family	and activation
CTSZ	cathenein 7	lysosomal cystein	unknown
0102		proteinase involved in	dikilown
		migration and adhesion	
CVCD4	chamaking (C X C matif)	<u>CVC characting receptor</u>	modiator of rear uitment, chemotovia
		<u>CAC chemokine receptor</u>	mediator of recruitment, chemotaxis,
	receptor 4	specific for stromal cell-	and activation
		derived factor-1	
<u>DDI14</u>	DNA-damage-inducible	negative regulation of mIOR	detense response to microbial signals
	transcript 4	<u>signalling upon cellular</u>	
		stress	
<u>FCER1A</u>	Fc fragment of IgE, high	alpha subunit of IgE-	positive regulation of type-I immune
	affinity I, receptor for;	mediated allergic response	response and macrophage
	alpha polypeptide		<u>differentiation</u>
<u>FKBP5</u>	FK506 binding protein 5	member of immunophilin	receptor for FK506 and rapamycin,
		protein family,	mediating calcineurin inhibition
		immunoregulation	
GPER	G protein-coupled	non-genomic signalling of	negative regulator of leukocyte
	estrogen receptor 1	estrogen stimulus	activation: innate immune response
HBA1	hemoglobin, alpha 1	alpha chain of hemoglobin	unknown
<u></u>	<u></u>	HbA	
HBB	hemoglobin beta	beta chain of hemoglobin	positive regulation of nitric oxide
<u></u>	<u>Homogroom, bota</u>	HbA	synthesis
	major histocompatibility	MHC class II recentor	antigen processing and presentation
	complex class IL DO	activity: pentide antigen	anigen proceeding and procentation
	alpha 1	binding	
	homo oxygonaso	homo catabolism	regulation of phagoevtosis and
	(doovoling) 1	neme catabolism	migration chomoking cynthosis and
	(decycling) T		hading and angiaganasia
	bontoglabin	proproprotoin of bonto alabia	nealing, and angiogenesis
<u> </u>	naplogiobin	preproprotein of naploglobin	acute-phase delense response
		subunit	
<u>IL1R2</u>	interleukin 1 receptor,	cytokine receptor for IL-1	cytokine-mediated immune response
	type II		
<u>LAPTM4A</u>	lysosomal protein	<u>unknown</u>	unknown
	transmembrane 4 alpha		
<u>LOC644063</u>	<u>heterogeneous nuclear</u>	<u>unknown</u>	unknown
	<u>ribonucleoprotein K</u>		
	<u>pseudogene 4</u>		
<u>S100P</u>	S100 calcium binding	cell cycle progression and	<u>unknown</u>
	<u>protein P</u>	<u>differentiation</u>	
<u>SLPI</u>	secretory leukocyte	secreted inhibitor of serin	negative regulation of endopeptidase
	peptidase inhibitor	proteinases	activity
<u>SOCS3</u>	suppressor of cytokine	negative regulator of	modulator of immune response,
	signaling 3	cytokine signalling	particularly IFN-y mediated
ΤΚΤ	transketolase	enzyme of pentose	metabolic modulator
		phosphate pathway	
		<u> </u>	

https://mc.manuscriptcentral.com/gut

Gut

SUPPLEMENTARY REFERENCES

- 1. Weitz J, Koch M, Debus J, et al. Colorectal cancer. Lancet 2005;**365**(9454):153-65.
- Nyugen J, Agrawal S, Gollapudi S, et al. Impaired functions of peripheral blood monocyte subpopulations in aged humans. Journal of clinical immunology 2010;**30**(6):806-13.
- Du P, Kibbe WA, Lin SM. lumi: a pipeline for processing Illumina microarray. Bioinformatics 2008;24(13):1547-8.
- 4. Lin SM, Du P, Huber W, et al. Model-based variance-stabilizing transformation for Illumina microarray data. Nucleic acids research 2008;**36**(2):e11.
- 5. Smyth GK. Linear models and empirical bayes methods for assessing differential expression in microarray experiments. Statistical applications in genetics and molecular biology 2004;**3**:Article3.
- 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(1):289-300.
- 7. Sample size for microarray experiments. Secondary Sample size for microarray experiments. <u>http://bioinformatics.mdanderson.org/MicroarraySampleSize/</u>.
- 8. Dietterich TG. Ensemble methods in machine learning. Lecture Notes in Computer Science 2000;**1857**:1-15.

9. Impute: Imputation for microarray data. [program]. 1.32.0 version, 2013.

- 10. Robin X, Turck N, Hainard A, et al. pROC: an open-source package for R and S+ to analyze and compare ROC curves. BMC bioinformatics 2011;**12**:77.
- Burges CJC. A Tutorial on Support Vector Machines for Pattern Recognition.
 Data Min Knowl Discov 1998;2(2):121-67.

- Liaw A, Wiener M. Classification and Regression by randomForest. R News: The Newsletter of the R Project 2002;2(3):18-22.
 - 13. Ambroise C, McLachlan GJ. Selection bias in gene extraction on the basis of microarray gene-expression data. Proceedings of the National Academy of Sciences of the United States of America 2002;**99**(10):6562-6.
- Piehler A, Grimholt R, Ovstebo R, et al. Gene expression results in lipopolysaccharide-stimulated monocytes depend significantly on the choice of reference genes. BMC Immunology 2010;11(1):21.
- 15. Guo C, Liu S, Wang J, et al. ACTB in cancer. Clinica chimica acta; international journal of clinical chemistry 2013;**417**:39-44.
- 16. Khatri P, Sirota M, Butte AJ. Ten Years of Pathway Analysis: Current Approaches and Outstanding Challenges. PLoS Comput Biol 2012;8(2):e1002375.
- 17. Jess P, Hansen IO, Gamborg M, et al. A nationwide Danish cohort study challenging the categorisation into right-sided and left-sided colon cancer. BMJ open 2013;3(5).