

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 protection by the applicable law.

(Article begins on next page)

1
2
3 **Tumour-Educated Circulating Monocytes are Powerful Candidate**
4
5 **Biomarkers for Diagnosis and Disease Follow-up of Colorectal Cancer**
6
7

8 Alexander Hamm^{1,2#}, Hans Prenen^{3#}, Wouter Van Delm^{4#}, Mario Di Matteo^{1,2}, Mathias Wenes^{1,2},
9 Estelle Delamarre^{1,2}, Thomas Schmidt⁵, Jürgen Weitz^{5,6}, Roberta Sarmiento⁷, Angelo Dezi⁷,
10 Giampietro Gasparini⁷, Françoise Rothé⁸, Robin Schmitz⁵, André D'Hoore⁹, Hannes Iserentant¹⁰,
11
12
13
14 Alain Hendlisz⁸ & Massimiliano Mazzone^{1,2}
15

16
17 ¹Lab of Molecular Oncology and Angiogenesis, Vesalius Research Center, VIB, Leuven, Belgium

18
19 ²Lab of Molecular Oncology and Angiogenesis, Vesalius Research Center, Department of Oncology,
20
21 KU Leuven, Leuven, Belgium

22
23 ³Digestive Oncology, University Hospitals Leuven and Department of Oncology, KU Leuven, Leuven,
24
25 Belgium

26
27 ⁴Nucleomics Core, VIB, Leuven, Belgium

28
29 ⁵Department of General, Visceral, and Transplantation Surgery, University of Heidelberg, Heidelberg,
30
31 Germany

32
33 ⁶Department of Visceral, Thoracic, and Vascular Surgery, University Hospital Carl Gustav Carus,
34
35 Technical University Dresden, Dresden, Germany

36
37 ⁷Department of Oncology, San Filippo Neri, Rome, Italy

38
39 ⁸Medical Oncology Clinic, Institut Jules Bordet, Brussels, Belgium

40
41 ⁹Department of Abdominal Surgery, University Hospitals Leuven, KU Leuven, Leuven, Belgium

42
43 ¹⁰VIB, Zwijnaarde, Belgium

44
45 #contributed equally to this study

46
47 **Correspondence:**

48
49 Massimiliano Mazzone, massimiliano.mazzone@vib-kuleuven.be,
50
51 Tel: +32-16-373213, Fax +32-16-372585
52
53 VIB Vesalius Research Center, KU Leuven, Herestraat 49, Bus 912, 3000 Leuven, Belgium

54
55 Hans Prenen, hans.prenen@uzleuven.be, Tel: +32-16-340238

56
57 Word Count: 4127

58
59 Key Words: Monocytes, colorectal cancer, screening, inflammation
60

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
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
P	non-metastatic CRC patient
P, PM	non-metastatic and metastatic CRC patients
PC	pancreatic cancer patient
PG	gastric cancer patient
PGT	gastritis patient
PM	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.

1
2
3 Conclusion: Our study is the first to demonstrate monocyte plasticity in response to
4
5 tumour-released soluble factors. The identified distinct signature in tumour-educated
6
7 monocytes might be used as candidate biomarker in CRC diagnosis and harbours
8
9 the potential for disease follow-up and therapeutic monitoring.
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
51
52
53
54
55
56
57
58
59
60

Confidential: For Review Only

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 gene signature may bring additive value to the established screening tools for CRC and early detection of recurrent disease, both offering patients better chances of cure. Moreover, the plasticity of monocytes may prove to be ideal for real-time follow-up of CRC treatment.

Confidential: For Review Only

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 type^{10 11}. 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.

1
2
3 Several novel accessible diagnostic tools share the major opportunity to make
4 frequent screening more appealing to a greater number of patients, as a less
5 invasive method is likely to increase compliance and allow for decreased screening
6 intervals (recently comprehensively reviewed⁶). While conventional blood-based
7 tumour markers (particularly carcino-embryonic antigen, CEA¹⁶) have been
8 established as supplemental markers in treatment monitoring, they have failed to
9 yield high diagnostic accuracy as primary screening tools. In addition to the
10 established FIT or FOBT⁵, other potential diagnostic markers include serum-
11 associated biomarkers (e.g. circulating tumour DNA¹⁷, micro-RNA¹⁸, methylation
12 markers like *SEPT9*¹⁹), genetic marker sets in white blood cells²⁰⁻²³, and, most
13 recently, fecal tumour DNA²⁴. However, all of these approaches display limited
14 sensitivity and specificity⁶. In this study, we therefore assess the sensitivity and
15 specificity of a novel gene signature in circulating monocytes for the diagnosis of
16 CRC in comparison to healthy individuals or to other cancer types, and assess its
17 robustness in prospectively obtained samples.
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
51
52
53
54
55
56
57
58
59
60

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

1
2
3 quantitative RT-PCR (qPCR) on the 7500Fast System (Applied Biosystems) using
4
5 intron-spanning PrimeTime qPCR Assays (Integrated DNA Technologies) listed in
6
7 Supplementary Table 1 as described in Supplementary Methods. For statistical
8
9 analysis, we followed a three-step top-down approach to construct a gene signature
10
11 for CRC, with details explained in Supplementary Methods.
12
13

14 15 16 **Multicentric validation study**

17
18 For validation of a diagnostic test, we used cohort II to train and validate a multi-gene
19
20 classifier. Splits in training and test sets for validation were performed by stratified
21
22 random sampling for centre of origin and class label as detailed out in Supplementary
23
24 Methods. Samples with missing values for more than 25% of the genes were
25
26 excluded from the analysis. We ruled out an effect of the class labeling on the
27
28 percentage of missing values with Fisher's exact test (Supplementary Table 2).
29
30

31
32 The training dataset was used to build three types of classifiers: a support vector
33
34 machine (SVM)²⁹ with linear kernel, a single-gene majority vote (SGMV) classifier,
35
36 and a random forest classifier (RF³⁰). Subsequently, we applied an ensemble
37
38 method³¹ that votes according to the majority of the three independent classifiers.
39
40 Performance was validated both with ranking (AUC) and classification (balanced
41
42 error rate, BER, Se, Sp) scores with 95% confidence intervals ([lower boundary;
43
44 upper boundary]). We explicitly opted for relatively simple computational models in
45
46 order to limit chances of over-fitting the training data and to maximize interpretability
47
48 of the models' internal decision-making process. Model flexibility was further
49
50 controlled through a Monte-Carlo cross-validation scheme (MCCV)³², before final
51
52 estimation of the model parameters. Validation of the predictive models was done on
53
54
55
56
57
58
59
60

1
2
3 the test set of cohort II, which were not included during development of the models.

4
5 Details on all classification methods are specified in Supplementary Methods.

6
7 In order to avoid biased conclusions, the analysis of the 23 genes was
8
9 complemented with a study by an independent team (DNAlytics, Belgium) that
10
11 adopted a slightly modified analysis protocol (see Supplementary Methods). All
12
13 complementary analyses were performed in R with scripts designed by DNAlytics,
14
15 fully independently from other analyses described in this paper.
16
17
18
19

20 21 **In vitro model system**

22
23 To study the effects of tumour-released soluble factors on gene expression in
24
25 monocytes, we established an *in vitro* model system, where monocytes from healthy
26
27 donors were challenged with tumour-released soluble factors and changes in gene
28
29 expression profile were analyzed by qPCR. See Supplementary Methods for details.
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

RESULTS

Establishment of putative biomarkers by genome-wide expression analysis

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.

TABLE 1: CLINICOPATHOLOGICAL CHARACTERISTICS OF PATIENTS AND HEALTHY VOLUNTEERS

Cohort	I		II								III		V	VI	
	P,PM	HV	P,PM				HV				P,PM	HV	P	P,PM	HV
			LEU ^a	HD ^b	SFN ^c	IJB ^d	LEU ^a	HD ^b	SFN ^c	IJB ^d					
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	/	16	3	2	2	/	/	/	/	16	/	0	6	/
non-metastatic	27	/	23	16	8	3	/	/	/	/	11	/	15	11	/
UICC stage															
1	3	/	7	2	1	1	/	/	/	/	2	/	4	2	/
2	12	/	8	8	2	0	/	/	/	/	3	/	7	7	/
3	12	/	8	6	5	2	/	/	/	/	6	/	4	2	/
4	28	/	16	3	2	2	/	/	/	/	16	/	0	6	/
Tumour localization															
Caecum	5	/	3	3	0	0	/	/	/	/	2	/	1	2	/
Ascendens	11	/	4	3	1	0	/	/	/	/	6	/	4	4	/
Transversum	0	/	4	3	0	0	/	/	/	/	2	/	1	0	/
Descendens	4	/	2	1	3	3	/	/	/	/	0	/	0	1	/
Sigmoid	28	/	15	3	2	0	/	/	/	/	10	/	5	6	/
Rectum	6	/	8	5	3	2	/	/	/	/	7	/	4	2	/
Double	1	/	3	1	1	0	/	/	/	/	0	/	0	2	/

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

1
2
3
4 The purity of the monocyte fraction was >90%, as assessed by FACS analysis in the
5
6 pilot phase (Supplementary Figure 1a) and verified by hemocytometric analysis for
7
8 each individual sample (Supplementary Figure 1b). Both absolute and relative
9
10 monocyte counts were not different between patients and healthy volunteers
11
12 (Supplementary Figure 1c). We therefore investigated differentially expressed genes
13
14 by genome-wide expression analysis using the Illumina HumanHT-12 v4 Expression
15
16 BeadChip Kit. The data discussed in this publication have been deposited in NCBI's
17
18 Gene Expression Omnibus³³ and are accessible through GEO Series accession
19
20 number
21
22 GSE47756
23
24 ([http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?token=hvmpvoswuqaeybc&acc=GSE](http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?token=hvmpvoswuqaeybc&acc=GSE47756)
25
26 [47756](http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?token=hvmpvoswuqaeybc&acc=GSE47756)). In first instance, we compared the average expression values of all CRC
27
28 patients (P,PM), comprised of non-metastatic (P) and metastatic (PM) patients, to
29
30 that of healthy volunteers (HV). The resulting gene signature of (P,PM) versus HV
31
32 consisted of 36 upregulated and 4 downregulated probes (Figure 2a, b, Table 2). In
33
34 second instance, we were interested if the gene signature in patients with
35
36 synchronous metastases *i.e.*, at the time of diagnosis (PM, n=28) was different from
37
38 that in non-metastatic patients (P, n=27). Interestingly, the number of up- and down-
39
40 regulated genes was comparable in both P and PM (in comparison to HV) (Table 2
41
42 and Supplementary Figure 2a, b), while there were no genes found to be differentially
43
44 expressed between the two patient groups (Supplementary Figure 2a, b), indicating
45
46 that the gene signature induced at early onset stays robust over disease progression.
47
48 Indeed, when post-hoc assessing those samples from patients with early stages (Tis
49
50 and T1), they clustered with the rest of the patient samples (data not shown). A
51
52 power analysis revealed that, for the given number of genes, samples and observed
53
54 variation, chances were very low ($<10^{-10}$) that truly differentially expressed genes with
55
56
57
58
59
60

1
2
3 fold changes larger than 1.5 had been missed. Therefore, adding more samples
4
5 would probably have changed little to the panel of candidate genes that our screen
6
7 returned.
8
9

10 11 **Confirmation of the gene signature in independently processed samples**

12
13 To validate the genetic signature, we performed quantitative RT-PCR (qPCR)
14
15 analysis on a random subset of PBM from 8 samples of each of the three groups (P,
16
17 PM, and HV), normalizing to reference gene *B2M*, which was selected after an
18
19 extensive screening procedure (Supplementary Note 1). To avoid bias in the
20
21 confirmation procedure, we freshly extracted RNA from independently stored
22
23 samples for confirmative expression analysis. In analyzing 43 putative marker genes
24
25 with probes listed in Supplementary Table 1, 23 genes showed differential
26
27 expression between (P, PM) and HV (Supplementary Figure 3b, Table 2, and
28
29 Supplementary Table 4). Thus, we were able to confirm a subset of the previously
30
31 established gene signature, independent of the RNA extraction and the platform used
32
33 for expression analysis. Information on the annotated biological function of the genes
34
35 of the diagnostic signature can be found in Supplementary Table 5 and
36
37 Supplementary Note 2.
38
39
40
41
42
43
44

45 46 **Confirmation of the gene signature in a multicentric validation set**

47
48 For a rigorous validation of the gene signature, we collected an independent
49
50 multicentric validation set (cohort II) from a total of 4 different European oncological
51
52 centres with stratified training and test sets as described in Supplementary Methods.
53
54 Using the panel of 23 genes confirmed previously, we found consistently differential
55
56 expression between all patients and the healthy volunteers (Figure 2c and
57
58
59
60

1
2
3 Supplementary Figure 4). In line with the findings from the screening phase, there
4
5 were no detectable differences in expression levels between P and PM
6
7 (Supplementary Figure 5), while either patient group alone compared to HV was
8
9 differentially expressed (data not shown).
10

11
12 In ROC analysis for single genes, we found that some, but not all of the genes that
13
14 displayed significantly differential expression were able to discriminate patient
15
16 samples from healthy individual samples with acceptable AUCs (Supplementary
17
18 Figure 6 and data not shown). We therefore hypothesized that a marker panel
19
20 consisting of multiple genes might yield better results in discriminating sample
21
22 identity. To address this question, we decided to test three different classification
23
24 algorithms on this data set, namely a support vector machine (SVM)²⁹ with linear
25
26 kernel, a single-gene majority vote (SGMV) classifier, a random forest classifier
27
28 (RF³⁰), and a combined classification by an ensemble method³¹, using the outcome
29
30 of the three classification algorithms for a final diagnostic decision. To limit over-
31
32 estimation of the performance by the particular training and test set, we performed a
33
34 MCCV as a conservative estimate with 1,000 cross-validations. Performance of all
35
36 classification algorithms in cohorts II – VI, including the conservative estimate of the
37
38 MCCV in cohort II, is given in detail in Table 2.
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

TABLE 2: PERFORMANCE SCORES OF MULTIGENE CLASSIFIER

	SGMV	SVM	RF	ENS
Cohort II (Validation)				
AUC [95% CI]	0.99 [0.99;1.00]	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]	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]
Cohort IV(pancreatic cancer)				
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]
Cohort IV (gastritis)				
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]
Cohort V (PR)				
Sensitivity [95% CI]	50.0 [20.0;80.0]	10.0 [0.0;30.0]	20.0 [0.0;50.0]	20.0 [0.0;50.0]
Cohort VI (PBMC)				
AUC [95% CI]	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

1
2
3 Strikingly, we achieved a remarkably high AUC of 0.99 [0.99;1.00] with a BER of
4
5 3.6% (Figure 2d and Table 2), translating into a sensitivity of 100.0%
6
7 [100.0%;100.0%] and a specificity of 92.9% [78.6%;100.0%] (Table 2). Neither of the
8
9 classification algorithms was capable of separating P from PM or detect differences
10
11 dependent on tumour localization (Supplementary Note 3).

12
13
14 In order to assess whether the diagnostic gene signature is actually suitable for
15
16 diagnosis of CRC in a screening setting, we have initiated a prospective sample
17
18 collection in both patients and healthy individuals who are subjected to colonoscopy.
19
20 In a pilot analysis in 27 patients (newly diagnosed with CRC by screening
21
22 colonoscopy) and 13 healthy individuals negative to screening colonoscopy (cohort
23
24 III), we found an AUC of 0.95 [0.85;1.00] with a BER of 7.6%, yielding a sensitivity of
25
26 92.6% [81.5%;100.0%] and a specificity of 92.3% [76.9%;100.0%] (Figure 2e and
27
28 Table 2). The complementary data analysis (independently performed by DNAlytics,
29
30 Belgium) on the same panel of 23 genes led to the matching conclusions in terms of
31
32 performance. The first experiment consisted in cross-validating a model on Cohort II
33
34 (BER: 8.4% [3.4%;13.4%]; AUC: 0.93 [0.88;0.98]). A second experiment consisted in
35
36 learning the same type of model on Cohort II and having it make predictions on
37
38 Cohort III (BER: 13.2%; AUC: 0.92).

45 **Soluble factors released by colorectal cancer cells induce an early, tumour** 46 47 **type-specific and reversible genetic fingerprint in monocytes**

48
49 We hypothesized that tumour-released soluble factors are the key players in inducing
50
51 the genetic signature in circulating monocytes. Thus, we established an *in vitro*
52
53 model system where we cultured freshly isolated human monocytes from healthy
54
55 donors in different conditions. In order to assess alterations in gene expression, we
56
57
58
59
60

1
2
3 first analyzed which of the 23 genes comprising the gene signature was up- or
4
5 downregulated in culture after 72 hours without any additional stimulus and excluded
6
7 these from the further *in vitro* studies (Supplementary Figure 7). Out of the remaining
8
9 gene signature, the majority (7/9) was specifically upregulated when culturing naïve
10
11 monocytes in medium conditioned by the CRC cell line HCT116, while expression
12
13 levels were not affected by mock medium (Figure 3a). Moreover, in line with the
14
15 coherent induction of the specific signature independent of the stage of the disease,
16
17 the induction *in vitro* was independent of hypoxic cues, as HCT116-conditioned
18
19 medium in hypoxia did not induce any different expression levels than medium
20
21 obtained in normoxia (Figure 3b). Likewise, the changes in expression levels of all
22
23 these genes occurred already 18 hours after stimulating monocytes with the
24
25 conditioned medium, consistent with the fact that already early stages are detectable
26
27 by the diagnostic signature.
28
29
30

31
32 To rule out an off-target effect of conditioned medium *i.e.*, unspecific cues from cell
33
34 metabolites, apoptotic bodies, pH, etc., we assessed the expression levels of the
35
36 genes upregulated by HCT116-conditioned medium in comparison to a benign colon
37
38 epithelium cell line, CCD 841 CoN (CCD), which did not induce alterations in gene
39
40 expression levels different from the Mock control (Figure 3a).
41
42

43
44 Prompted by this finding, we investigated if the induction of the genetic signature was
45
46 a general effect of malignant transformation or might be specific to the histotype of
47
48 cancer. To address this question, we conditioned medium with a gastric cancer cell
49
50 line, MKN-45 (MKN), to compare CRC to another frequent gastrointestinal solid
51
52 neoplasm. Remarkably, when comparing the expression levels in naïve monocytes
53
54 upon stimulation with the different conditioned media, we found that MKN-45
55
56
57
58
59
60

1
2
3 conditioned medium did not induce the same upregulation of the genes of interest as
4
5 HCT116 conditioned medium (Figure 3c).
6

7
8 As immune cells are highly versatile and plastic cells mirroring the microenvironment,
9
10 where they are embedded, we reasoned that the genetic signature induced by CRC
11
12 in monocytes might be dependent on the continuous presence of the stimulating
13
14 agents and thus be reversible upon inversion of the conditions. We therefore
15
16 incubated naïve monocytes first with HCT116-conditioned medium for 18 hours and
17
18 then refreshed the medium with plain culture medium, thus withdrawing the tumour-
19
20 released soluble factors. Strikingly, the previously elevated expression levels of a set
21
22 of marker genes were almost entirely reverted to the original (and to the mock
23
24 control) expression levels 72 hours after withdrawing the tumour-cell conditioned
25
26 medium (Figure 3d), whereas they remained constantly overexpressed when the
27
28 conditioned medium was maintained (data not shown).
29
30
31
32
33

34 **The monocyte signature is specific for CRC and might serve as a candidate** 35 **biomarker of disease follow-up** 36

37
38 Based on the *in vitro* results showing that the genetic signature is specific to CRC,
39
40 we sought to confirm these findings *in vivo*. We therefore assessed the diagnostic
41
42 signature in patients with i. cancer of the stomach and gastro-esophageal junction
43
44 (PG, n=15) and ii. pancreatic ductal adenocarcinoma (PC, n=16), two other frequent
45
46 cancers of the gastrointestinal tract¹. In addition, we analysed iii. patients with
47
48 gastritis (PGT, n=10) in order to compare the gene signature in CRC to a benign
49
50 inflammatory condition of the gastrointestinal tract (cohort IV). In line with the *in vitro*
51
52 results we saw that the vast majority of all genes were not significantly different
53
54 between either of the patient groups and healthy volunteers, indicating the specificity
55
56
57
58
59
60

1
2
3 of this monocyte imprinting by colorectal cancer cells (Figure 4a and data not
4 shown). Moreover, the classifier established to diagnose CRC could not separate
5 patients with gastric cancer (AUC 0.63 [0.48;0.77]), pancreatic cancer (AUC 0.41
6 [0.27;0.50]), or gastritis (AUC 0.52 [0.35;0.68]) from healthy individuals (Figure 4c, d
7 and Table 2).
8
9

10
11
12
13
14 The finding that the genetic signature is reverted upon withdrawal of the stimulating
15 agents prompted us to investigate in a pilot phase the behaviour of the entire
16 diagnostic signature in patients upon curative treatment *i.e.*, patients with surgically
17 removed tumours without any evidence of residual disease. To this end, we isolated
18 monocytes from 15 patients of stages I to III treated with curative intent (with or
19 without adjuvant treatment) and presenting at follow-up without detectable residual
20 disease (PR) (cohort V). Here, we found that virtually all of the previously
21 upregulated genes were reverted to expression levels comparable to those of healthy
22 volunteers (Figure 4b). Consequently, when applying the previously established
23 classifier, we found that it was able to distinguish accurately between patients in
24 remission and patients with tumour, while it could not detect differences between
25 patients in remission and healthy volunteers (Table 2).
26
27
28
29
30
31
32
33
34
35
36
37
38
39

40
41 Finally, as the plasticity of the signature offers the perspective to use the gene
42 signature for follow-up of treated patients, we became interested if the same
43 signature could be used to diagnose relapse (frequently as metachronous
44 metastases rather than local recurrence²). Although our dataset was not powered to
45 address this question with sufficient significance, we post-hoc identified four patients
46 from cohorts I and II included at presentation with metachronous metastases. All four
47 clustered clearly in the group of patients, separately from the healthy volunteers
48 (Figure 4e), suggesting that the signature might be used to detect disease relapse in
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 line with the previous results that show coherent expression over disease
4
5 progression.
6
7

8
9 **The gene signature is specific to monocytes in comparison to all peripheral**
10 **blood mononuclear cells (PBMCs)**
11
12

13
14 To rigorously assess if the genetic fingerprint identified in monocytes was specific to
15 this cell type or an epiphenomenon of genetic shifts in the entire population of
16 PBMCs, we isolated both monocytes and full PBMC fractions from 17 patients and 7
17
18 healthy volunteers for a comparative analysis (cohort VI). Interestingly, we found that
19
20 while in the monocyte population, the diagnostic marker set of 23 genes was
21
22 upregulated in all patients (both P and PM) in accordance with our previous results
23
24 (Figure 5a), there were no significant differences in the expression levels of the
25
26 analyzed genes in the full PBMC compartment when comparing patients to healthy
27
28 volunteers (Figure 5a). Consistently, applying the previously established classifier
29
30 with the defined cut-off values, it was impossible to separate the patient group from
31
32 the healthy volunteer group in PBMC (Figure 5b and Table 2), while the classifier
33
34 confirmed its accuracy in PBM (Figure 5c and Table 2). Thus, the differential
35
36 regulation of the gene signature in PBM used for CRC diagnosis is specific to the
37
38 monocytic lineage, reinforcing our initial working hypothesis that these cells are
39
40 specifically affected by tumour-secreted factors.
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

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 strategies^{2,6}. 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 biology-bound 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

1
2
3 most likely depends on the development of a one-step assay with capture of
4
5 monocytes from whole blood and gene expression analysis in a multiplex qPCR
6
7 assay with absolute quantification, avoiding extensive preanalytical processing steps.
8
9 However, the analytical reliability of this assay needs to be thoroughly established,
10
11 most likely requiring centralization of the analysis during the first phase of
12
13 distribution.
14

15
16 The finding that the specific gene signature is reversible if the stimulating cues are
17
18 withdrawn, was not only demonstrated *in vitro*, but also in a pilot analysis *in vivo* in
19
20 samples of patients after curative treatment. Although not completely unexpected in
21
22 view of the plastic nature of the monocyte-macrophage lineage, this analysis opens
23
24 avenues for treatment monitoring and companion diagnostics and will be assessed in
25
26 detail in a prospective study during patient follow-up.
27
28

29
30
31 If supported by further prospective validation studies, this gene signature may
32
33 outperform other published non-invasive test for CRC diagnosis⁶ (including single
34
35 surface markers in monocytes^{37 38}) or score similar to the most recent evaluation of
36
37 fecal tumour DNA²⁴. Moreover, we are the first to demonstrate that a potential
38
39 diagnostic biomarker obtained in patients at the time of primary diagnosis might also
40
41 be suitable for disease follow-up and thus assessment of treatment response, owing
42
43 to its high plasticity.
44
45

46
47 We acknowledge the limited conclusions that can be drawn from our case-control
48
49 study. Despite the confirmation in independent samples, we cannot fully exclude
50
51 possible confounders that can only be unveiled by a blinded, prospective sample
52
53 collection in screening individuals. These include, but are not limited to, the bias of
54
55 selecting patients that underwent colonoscopy for a clinical indication; the differences
56
57
58
59
60

1
2
3 in age, nutrition status, diet, and potentially lifestyle between patients and healthy
4
5 volunteers; the unblinded sample collection and processing. It is therefore of
6
7 paramount importance that a prospective validation study initiated by our group
8
9 includes screening individuals prospectively with blinded sample processing. In
10
11 addition, strategies to minimize false negatives and false positives (with potential
12
13 morbidity resulting from colonoscopy and treatment) will need to be developed. This
14
15 can be achieved by calculating a risk ratio on the basis of the individual expression
16
17 profile, which could replace the current binary output (cancer vs. healthy) and thus
18
19 define groups at risk that need to be subjected to colonoscopy as the gold standard.
20
21 An informed choice on the thresholds would, at least in first instance, emphasize a
22
23 high sensitivity at the expense of specificity. The resulting morbidity has to be
24
25 correlated to the morbidity of screening colonoscopy.
26
27

28
29 Our study raises important questions, which will need to be addressed in further
30
31 studies. First, the biological mechanisms and pivotal regulatory pathways in directing
32
33 the fate of the monocyte gene signature are still unexplored. Of note, only a few
34
35 genes appear to be commonly upregulated in CRC in comparison to gastric cancer
36
37 and pancreatic cancer. While this demonstrates specificity for CRC, it also means
38
39 further studies will be required to identify gene signatures specific to other tumours
40
41 and possibly benign pathological conditions. Second, we will need to assess if the
42
43 gene signature is already imprinted in pre-neoplastic lesions (*i.e.*, polyps) and
44
45 determine the transformation steps at which the specific upregulation occurs. Third,
46
47 as monocyte plasticity is the starting hypothesis of this study, we will need to assess
48
49 if treatment regimes (e.g., steroids, chemotherapy, irradiation, postoperative stress
50
51 conditions) affect the behaviour of the gene expression profile or interfere with its
52
53 diagnostic capabilities. Fourth, we are currently investigating in a prospective setting
54
55
56
57
58
59
60

1
2
3 if the gene signature is suitable for detection of relapse, as suggested by our
4
5 preliminary data. Last, future prospective studies will also reveal the significance of
6
7 this gene signature in early monitoring of treatment efficacy in metastatic disease.
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
51
52
53
54
55
56
57
58
59
60

Confidential: For Review Only

CONCLUSIONS

1
2
3
4
5
6 Taken together, these data provide unprecedented evidence that tumour-educated
7
8 monocytes exhibit a distinct and plastic gene signature, which may not only be
9
10 suitable for diagnosis of CRC, but potentially allows to monitor for success of therapy
11
12 or for relapse. As monocytes can be obtained in a non-invasive way, these findings
13
14 offer exciting new opportunities for both improving CRC diagnosis and enriching the
15
16 armamentarium of therapeutic strategies, provided that the data obtained here can
17
18 be replicated in an independent broad screening setting.
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
51
52
53
54
55
56
57
58
59
60

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

1
2
3 **a-d**, Stimulating freshly isolated, naïve monocytes with medium containing soluble
4 factors demonstrates that the genetic fingerprint in monocytes used for the diagnostic
5 gene signature is specifically induced by the transformed colon epithelium (HCT) in
6 comparison to a benign cell line (CCD), as demonstrated by expression analysis
7 comparing selective marker genes in stimulated monocytes to mock control (**a**).
8 Genetic alterations are independent of hypoxic cues (**b**). The gene signature is
9 specific to CRC in comparison to monocytes stimulated by a gastric cancer cell line
10 (MKN) (**c**). The gene signature is reverted after withdrawal of the stimulus *i.e.*, the
11 conditioned medium (**d**). n=6 (biological replicates from 6 different healthy donors);
12 bars, mean with SEM; *, p<0.05; **, p<0.01; ***, p<0.001; ****, p<0.0001; #, p<0.05
13 towards mock control, assessed by ANOVA with Bonferroni correction. All
14 experiments were repeated at least twice.

15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32 **Figure 4: The diagnostic gene signature is specific for CRC of all stages and**
33 **reverts upon curative treatment**

34
35
36 **a**, Expression of the gene signature in patients with cancer of the gastro-esophageal
37 junction (PG), demonstrating no upregulation and thus specificity of the diagnostic
38 signature for CRC. See Figure 2 for details on graphic elements. **b**, Gene signature
39 in patients after curative treatment (patients in remission, PR), in which the
40 expression levels revert to those of healthy volunteers in comparison to CRC
41 patients. **c, d**, ROC analyses corresponding to Figure 4a. **e**, Four patients with
42 isolated metastatic recurrence at the time of analysis (black dots) in a 2D-projection
43 of the multi-gene expression levels. The gene signature of metachronously
44 metastasized patients clusters with those patients with primary tumours (red), distinct
45 from healthy individuals (blue). *, p<0.05; **, p<0.01; ***, p<0.001

1
2
3
4
5 **Figure 5: Specificity of the gene signature to monocytes in comparison to**
6
7 **PBMCs**

8
9
10 **a**, Expression study assessing the gene signature in PBMCs in comparison to
11 monocytes (PBMs). While the entire signature is confirmed in PBMs in this
12 independent sample set, it is impossible to detect robust genetic alterations in
13 PBMCs, demonstrating specificity to PBMs. See Figure 2 for details on graphic
14 elements. **b**, Corresponding ROC analysis in PBMCs. **c**, ROC analysis of P,PM
15 versus HV in monocytes, confirming the previously established classification
16 performance. *, $p < 0.05$; **, $p < 0.01$; ***, $p < 0.001$
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
51
52
53
54
55
56
57
58
59
60

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.
3. Lieberman DA. Clinical practice. Screening for colorectal cancer. *The New England journal of medicine* 2009;**361**(12):1179-87.
4. 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.
5. 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.
6. Pawa N, Arulampalam T, Norton JD. Screening for colorectal cancer: established and emerging modalities. *Nature reviews Gastroenterology & hepatology* 2011;**8**(12):711-22.
7. 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.
8. Shi C, Pamer EG. Monocyte recruitment during infection and inflammation. *Nature reviews Immunology* 2011;**11**(11):762-74.
9. 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.

- 1
2
3 10. Sica A, Mantovani A. Macrophage plasticity and polarization: in vivo veritas. The
4
5 Journal of clinical investigation 2012;**122**(3):787-95.
6
- 7 11. Wynn TA, Chawla A, Pollard JW. Macrophage biology in development,
8
9 homeostasis and disease. Nature 2013;**496**(7446):445-55.
10
- 11 12. Irvine KM, Gallego P, An X, et al. Peripheral blood monocyte gene expression
12
13 profile clinically stratifies patients with recent-onset type 1 diabetes. Diabetes
14
15 2012;**61**(5):1281-90.
16
- 17 13. Zawada AM, Rogacev KS, Schirmer SH, et al. Monocyte heterogeneity in human
18
19 cardiovascular disease. Immunobiology 2012;**217**(12):1273-84.
20
- 21 14. Ma H, Hong M, Duan J, et al. Altered cytokine gene expression in peripheral
22
23 blood monocytes across the menstrual cycle in primary dysmenorrhea: a
24
25 case-control study. PloS one 2013;**8**(2):e55200.
26
27
- 28 15. Honda T, Inagawa H, Yamamoto I. Differential expression of mRNA in human
29
30 monocytes following interaction with human colon cancer cells. Anticancer
31
32 research 2011;**31**(7):2493-7.
33
34
- 35 16. Fletcher RH. Carcinoembryonic antigen. Annals of internal medicine
36
37 1986;**104**(1):66-73.
38
39
- 40 17. Schwarzenbach H, Hoon DS, Pantel K. Cell-free nucleic acids as biomarkers in
41
42 cancer patients. Nature reviews Cancer 2011;**11**(6):426-37.
43
44
- 45 18. Huang Z, Huang D, Ni S, et al. Plasma microRNAs are promising novel
46
47 biomarkers for early detection of colorectal cancer. International journal of
48
49 cancer Journal international du cancer 2010;**127**(1):118-26.
50
51
- 52 19. Church TR, Wandell M, Lofton-Day C, et al. Prospective evaluation of methylated
53
54 SEPT9 in plasma for detection of asymptomatic colorectal cancer. Gut 2013.
55
56
57
58
59
60

- 1
2
3 20. Xu Y, Xu Q, Yang L, et al. Gene expression analysis of peripheral blood cells
4
5 reveals toll-like receptor pathway deregulation in colorectal cancer. PloS one
6
7 2013;**8**(5):e62870.
8
9
10 21. Han M, Liew CT, Zhang HW, et al. Novel blood-based, five-gene biomarker set
11
12 for the detection of colorectal cancer. Clinical cancer research : an official
13
14 journal of the American Association for Cancer Research 2008;**14**(2):455-60.
15
16 22. Marshall KW, Mohr S, Khettabi FE, et al. A blood-based biomarker panel for
17
18 stratifying current risk for colorectal cancer. International journal of cancer
19
20 Journal international du cancer 2010;**126**(5):1177-86.
21
22
23 23. Nichita C, Ciarloni L, Monnier-Benoit S, et al. A novel gene expression signature
24
25 in peripheral blood mononuclear cells for early detection of colorectal cancer.
26
27 Alimentary pharmacology & therapeutics 2014;**39**(5):507-17.
28
29
30 24. Imperiale TF, Ransohoff DF, Itzkowitz SH, et al. Multitarget stool DNA testing for
31
32 colorectal-cancer screening. The New England journal of medicine
33
34 2014;**370**(14):1287-97.
35
36 25. Nyugen J, Agrawal S, Gollapudi S, et al. Impaired functions of peripheral blood
37
38 monocyte subpopulations in aged humans. Journal of clinical immunology
39
40 2010;**30**(6):806-13.
41
42
43 26. Smyth GK. Linear models and empirical bayes methods for assessing differential
44
45 expression in microarray experiments. Statistical applications in genetics and
46
47 molecular biology 2004;**3**:Article3.
48
49
50 27. Benjamini Y, Hochberg Y. Controlling the False Discovery Rate: A Practical and
51
52 Powerful Approach to Multiple Testing. Journal of the Royal Statistical Society
53
54 Series B (Methodological) 1995;**57**(1):289-300.
55
56
57
58
59
60

- 1
2
3 28. Sample size for microarray experiments. Secondary Sample size for microarray
4
5 experiments. <http://bioinformatics.mdanderson.org/MicroarraySampleSize/>.
6
- 7 29. Burges CJC. A Tutorial on Support Vector Machines for Pattern Recognition.
8
9 Data Min Knowl Discov 1998;**2**(2):121-67.
- 10 30. Breiman L. Random Forests. Mach Learn 2001;**45**(1):5-32.
- 11
12 31. Dietterich TG. Ensemble Methods in Machine Learning. Proceedings of the First
13
14 International Workshop on Multiple Classifier Systems: Springer-Verlag,
15
16 2000:1-15.
- 17
18 32. Wessels LF, Reinders MJ, Hart AA, et al. A protocol for building and evaluating
19
20 predictors of disease state based on microarray data. Bioinformatics
21
22 2005;**21**(19):3755-62.
- 23
24 33. Edgar R, Domrachev M, Lash AE. Gene Expression Omnibus: NCBI gene
25
26 expression and hybridization array data repository. Nucleic acids research
27
28 2002;**30**(1):207-10.
- 29
30 34. Piehler A, Grimholt R, Ovstebo R, et al. Gene expression results in
31
32 lipopolysaccharide-stimulated monocytes depend significantly on the choice of
33
34 reference genes. BMC Immunology 2010;**11**(1):21.
- 35
36 35. Guo C, Liu S, Wang J, et al. ACTB in cancer. Clinica chimica acta; international
37
38 journal of clinical chemistry 2013;**417**:39-44.
- 39
40 36. Jones S, Chen WD, Parmigiani G, et al. Comparative lesion sequencing provides
41
42 insights into tumor evolution. Proceedings of the National Academy of
43
44 Sciences of the United States of America 2008;**105**(11):4283-8.
- 45
46 37. Goede V, Coutelle O, Shimabukuro-Vornhagen A, et al. Analysis of Tie2-
47
48 expressing monocytes (TEM) in patients with colorectal cancer. Cancer
49
50 investigation 2012;**30**(3):225-30.
51
52
53
54
55
56
57
58
59
60

1
2
3 38. Schauer D, Starlinger P, Reiter C, et al. Intermediate monocytes but not TIE2-
4
5 expressing monocytes are a sensitive diagnostic indicator for colorectal
6
7 cancer. PloS one 2012;7(9):e44450.

8
9
10 39. Khatri P, Sirota M, Butte AJ. Ten Years of Pathway Analysis: Current Approaches
11
12 and Outstanding Challenges. PLoS Comput Biol 2012;8(2):e1002375.
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
51
52
53
54
55
56
57
58
59
60

1
2
3 **Tumour-Educated Circulating Monocytes are Powerful Candidate**
4
5 **Biomarkers for Diagnosis and Disease Follow-up of Colorectal Cancer**
6
7

8 Alexander Hamm^{1,2#}, Hans Prenen^{3#}, Wouter Van Delm^{4#}, Mario Di Matteo^{1,2}, Mathias Wenes^{1,2},
9 Estelle Delamarre^{1,2}, Thomas Schmidt⁵, Jürgen Weitz^{5,6}, Roberta Sarmiento⁷, Angelo Dezi⁷,
10 Giampietro Gasparini⁷, Françoise Rothé⁸, Robin Schmitz⁵, André D'Hoore⁹, Hannes Iserentant¹⁰,
11
12
13
14 Alain Hendlisz⁸ & Massimiliano Mazzone^{1,2}
15

16
17 ¹Lab of Molecular Oncology and Angiogenesis, Vesalius Research Center, VIB, Leuven, Belgium

18
19 ²Lab of Molecular Oncology and Angiogenesis, Vesalius Research Center, Department of Oncology,
20
21 KU Leuven, Leuven, Belgium

22
23 ³Digestive Oncology, University Hospitals Leuven and Department of Oncology, KU Leuven, Leuven,
24
25 Belgium

26
27 ⁴Nucleomics Core, VIB, Leuven, Belgium

28
29 ⁵Department of General, Visceral, and Transplantation Surgery, University of Heidelberg, Heidelberg,
30
31 Germany

32
33 ⁶Department of Visceral, Thoracic, and Vascular Surgery, University Hospital Carl Gustav Carus,
34
35 Technical University Dresden, Dresden, Germany

36
37 ⁷Department of Oncology, San Filippo Neri, Rome, Italy

38
39 ⁸Medical Oncology Clinic, Institut Jules Bordet, Brussels, Belgium

40
41 ⁹Department of Abdominal Surgery, University Hospitals Leuven, KU Leuven, Leuven, Belgium

42
43 ¹⁰VIB, Zwijnaarde, Belgium

44
45 #contributed equally to this study

46
47 **Correspondence:**

48
49 Massimiliano Mazzone, massimiliano.mazzone@vib-kuleuven.be,
50
51 Tel: +32-16-373213, Fax +32-16-372585
52
53 VIB Vesalius Research Center, KU Leuven, Herestraat 49, Bus 912, 3000 Leuven, Belgium

54
55 Hans Prenen, hans.prenen@uzleuven.be, Tel: +32-16-340238

56
57 Word Count: 4127

58
59 Key Words: Monocytes, colorectal cancer, screening, inflammation
60

LIST OF ABBREVIATIONS

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

Labels of patient groups:

<u>HV</u>	<u>healthy volunteer</u>
<u>P</u>	<u>non-metastatic CRC patient</u>
<u>P, PM</u>	<u>non-metastatic and metastatic CRC patients</u>
<u>PC</u>	<u>pancreatic cancer patient</u>
<u>PG</u>	<u>gastric cancer patient</u>
<u>PGT</u>	<u>gastritis patient</u>
<u>PM</u>	<u>metastatic CRC patient</u>
<u>PR</u>	<u>patient in remission from CRC</u>

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.

1
2
3 Conclusion: Our study is the first to demonstrate monocyte plasticity in response to
4
5 tumour-released soluble factors. The identified distinct signature in tumour-educated
6
7 monocytes might be used as candidate biomarker in CRC diagnosis and harbours
8
9 the potential for disease follow-up and therapeutic monitoring.
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
51
52
53
54
55
56
57
58
59
60

Confidential: For Review Only

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 gene signature may bring additive value to the established screening tools for CRC and early detection of recurrent disease, both offering patients better chances of cure. Moreover, the plasticity of monocytes may prove to be ideal for real-time follow-up of CRC treatment.

Confidential: For Review Only

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 type^{10 11}. 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.

1
2
3 Several novel accessible diagnostic tools share the major opportunity to make
4
5 frequent screening more appealing to a greater number of patients, as a less
6
7 invasive method is likely to increase compliance and allow for decreased screening
8
9 intervals (recently comprehensively reviewed⁶). While conventional blood-based
10
11 tumour markers (particularly carcino-embryonic antigen, CEA¹⁶) have been
12
13 established as supplemental markers in treatment monitoring, they have failed to
14
15 yield high diagnostic accuracy as primary screening tools. In addition to the
16
17 established FIT or FOBT⁵, other potential diagnostic markers include serum-
18
19 associated biomarkers (e.g. circulating tumour DNA¹⁷, micro-RNA¹⁸, methylation
20
21 markers like *SEPT9*¹⁹), genetic marker sets in white blood cells²⁰⁻²³, and, most
22
23 recently, fecal tumour DNA²⁴. However, all of these approaches display limited
24
25 sensitivity and specificity⁶. In this study, we therefore assess the sensitivity and
26
27 specificity of a novel gene signature in circulating monocytes for the diagnosis of
28
29 CRC in comparison to healthy individuals or to other cancer types, and assess its
30
31 robustness in prospectively obtained samples.
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

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

1
2
3 quantitative RT-PCR (qPCR) on the 7500Fast System (Applied Biosystems) using
4
5 intron-spanning PrimeTime qPCR Assays (Integrated DNA Technologies) listed in
6
7 Supplementary Table 1 as described in Supplementary Methods. For statistical
8
9 analysis, we followed a three-step top-down approach to construct a gene signature
10
11 for CRC, with details explained in Supplementary Methods.

16 **Multicentric validation study**

17
18 For validation of a diagnostic test, we used cohort II to train and validate a multi-gene
19
20 classifier. Splits in training and test sets for validation were performed by stratified
21
22 random sampling for centre of origin and class label as detailed out in Supplementary
23
24 Methods. Samples with missing values for more than 25% of the genes were
25
26 excluded from the analysis. We ruled out an effect of the class labeling on the
27
28 percentage of missing values with Fisher's exact test (Supplementary Table 2).
29

30
31 The training dataset was used to build three types of classifiers: a support vector
32
33 machine (SVM)²⁹ with linear kernel, a single-gene majority vote (SGMV) classifier,
34
35 and a random forest classifier (RF³⁰). Subsequently, we applied an ensemble
36
37 method³¹ that votes according to the majority of the three independent classifiers.
38
39 Performance was validated both with ranking (AUC) and classification (balanced
40
41 error rate, BER, Se, Sp) scores with 95% confidence intervals ([lower boundary;
42
43 upper boundary]). We explicitly opted for relatively simple computational models in
44
45 order to limit chances of over-fitting the training data and to maximize interpretability
46
47 of the models' internal decision-making process. Model flexibility was further
48
49 controlled through a Monte-Carlo cross-validation scheme (MCCV)³², before final
50
51 estimation of the model parameters. Validation of the predictive models was done on
52
53
54
55
56
57
58
59
60

1
2
3 the test set of cohort II, which were not included during development of the models.

4
5 Details on all classification methods are specified in Supplementary Methods.

6
7 In order to avoid biased conclusions, the analysis of the 23 genes was
8 complemented with a study by an independent team (DNAlytics, Belgium) that
9 adopted a slightly modified analysis protocol (see Supplementary Methods). All
10 complementary analyses were performed in R with scripts designed by DNAlytics,
11 fully independently from other analyses described in this paper.
12
13
14
15
16
17
18
19

20 **In vitro model system**

21
22 To study the effects of tumour-released soluble factors on gene expression in
23 monocytes, we established an *in vitro* model system, where monocytes from healthy
24 donors were challenged with tumour-released soluble factors and changes in gene
25 expression profile were analyzed by qPCR. See Supplementary Methods for details.
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
51
52
53
54
55
56
57
58
59
60

RESULTS

Establishment of putative biomarkers by genome-wide expression analysis

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.

TABLE 1: CLINICOPATHOLOGICAL CHARACTERISTICS OF PATIENTS AND HEALTHY VOLUNTEERS

Cohort	I		II								III		V	VI	
	P,PM	HV	P,PM				HV				P,PM	HV	P	P,PM	HV
			LEU ^a	HD ^b	SFN ^c	IJB ^d	LEU ^a	HD ^b	SFN ^c	IJB ^d					
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	/	16	3	2	2	/	/	/	/	16	/	0	6	/
non-metastatic	27	/	23	16	8	3	/	/	/	/	11	/	15	11	/
UICC stage															
1	3	/	7	2	1	1	/	/	/	/	2	/	4	2	/
2	12	/	8	8	2	0	/	/	/	/	3	/	7	7	/
3	12	/	8	6	5	2	/	/	/	/	6	/	4	2	/
4	28	/	16	3	2	2	/	/	/	/	16	/	0	6	/
Tumour localization															
Caecum	5	/	3	3	0	0	/	/	/	/	2	/	1	2	/
Ascendens	11	/	4	3	1	0	/	/	/	/	6	/	4	4	/
Transversum	0	/	4	3	0	0	/	/	/	/	2	/	1	0	/
Descendens	4	/	2	1	3	3	/	/	/	/	0	/	0	1	/
Sigmoid	28	/	15	3	2	0	/	/	/	/	10	/	5	6	/
Rectum	6	/	8	5	3	2	/	/	/	/	7	/	4	2	/
Double	1	/	3	1	1	0	/	/	/	/	0	/	0	2	/

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

1
2
3
4 The purity of the monocyte fraction was >90%, as assessed by FACS analysis in the
5
6 pilot phase (Supplementary Figure 1a) and verified by hemocytometric analysis for
7
8 each individual sample (Supplementary Figure 1b). Both absolute and relative
9
10 monocyte counts were not different between patients and healthy volunteers
11
12 (Supplementary Figure 1c). We therefore investigated differentially expressed genes
13
14 by genome-wide expression analysis using the Illumina HumanHT-12 v4 Expression
15
16 BeadChip Kit. The data discussed in this publication have been deposited in NCBI's
17
18 Gene Expression Omnibus³³ and are accessible through GEO Series accession
19
20 number
21
22 GSE47756
23
24 ([http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?token=hvmpvoswuqaeybc&acc=GSE](http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?token=hvmpvoswuqaeybc&acc=GSE47756)
25
26 [47756](http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?token=hvmpvoswuqaeybc&acc=GSE47756)). In first instance, we compared the average expression values of all CRC
27
28 patients (P,PM), comprised of non-metastatic (P) and metastatic (PM) patients, to
29
30 that of healthy volunteers (HV). The resulting gene signature of (P,PM) versus HV
31
32 consisted of 36 upregulated and 4 downregulated probes (Figure 2a, b, Table 2). In
33
34 second instance, we were interested if the gene signature in patients with
35
36 synchronous metastases *i.e.*, at the time of diagnosis (PM, n=28) was different from
37
38 that in non-metastatic patients (P, n=27). Interestingly, the number of up- and down-
39
40 regulated genes was comparable in both P and PM (in comparison to HV) (Table 2
41
42 and Supplementary Figure 2a, b), while there were no genes found to be differentially
43
44 expressed between the two patient groups (Supplementary Figure 2a, b), indicating
45
46 that the gene signature induced at early onset stays robust over disease progression.
47
48 Indeed, when post-hoc assessing those samples from patients with early stages (Tis
49
50 and T1), they clustered with the rest of the patient samples (data not shown). A
51
52 power analysis revealed that, for the given number of genes, samples and observed
53
54 variation, chances were very low ($<10^{-10}$) that truly differentially expressed genes with
55
56
57
58
59
60

1
2
3 fold changes larger than 1.5 had been missed. Therefore, adding more samples
4
5 would probably have changed little to the panel of candidate genes that our screen
6
7 returned.
8
9

10 **Confirmation of the gene signature in independently processed samples**

11
12 To validate the genetic signature, we performed quantitative RT-PCR (qPCR)
13
14 analysis on a random subset of PBM from 8 samples of each of the three groups (P,
15
16 PM, and HV), normalizing to reference gene *B2M*, which was selected after an
17
18 extensive screening procedure (Supplementary Note 1). To avoid bias in the
19
20 confirmation procedure, we freshly extracted RNA from independently stored
21
22 samples for confirmative expression analysis. In analyzing 43 putative marker genes
23
24 with probes listed in Supplementary Table 1, 23 genes showed differential
25
26 expression between (P, PM) and HV (Supplementary Figure 3b, Table 2, and
27
28 Supplementary Table 4). Thus, we were able to confirm a subset of the previously
29
30 established gene signature, independent of the RNA extraction and the platform used
31
32 for expression analysis. Information on the annotated biological function of the genes
33
34 of the diagnostic signature can be found in Supplementary Table 5 and
35
36 Supplementary Note 2.
37
38
39
40
41
42
43
44

45 **Confirmation of the gene signature in a multicentric validation set**

46
47 For a rigorous validation of the gene signature, we collected an independent
48
49 multicentric validation set (cohort II) from a total of 4 different European oncological
50
51 centres with stratified training and test sets as described in Supplementary Methods.
52
53 Using the panel of 23 genes confirmed previously, we found consistently differential
54
55 expression between all patients and the healthy volunteers (Figure 2c and
56
57
58
59
60

1
2
3 Supplementary Figure 4). In line with the findings from the screening phase, there
4
5 were no detectable differences in expression levels between P and PM
6
7 (Supplementary Figure 5), while either patient group alone compared to HV was
8
9 differentially expressed (data not shown).

10
11 In ROC analysis for single genes, we found that some, but not all of the genes that
12
13 displayed significantly differential expression were able to discriminate patient
14
15 samples from healthy individual samples with acceptable AUCs (Supplementary
16
17 Figure 6 and data not shown). We therefore hypothesized that a marker panel
18
19 consisting of multiple genes might yield better results in discriminating sample
20
21 identity. To address this question, we decided to test three different classification
22
23 algorithms on this data set, namely a support vector machine (SVM)²⁹ with linear
24
25 kernel, a single-gene majority vote (SGMV) classifier, a random forest classifier
26
27 (RF³⁰), and a combined classification by an ensemble method³¹, using the outcome
28
29 of the three classification algorithms for a final diagnostic decision. To limit over-
30
31 estimation of the performance by the particular training and test set, we performed a
32
33 MCCV as a conservative estimate with 1,000 cross-validations. Performance of all
34
35 classification algorithms in cohorts II – VI, including the conservative estimate of the
36
37 MCCV in cohort II, is given in detail in Table 2.
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

TABLE 2: PERFORMANCE SCORES OF MULTIGENE CLASSIFIER

	SGMV	SVM	RF	ENS
Cohort II (Validation)				
AUC [95% CI]	0.99 [0.99;1.00]	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]	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]
Cohort IV(pancreatic cancer)				
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]
Cohort IV (gastritis)				
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]
Cohort V (PR)				
Sensitivity [95% CI]	50.0 [20.0;80.0]	10.0 [0.0;30.0]	20.0 [0.0;50.0]	20.0 [0.0;50.0]
Cohort VI (PBMC)				
AUC [95% CI]	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

1
2
3 Strikingly, we achieved a remarkably high AUC of 0.99 [0.99;1.00] with a BER of
4
5 3.6% (Figure 2d and Table 2), translating into a sensitivity of 100.0%
6
7 [100.0%;100.0%] and a specificity of 92.9% [78.6%;100.0%] (Table 2). Neither of the
8
9 classification algorithms was capable of separating P from PM or detect differences
10
11 dependent on tumour localization (Supplementary Note 3).

12
13
14 In order to assess whether the diagnostic gene signature is actually suitable for
15
16 diagnosis of CRC in a screening setting, we have initiated a prospective sample
17
18 collection in both patients and healthy individuals who are subjected to colonoscopy.
19
20 In a pilot analysis in 27 patients (newly diagnosed with CRC by screening
21
22 colonoscopy) and 13 healthy individuals negative to screening colonoscopy (cohort
23
24 III), we found an AUC of 0.95 [0.85;1.00] with a BER of 7.6%, yielding a sensitivity of
25
26 92.6% [81.5%;100.0%] and a specificity of 92.3% [76.9%;100.0%] (Figure 2e and
27
28 Table 2). The complementary data analysis (independently performed by DNAlytics,
29
30 Belgium) on the same panel of 23 genes led to the matching conclusions in terms of
31
32 performance. The first experiment consisted in cross-validating a model on Cohort II
33
34 (BER: 8.4% [3.4%;13.4%]; AUC: 0.93 [0.88;0.98]). A second experiment consisted in
35
36 learning the same type of model on Cohort II and having it make predictions on
37
38 Cohort III (BER: 13.2%; AUC: 0.92).

39
40
41
42
43
44
45 **Soluble factors released by colorectal cancer cells induce an early, tumour**
46
47 **type-specific and reversible genetic fingerprint in monocytes**

48
49 We hypothesized that tumour-released soluble factors are the key players in inducing
50
51 the genetic signature in circulating monocytes. Thus, we established an *in vitro*
52
53 model system where we cultured freshly isolated human monocytes from healthy
54
55 donors in different conditions. In order to assess alterations in gene expression, we
56
57
58
59
60

1
2
3 first analyzed which of the 23 genes comprising the gene signature was up- or
4
5 downregulated in culture after 72 hours without any additional stimulus and excluded
6
7 these from the further *in vitro* studies (Supplementary Figure 7). Out of the remaining
8
9 gene signature, the majority (7/9) was specifically upregulated when culturing naïve
10
11 monocytes in medium conditioned by the CRC cell line HCT116, while expression
12
13 levels were not affected by mock medium (Figure 3a). Moreover, in line with the
14
15 coherent induction of the specific signature independent of the stage of the disease,
16
17 the induction *in vitro* was independent of hypoxic cues, as HCT116-conditioned
18
19 medium in hypoxia did not induce any different expression levels than medium
20
21 obtained in normoxia (Figure 3b). Likewise, the changes in expression levels of all
22
23 these genes occurred already 18 hours after stimulating monocytes with the
24
25 conditioned medium, consistent with the fact that already early stages are detectable
26
27 by the diagnostic signature.
28
29
30

31
32 To rule out an off-target effect of conditioned medium *i.e.*, unspecific cues from cell
33
34 metabolites, apoptotic bodies, pH, etc., we assessed the expression levels of the
35
36 genes upregulated by HCT116-conditioned medium in comparison to a benign colon
37
38 epithelium cell line, CCD 841 CoN (CCD), which did not induce alterations in gene
39
40 expression levels different from the Mock control (Figure 3a).
41
42

43 Prompted by this finding, we investigated if the induction of the genetic signature was
44
45 a general effect of malignant transformation or might be specific to the histotype of
46
47 cancer. To address this question, we conditioned medium with a gastric cancer cell
48
49 line, MKN-45 (MKN), to compare CRC to another frequent gastrointestinal solid
50
51 neoplasm. Remarkably, when comparing the expression levels in naïve monocytes
52
53 upon stimulation with the different conditioned media, we found that MKN-45
54
55
56
57
58
59
60

1
2
3 conditioned medium did not induce the same upregulation of the genes of interest as
4
5 HCT116 conditioned medium (Figure 3c).
6

7
8 As immune cells are highly versatile and plastic cells mirroring the microenvironment,
9
10 where they are embedded, we reasoned that the genetic signature induced by CRC
11
12 in monocytes might be dependent on the continuous presence of the stimulating
13
14 agents and thus be reversible upon inversion of the conditions. We therefore
15
16 incubated naïve monocytes first with HCT116-conditioned medium for 18 hours and
17
18 then refreshed the medium with plain culture medium, thus withdrawing the tumour-
19
20 released soluble factors. Strikingly, the previously elevated expression levels of a set
21
22 of marker genes were almost entirely reverted to the original (and to the mock
23
24 control) expression levels 72 hours after withdrawing the tumour-cell conditioned
25
26 medium (Figure 3d), whereas they remained constantly overexpressed when the
27
28 conditioned medium was maintained (data not shown).
29
30
31
32
33

34 **The monocyte signature is specific for CRC and might serve as a candidate** 35 **biomarker of disease follow-up** 36

37
38 Based on the *in vitro* results showing that the genetic signature is specific to CRC,
39
40 we sought to confirm these findings *in vivo*. We therefore assessed the diagnostic
41
42 signature in patients with i. cancer of the stomach and gastro-esophageal junction
43
44 (PG, n=15) and ii. pancreatic ductal adenocarcinoma (PC, n=16), two other frequent
45
46 cancers of the gastrointestinal tract¹. In addition, we analysed iii. patients with
47
48 gastritis (PGT, n=10) in order to compare the gene signature in CRC to a benign
49
50 inflammatory condition of the gastrointestinal tract (cohort IV). In line with the *in vitro*
51
52 results we saw that the vast majority of all genes were not significantly different
53
54 between either of the patient groups and healthy volunteers, indicating the specificity
55
56
57
58
59
60

1
2
3 of this monocyte imprinting by colorectal cancer cells (Figure 4a and data not
4 shown). Moreover, the classifier established to diagnose CRC could not separate
5 patients with gastric cancer (AUC 0.63 [0.48;0.77]), pancreatic cancer (AUC 0.41
6 [0.27;0.50]), or gastritis (AUC 0.52 [0.35;0.68]) from healthy individuals (Figure 4c, d
7 and Table 2).

8
9
10
11
12
13
14 The finding that the genetic signature is reverted upon withdrawal of the stimulating
15 agents prompted us to investigate in a pilot phase the behaviour of the entire
16 diagnostic signature in patients upon curative treatment *i.e.*, patients with surgically
17 removed tumours without any evidence of residual disease. To this end, we isolated
18 monocytes from 15 patients of stages I to III treated with curative intent (with or
19 without adjuvant treatment) and presenting at follow-up without detectable residual
20 disease (PR) (cohort V). Here, we found that virtually all of the previously
21 upregulated genes were reverted to expression levels comparable to those of healthy
22 volunteers (Figure 4b). Consequently, when applying the previously established
23 classifier, we found that it was able to distinguish accurately between patients in
24 remission and patients with tumour, while it could not detect differences between
25 patients in remission and healthy volunteers (Table 2).

26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41 Finally, as the plasticity of the signature offers the perspective to use the gene
42 signature for follow-up of treated patients, we became interested if the same
43 signature could be used to diagnose relapse (frequently as metachronous
44 metastases rather than local recurrence²). Although our dataset was not powered to
45 address this question with sufficient significance, we post-hoc identified four patients
46 from cohorts I and II included at presentation with metachronous metastases. All four
47 clustered clearly in the group of patients, separately from the healthy volunteers
48 (Figure 4e), suggesting that the signature might be used to detect disease relapse in
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 line with the previous results that show coherent expression over disease
4
5 progression.
6
7

8
9 **The gene signature is specific to monocytes in comparison to all peripheral**
10 **blood mononuclear cells (PBMCs)**
11

12
13
14 To rigorously assess if the genetic fingerprint identified in monocytes was specific to
15
16 this cell type or an epiphenomenon of genetic shifts in the entire population of
17
18 PBMCs, we isolated both monocytes and full PBMC fractions from 17 patients and 7
19
20 healthy volunteers for a comparative analysis (cohort VI). Interestingly, we found that
21
22 while in the monocyte population, the diagnostic marker set of 23 genes was
23
24 upregulated in all patients (both P and PM) in accordance with our previous results
25
26 (Figure 5a), there were no significant differences in the expression levels of the
27
28 analyzed genes in the full PBMC compartment when comparing patients to healthy
29
30 volunteers (Figure 5a). Consistently, applying the previously established classifier
31
32 with the defined cut-off values, it was impossible to separate the patient group from
33
34 the healthy volunteer group in PBMC (Figure 5b and Table 2), while the classifier
35
36 confirmed its accuracy in PBM (Figure 5c and Table 2). Thus, the differential
37
38 regulation of the gene signature in PBM used for CRC diagnosis is specific to the
39
40 monocytic lineage, reinforcing our initial working hypothesis that these cells are
41
42 specifically affected by tumour-secreted factors.
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

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 strategies^{2,6}. 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 biology-bound 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

1
2
3 most likely depends on the development of a one-step assay with capture of
4 monocytes from whole blood and gene expression analysis in a multiplex qPCR
5 assay with absolute quantification, avoiding extensive preanalytical processing steps.
6
7 However, the analytical reliability of this assay needs to be thoroughly established,
8 most likely requiring centralization of the analysis during the first phase of
9 distribution.

10
11
12
13
14
15
16 The finding that the specific gene signature is reversible if the stimulating cues are
17 withdrawn, was not only demonstrated *in vitro*, but also in a pilot analysis *in vivo* in
18 samples of patients after curative treatment. Although not completely unexpected in
19 view of the plastic nature of the monocyte-macrophage lineage, this analysis opens
20 avenues for treatment monitoring and companion diagnostics and will be assessed in
21 detail in a prospective study during patient follow-up.
22
23
24
25
26
27
28
29
30

31
32 If supported by further prospective validation studies, this gene signature may
33 outperform other published non-invasive test for CRC diagnosis⁶ (including single
34 surface markers in monocytes^{37 38}) or score similar to the most recent evaluation of
35 fecal tumour DNA²⁴. Moreover, we are the first to demonstrate that a potential
36 diagnostic biomarker obtained in patients at the time of primary diagnosis might also
37 be suitable for disease follow-up and thus assessment of treatment response, owing
38 to its high plasticity.
39
40
41
42
43
44
45
46

47 We acknowledge the limited conclusions that can be drawn from our case-control
48 study. Despite the confirmation in independent samples, we cannot fully exclude
49 possible confounders that can only be unveiled by a blinded, prospective sample
50 collection in screening individuals. These include, but are not limited to, the bias of
51 selecting patients that underwent colonoscopy for a clinical indication; the differences
52
53
54
55
56
57
58
59
60

1
2
3 in age, nutrition status, diet, and potentially lifestyle between patients and healthy
4 volunteers; the unblinded sample collection and processing. It is therefore of
5 paramount importance that a prospective validation study initiated by our group
6 includes screening individuals prospectively with blinded sample processing. In
7 addition, strategies to minimize false negatives and false positives (with potential
8 morbidity resulting from colonoscopy and treatment) will need to be developed. This
9 can be achieved by calculating a risk ratio on the basis of the individual expression
10 profile, which could replace the current binary output (cancer vs. healthy) and thus
11 define groups at risk that need to be subjected to colonoscopy as the gold standard.
12 An informed choice on the thresholds would, at least in first instance, emphasize a
13 high sensitivity at the expense of specificity. The resulting morbidity has to be
14 correlated to the morbidity of screening colonoscopy.

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
51
52
53
54
55
56
57
58
59
60

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

1
2
3 if the gene signature is suitable for detection of relapse, as suggested by our
4 preliminary data. Last, future prospective studies will also reveal the significance of
5 this gene signature in early monitoring of treatment efficacy in metastatic disease.
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
51
52
53
54
55
56
57
58
59
60

Confidential: For Review Only

CONCLUSIONS

1
2
3
4
5
6 Taken together, these data provide unprecedented evidence that tumour-educated
7
8 monocytes exhibit a distinct and plastic gene signature, which may not only be
9
10 suitable for diagnosis of CRC, but potentially allows to monitor for success of therapy
11
12 or for relapse. As monocytes can be obtained in a non-invasive way, these findings
13
14 offer exciting new opportunities for both improving CRC diagnosis and enriching the
15
16 armamentarium of therapeutic strategies, provided that the data obtained here can
17
18 be replicated in an independent broad screening setting.
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
51
52
53
54
55
56
57
58
59
60

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

1
2
3 **a-d**, Stimulating freshly isolated, naïve monocytes with medium containing soluble
4 factors demonstrates that the genetic fingerprint in monocytes used for the diagnostic
5 gene signature is specifically induced by the transformed colon epithelium (HCT) in
6 comparison to a benign cell line (CCD), as demonstrated by expression analysis
7 comparing selective marker genes in stimulated monocytes to mock control (**a**).
8 Genetic alterations are independent of hypoxic cues (**b**). The gene signature is
9 specific to CRC in comparison to monocytes stimulated by a gastric cancer cell line
10 (MKN) (**c**). The gene signature is reverted after withdrawal of the stimulus *i.e.*, the
11 conditioned medium (**d**). n=6 (biological replicates from 6 different healthy donors);
12 bars, mean with SEM; *, p<0.05; **, p<0.01; ***, p<0.001; ****, p<0.0001; #, p<0.05
13 towards mock control, assessed by ANOVA with Bonferroni correction. All
14 experiments were repeated at least twice.

15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32 **Figure 4: The diagnostic gene signature is specific for CRC of all stages and**
33 **reverts upon curative treatment**

34
35
36 **a**, Expression of the gene signature in patients with cancer of the gastro-esophageal
37 junction (PG), demonstrating no upregulation and thus specificity of the diagnostic
38 signature for CRC. See Figure 2 for details on graphic elements. **b**, Gene signature
39 in patients after curative treatment (patients in remission, PR), in which the
40 expression levels revert to those of healthy volunteers in comparison to CRC
41 patients. **c, d, ROC analyses corresponding to Figure 4a.** **e**, Four patients with
42 isolated metastatic recurrence at the time of analysis (black dots) in a 2D-projection
43 of the multi-gene expression levels. The gene signature of metachronously
44 metastasized patients clusters with those patients with primary tumours (red), distinct
45 from healthy individuals (blue). *, p<0.05; **, p<0.01; ***, p<0.001

1
2
3
4
5 **Figure 5: Specificity of the gene signature to monocytes in comparison to**
6
7 **PBMCs**

8
9
10 **a**, Expression study assessing the gene signature in PBMCs in comparison to
11 monocytes (PBMs). While the entire signature is confirmed in PBMs in this
12 independent sample set, it is impossible to detect robust genetic alterations in
13 PBMCs, demonstrating specificity to PBMs. See Figure 2 for details on graphic
14 elements. **b**, Corresponding ROC analysis in PBMCs. **c**, ROC analysis of P,PM
15 versus HV in monocytes, confirming the previously established classification
16 performance. *, $p < 0.05$; **, $p < 0.01$; ***, $p < 0.001$
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
51
52
53
54
55
56
57
58
59
60

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.
3. Lieberman DA. Clinical practice. Screening for colorectal cancer. *The New England journal of medicine* 2009;**361**(12):1179-87.
4. 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.
5. 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.
6. Pawa N, Arulampalam T, Norton JD. Screening for colorectal cancer: established and emerging modalities. *Nature reviews Gastroenterology & hepatology* 2011;**8**(12):711-22.
7. 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.
8. Shi C, Pamer EG. Monocyte recruitment during infection and inflammation. *Nature reviews Immunology* 2011;**11**(11):762-74.
9. 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.

- 1
2
3 10. Sica A, Mantovani A. Macrophage plasticity and polarization: in vivo veritas. The
4
5 Journal of clinical investigation 2012;**122**(3):787-95.
6
- 7 11. Wynn TA, Chawla A, Pollard JW. Macrophage biology in development,
8
9 homeostasis and disease. Nature 2013;**496**(7446):445-55.
10
- 11 12. Irvine KM, Gallego P, An X, et al. Peripheral blood monocyte gene expression
12
13 profile clinically stratifies patients with recent-onset type 1 diabetes. Diabetes
14
15 2012;**61**(5):1281-90.
16
- 17 13. Zawada AM, Rogacev KS, Schirmer SH, et al. Monocyte heterogeneity in human
18
19 cardiovascular disease. Immunobiology 2012;**217**(12):1273-84.
20
- 21 14. Ma H, Hong M, Duan J, et al. Altered cytokine gene expression in peripheral
22
23 blood monocytes across the menstrual cycle in primary dysmenorrhea: a
24
25 case-control study. PloS one 2013;**8**(2):e55200.
26
27
- 28 15. Honda T, Inagawa H, Yamamoto I. Differential expression of mRNA in human
29
30 monocytes following interaction with human colon cancer cells. Anticancer
31
32 research 2011;**31**(7):2493-7.
33
- 34 16. Fletcher RH. Carcinoembryonic antigen. Annals of internal medicine
35
36 1986;**104**(1):66-73.
37
- 38 17. Schwarzenbach H, Hoon DS, Pantel K. Cell-free nucleic acids as biomarkers in
39
40 cancer patients. Nature reviews Cancer 2011;**11**(6):426-37.
41
- 42 18. Huang Z, Huang D, Ni S, et al. Plasma microRNAs are promising novel
43
44 biomarkers for early detection of colorectal cancer. International journal of
45
46 cancer Journal international du cancer 2010;**127**(1):118-26.
47
48
- 49 19. Church TR, Wandell M, Lofton-Day C, et al. Prospective evaluation of methylated
50
51 SEPT9 in plasma for detection of asymptomatic colorectal cancer. Gut 2013.
52
53
54
55
56
57
58
59
60

- 1
2
3 20. Xu Y, Xu Q, Yang L, et al. Gene expression analysis of peripheral blood cells
4 reveals toll-like receptor pathway deregulation in colorectal cancer. PloS one
5 2013;**8**(5):e62870.
6
7
8
9
10 21. Han M, Liew CT, Zhang HW, et al. Novel blood-based, five-gene biomarker set
11 for the detection of colorectal cancer. Clinical cancer research : an official
12 journal of the American Association for Cancer Research 2008;**14**(2):455-60.
13
14
15
16 22. Marshall KW, Mohr S, Khettabi FE, et al. A blood-based biomarker panel for
17 stratifying current risk for colorectal cancer. International journal of cancer
18 Journal international du cancer 2010;**126**(5):1177-86.
19
20
21
22
23 23. Nichita C, Ciarloni L, Monnier-Benoit S, et al. A novel gene expression signature
24 in peripheral blood mononuclear cells for early detection of colorectal cancer.
25 Alimentary pharmacology & therapeutics 2014;**39**(5):507-17.
26
27
28
29
30 24. Imperiale TF, Ransohoff DF, Itzkowitz SH, et al. Multitarget stool DNA testing for
31 colorectal-cancer screening. The New England journal of medicine
32 2014;**370**(14):1287-97.
33
34
35
36 25. Nyugen J, Agrawal S, Gollapudi S, et al. Impaired functions of peripheral blood
37 monocyte subpopulations in aged humans. Journal of clinical immunology
38 2010;**30**(6):806-13.
39
40
41
42
43 26. Smyth GK. Linear models and empirical bayes methods for assessing differential
44 expression in microarray experiments. Statistical applications in genetics and
45 molecular biology 2004;**3**:Article3.
46
47
48
49 27. Benjamini Y, Hochberg Y. Controlling the False Discovery Rate: A Practical and
50 Powerful Approach to Multiple Testing. Journal of the Royal Statistical Society
51 Series B (Methodological) 1995;**57**(1):289-300.
52
53
54
55
56
57
58
59
60

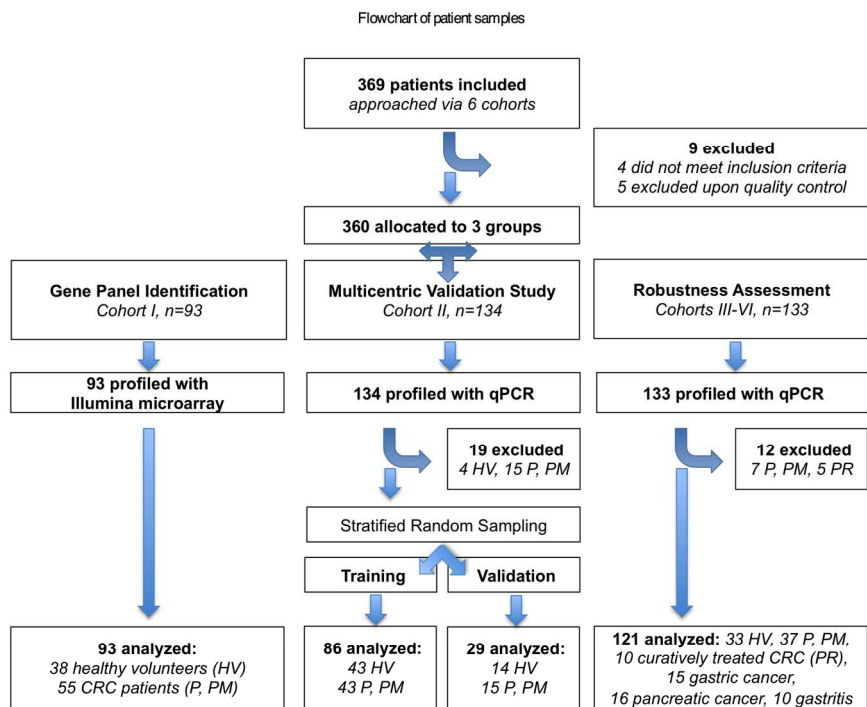
- 1
2
3 28. Sample size for microarray experiments. Secondary Sample size for microarray
4 experiments. <http://bioinformatics.mdanderson.org/MicroarraySampleSize/>.
5
6
7 29. Burges CJC. A Tutorial on Support Vector Machines for Pattern Recognition.
8 Data Min Knowl Discov 1998;**2**(2):121-67.
9
10 30. Breiman L. Random Forests. Mach Learn 2001;**45**(1):5-32.
11
12 31. Dietterich TG. Ensemble Methods in Machine Learning. Proceedings of the First
13 International Workshop on Multiple Classifier Systems: Springer-Verlag,
14 2000:1-15.
15
16 32. Wessels LF, Reinders MJ, Hart AA, et al. A protocol for building and evaluating
17 predictors of disease state based on microarray data. Bioinformatics
18 2005;**21**(19):3755-62.
19
20 33. Edgar R, Domrachev M, Lash AE. Gene Expression Omnibus: NCBI gene
21 expression and hybridization array data repository. Nucleic acids research
22 2002;**30**(1):207-10.
23
24 34. Piehler A, Grimholt R, Ovstebo R, et al. Gene expression results in
25 lipopolysaccharide-stimulated monocytes depend significantly on the choice of
26 reference genes. BMC Immunology 2010;**11**(1):21.
27
28 35. Guo C, Liu S, Wang J, et al. ACTB in cancer. Clinica chimica acta; international
29 journal of clinical chemistry 2013;**417**:39-44.
30
31 36. Jones S, Chen WD, Parmigiani G, et al. Comparative lesion sequencing provides
32 insights into tumor evolution. Proceedings of the National Academy of
33 Sciences of the United States of America 2008;**105**(11):4283-8.
34
35 37. Goede V, Coutelle O, Shimabukuro-Vornhagen A, et al. Analysis of Tie2-
36 expressing monocytes (TEM) in patients with colorectal cancer. Cancer
37 investigation 2012;**30**(3):225-30.
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 38. Schauer D, Starlinger P, Reiter C, et al. Intermediate monocytes but not TIE2-
4
5 expressing monocytes are a sensitive diagnostic indicator for colorectal
6
7 cancer. PloS one 2012;7(9):e44450.

8
9
10 39. Khatri P, Sirota M, Butte AJ. Ten Years of Pathway Analysis: Current Approaches
11
12 and Outstanding Challenges. PLoS Comput Biol 2012;8(2):e1002375.
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
51
52
53
54
55
56
57
58
59
60

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
51
52
53
54
55
56
57
58
59
60

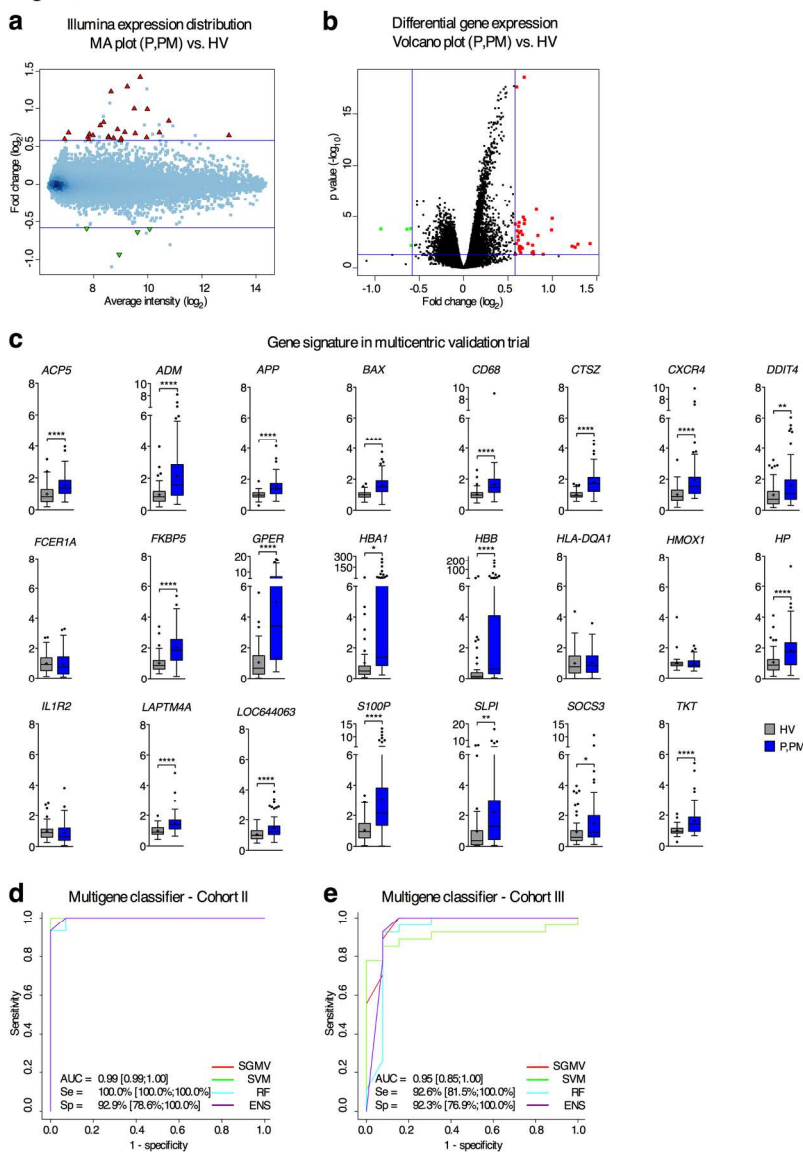
Figure 1



193x153mm (300 x 300 DPI)

Review Only

Figure 2

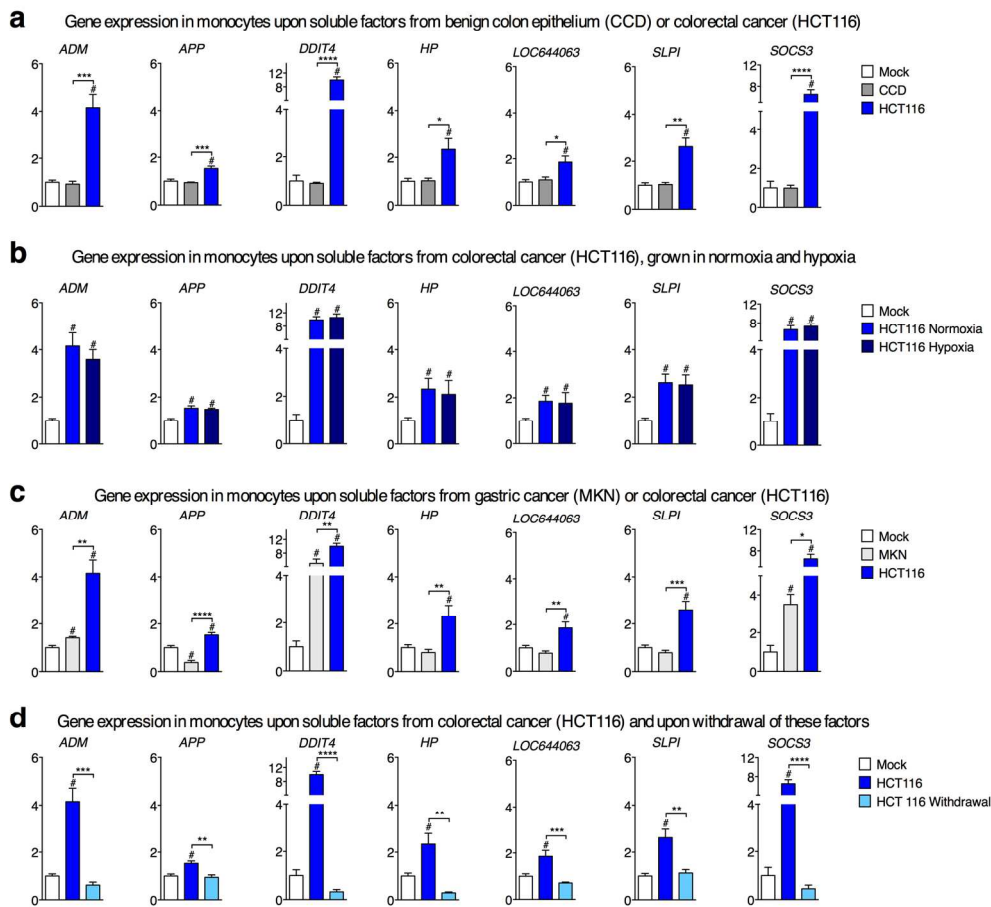


154x226mm (300 x 300 DPI)

only

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
51
52
53
54
55
56
57
58
59
60

Figure 3

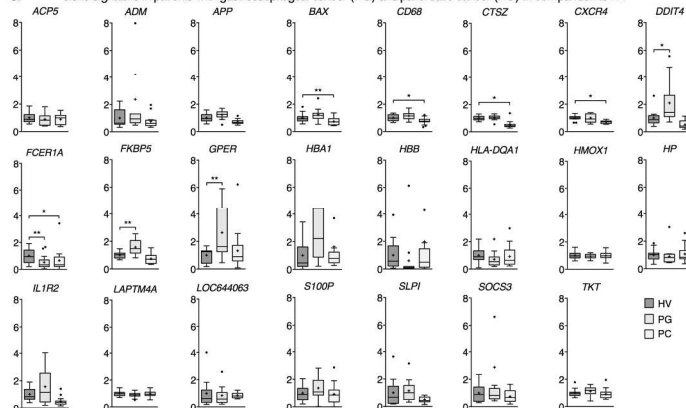


169x164mm (300 x 300 DPI)

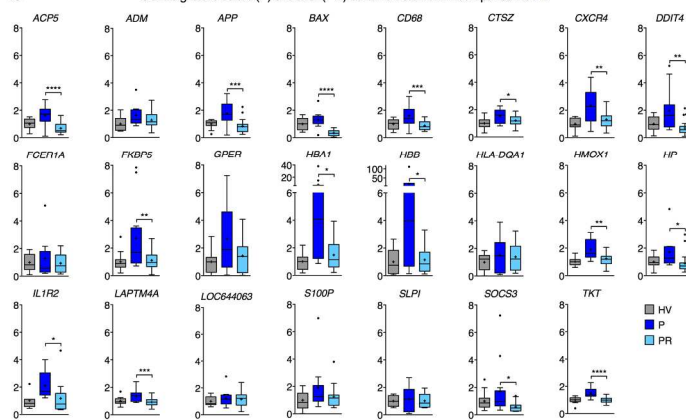
ew Only

Figure 4

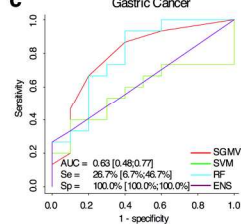
a Gene signature in patients with gastroesophageal cancer (PG) and pancreatic cancer (PC) in comparison to HV



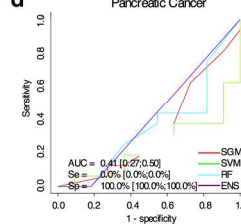
b Gene signature before (P) and after (PR) curative treatment in comparison to HV



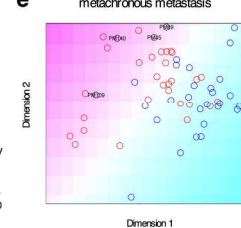
c Multigene Classifier Gastric Cancer



d Multigene Classifier Pancreatic Cancer

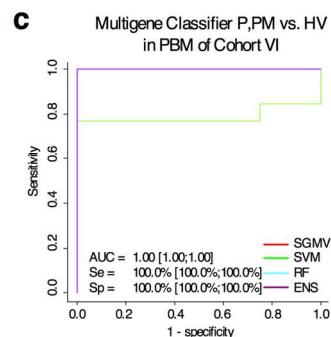
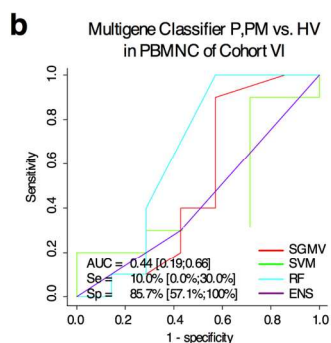
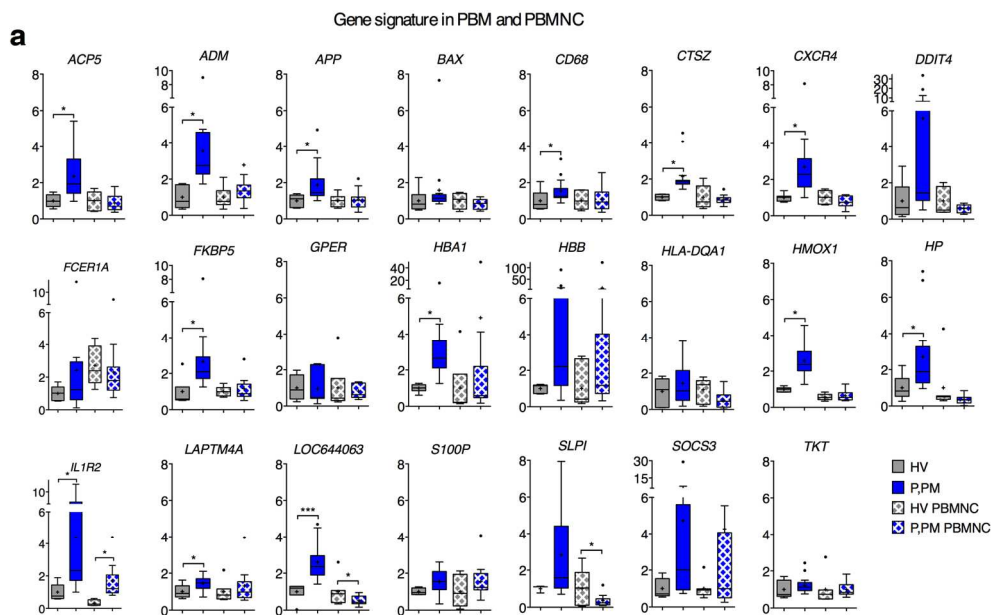


e Clustering of patients with metachronous metastasis



164x264mm (300 x 300 DPI)

Figure 5



164x179mm (300 x 300 DPI)

Only

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

CONTENTS

Supplementary Methods	Page 3
Supplementary Notes	Page 14
Supplementary Figures	Page 17
Supplementary Tables	Page 28
Supplementary References	Page 35

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
51
52
53
54
55
56
57
58
59
60

Confidential: For Review Only

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

1
2
3 gut flora. Samples were excluded from further analysis when final histology of the
4
5 surgical specimen did not confirm adenocarcinoma of the large intestine (assessed
6
7 by board-certified pathologists within clinical routine procedures).
8

9
10 The following four oncological centres contributed samples to this study: Digestive
11
12 Oncology, University Hospitals Leuven and Department of Oncology, KU Leuven,
13
14 Leuven, Belgium; Department of General, Visceral, and Transplantation Surgery,
15
16 University of Heidelberg, Heidelberg, Germany; Department of Oncology, San Filippo
17
18 Neri, Rome, Italy; Medical Oncology Clinic, Institut Jules Bordet, Brussels, Belgium.
19
20 The responsible scientists in each centre (1-2 per centre) were trained in the protocol
21
22 for isolation of PBM to ensure uniformity of the procedure. All participants gave
23
24 written informed consent, and the study was approved by the respective institutional
25
26 review boards (Leuven: B322201215873, Brussels: CE1950, Heidelberg: 323/2004,
27
28 Rome: 319/51). No adverse events from blood collection or colonoscopy were
29
30 recorded in included participants.
31
32
33
34
35

36 **Isolation of PBM**

37
38 20ml of EDTA-anticoagulated peripheral venous blood was collected following clinical
39
40 routine procedure, stored at 4°C and processed within 2 hours of blood collection.
41
42 For further isolation, blood was diluted 1:2 with DPBS (free of Ca²⁺ and Mg²⁺) and
43
44 layered carefully on Lymphoprep (Axis-Shield) in two separate tubes. All blood
45
46 collection and isolation steps were performed identical for samples of all origin.
47
48 Density gradient centrifugation was performed at 1,200g for 20 minutes at low
49
50 acceleration and no brake. Samples with macroscopically visible hemolysis were
51
52 excluded from further analysis. The PBMC interphase was collected carefully and
53
54 washed twice for 12 minutes at 250g and 175g with PBS. Hemocytometric analysis
55
56
57
58
59
60

1
2
3 was performed to ensure purity of PBMCs, and the pellet was pooled for further
4
5 processing and washed once for 10 minutes at 300g. Cells were then incubated with
6
7 CD14 magnetically-conjugated beads (BD) for 15 minutes at 4°C, washed 10
8
9 minutes at 300g and positively separated with the MACS system (Miltenyi) following
10
11 the manufacturer's instructions. The CD14+ fraction was flushed out and washed
12
13 once 10 minutes at 300g. Purity was assessed by FACS analysis for CD14 in the
14
15 pilot phase and by hemocytometric analysis (CellDyn 3700, Abbott) in every further
16
17 sample. Only samples with purity of >90% and viability >95% (assessed by Trypan
18
19 Blue staining) were retained for further analysis. Cell pellets were lysed in Buffer RLT
20
21 (Qiagen) at 10^6 monocytes in 350 μ l of Buffer RLT and stored at -80°C. For each
22
23 respective expression study, all samples were extracted simultaneously with the
24
25 RNeasy Mini Kit (Qiagen) following the manufacturer's instructions. Quality control
26
27 was performed by checking RNA quality on the Nanodrop system, and RNA integrity
28
29 was checked for microarray samples on the Agilent Bio-Analyzer. Only samples with
30
31 an extinction fraction 260/280 > 1.8 and 260/230 > 1.5, and an RNA integrity index of
32
33 >6 were retained for further analysis.
34
35
36
37
38
39
40

41 **Genome-wide expression analysis**

42
43 For genome-wide expression analysis, RNA was amplified and biotinylated using
44
45 Illumina TotalPrep RNA Amplification Kit (Ambion) following the manufacturer's
46
47 instructions to obtain biotinylated cRNA, which was hybridized to Illumina HumanHT-
48
49 12 v4 Expression BeadChips (Illumina) with the Illumina Whole-Genome Gene
50
51 Expression Direct Hybridization Assay (Illumina) following the manufacturer's
52
53 instructions. The Illumina HumanHT-12 v4 Expression BeadChip Kit contains 47,323
54
55 probes and 887 controls. After scanning, background-corrected expression values
56
57
58
59
60

1
2
3 and detection scores were extracted with GenomeStudio GX (version 1.5.4). For
4 each array, we used the summarized expression level (AVG_Signal), standard error
5 of the bead replicates (BEAD_STERR), number of beads used (AVG_NBEADS) and
6 a detection score, which estimates the probability of a gene being detected above the
7 background. Resulting expression data was analyzed with R, using the lumi
8 package³. A variance stabilizing transformation⁴ was applied, followed by quantile
9 normalization to compensate for batch effects of the individual bead chips. For each
10 probe, the number of present calls over all samples was determined (the threshold
11 on the detection was $p < 0.01$), and probes absent in all samples were omitted in the
12 analysis. This omitted subset consisted of 18,396 probes. Hence, analysis was
13 performed for 28,927 probes. Differential expression was assessed with the limma
14 package of R⁵.

31 **Quantitative RT-PCR (qPCR)**

32 For qPCR analyses, 400ng of RNA was reverse transcribed with SuperScript III First
33 Strand Kit (Invitrogen) following the manufacturer's instructions, and qPCR was
34 performed in duplicates on a 7500Fast System (Applied Biosystems) using intron-
35 spanning PrimeTime qPCR Assays (Integrated DNA Technologies) listed in
36 Supplementary Table 2. Wherever possible, qPCR assays were selected that
37 covered the exon in which the Illumina Expression BeadChip probe was located.
38 Raw data was analyzed with SDS v1.4 (Applied Biosystems), and expression was
39 normalized within samples with the $\Delta\Delta$ CT method to reference gene *B2M*. Data was
40 expressed relative to the average expression of that gene in the healthy volunteers in
41 the dataset. Data points where duplicates differed by more than 1 CT were
42 discarded. Inter-run validity was verified by both processing and running previously
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 analyzed samples as internal controls and ensuring correct clustering within their
4
5 respective groups. Where necessary for normalization purposes, stored and
6
7 validated healthy volunteer samples were re-profiled along with samples from cohorts
8
9 IV and V.
10

14 Identification of a gene signature

16 For each pair-wise comparison between HV, P and PM, we evaluated all probes with
17
18 a moderated t-test, as implemented in the limma-package⁵ of R. P-values were
19
20 adjusted for multiple testing with Benjamini-Hochberg to control the false discovery
21
22 rate⁶. A probe was selected as being differentially expressed between two groups
23
24 when the adjusted p-value was smaller than 0.05 and the fold change exceeded 1.5
25
26 times up- or down-regulation ($\log_2 > 0.58$ or < -0.58 , respectively). For the
27
28 comparison between PM/P and HV, differential expression of the selected genes was
29
30 further validated with qPCR in 8 randomly selected individuals from each of the
31
32 groups in cohort I. The panel of 35 candidate genes derived from the 40 Illumina
33
34 probes differentially expressed in cohort I was augmented by 8 genes which
35
36 marginally missed the applied cutoff criteria and had been identified in unpublished *in*
37
38 *vitro* and *in vivo* screens during the pilot phase. Minimal sample size for further
39
40 cohorts was chosen to be 15 after conducting a statistical power analysis with the
41
42 data from cohort I to estimate the expected variation in gene expression. Sample size
43
44 was chosen to achieve a statistical power of 0.9 with an ordinary t-test when fold
45
46 changes of 1.5 are considered and 5% false positives are accepted. Power
47
48 calculations were done with the online tool from the Department of Bioinformatics
49
50 and Computational Biology of MD Anderson Cancer Center⁷. Differential expression
51
52 was considered to be confirmed by qPCR when the p-value after a two-tailed
53
54
55
56
57
58
59
60

1
2
3 unpaired t-test was smaller than 0.1 and/or the associated area under the ROC curve
4
5 (AUC) was larger than 0.7. as calculated with Prism (GraphPad, Inc.). We chose
6
7 deliberately for loose cut-offs on p-value and AUC for the confirmation, since less
8
9 distinctly differentially expressed genes could in theory still add value to a (later
10
11 developed) multiple-gene classification strategy.
12
13

14 15 16 **Multicentric validation study**

17
18 *Overview.* The diagnostic test consists of a gene panel assay in combination with
19
20 software for decision support. The software implements an algorithm that takes the
21
22 data from the assay as input and outputs a binary decision: whether the profiled
23
24 sample comes from a CRC patient or not. The algorithm is an ensemble method
25
26 (ENS)⁸ that consults 3 subroutines, then counts the number of votes in favor of CRC
27
28 and finally proposes the decision that is supported by at least 2 subroutines. The 3
29
30 subroutines form a heterogeneous set of alternative classification algorithms: an
31
32 easily interpretable ensemble stump classifier (SGMV – single gene majority vote), a
33
34 linear support vector machine (SVM) and a more complex random forest (RF). The
35
36 parameters of the 3 subroutines were fitted in parallel to a subset of samples from
37
38 the multi-centric cohort II. This training subset was constructed via stratified random
39
40 sampling. Performance of the algorithm was assessed through a Monte Carlo cross-
41
42 validation (MCCV) procedure on the training data and further validated on the
43
44 samples from cohort II that were excluded during training.
45
46
47
48

49
50 *Stratified random sampling.* We identified combinations of the four oncology centres
51
52 and two sample classes (i.e. HV or CRC) as 8 strata. From each stratum, we
53
54 sampled 2 times as much training samples as validation samples. The actual number
55
56 of samples per stratum was chosen so that i. there was no evidence of dependence
57
58
59
60

1
2
3 of class labeling on centre in either validation or training dataset, ii. the final datasets
4
5 were balanced (i.e. as much HV as CRC). Dependence between class labeling and
6
7 centre of origin was excluded by testing with a Fisher's exact test ($p > 0.93$). The
8
9 random split was performed prior to fitting parameters and retained for all further
10
11 analyses to obtain realistic measures of classification performance. Since our
12
13 subroutines required complete data, we imputed missing values after assembling the
14
15 training and validation datasets for each dataset separately using nearest neighbor
16
17 averaging, as implemented in the impute-package in R⁹.

20
21 *Subroutines.* The SGMV compares the expression value of each input gene first to a
22
23 gene-specific cut-off and then assigns a defined class to an unknown sample
24
25 depending on whether the cut-offs are exceeded for at least half of the genes (i.e.
26
27 majority vote). The SGMV parameters hence consist of gene-specific cut-offs. The
28
29 gene-specific cut-offs are fitted by taking that value that corresponds to the point
30
31 closest to the top-left corner of the gene-associated ROC curve, using the pROC-
32
33 package in R¹⁰. The SVM with linear kernel is similar to linear discriminant analysis,
34
35 taking as input the expression values of a set of genes and comparing a linear
36
37 combination of the input values to a threshold in order to assign a defined class to an
38
39 unknown sample, thereby giving higher weight to more informative genes. The SVM
40
41 parameters hence consist of gene-specific weights and one threshold. We fitted the
42
43 parameters with the kernlab-package in R¹¹. The RF pushes the expression values of
44
45 a set of genes through a multitude of decision trees (each looking at a random subset
46
47 of genes and built from a random subset of samples from the training data), notes
48
49 down for each class the proportion of supporting individual trees and finally assigns
50
51 the class with highest support. The RF parameters hence consist of individual
52
53 decision trees. We fitted the parameters with the randomForest-package in R¹².

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

16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32 *Performance metrics.* The classifiers were validated on the qPCR test dataset,
33 constructed from healthy volunteers and patients of multi-centric cohort II who were
34 not included during development of the models (see above). To verify the similarity of
35 the test set to the training set, a Spearman-correlation between all assays was
36 performed, ensuring that test assays did not cluster separately from training assays.
37 A separate clustering would have been an indication that the training dataset was not
38 representative for the test samples. Two types of performance were finally reported:
39 ranking performance and classification performance. Ranking performance is the
40 capability of an algorithm to give a higher score to an individual from class CRC than
41 to an individual from class HV. We measured ranking performance by the area under
42 the ROC curve (AUC). For all 4 routines (SGMV, SVM, RF and ENS), we provided
43 the AUC as well as the lower bound and upper bound of its 95% confidence interval,
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 as computed after 2,000 bootstraps with the pROC-package in R¹⁰. Classification
4
5 performance measures the capability of an algorithm to assign an individual to the
6
7 correct class. We reported for all routines the balanced error rate (BER), sensitivity
8
9 (Se) and specificity (Sp). For Se and Sp, we also computed the lower bound and
10
11 upper bound of the 95% confidence interval after 2,000 bootstraps.
12
13

14 15 16 **Complementary data analysis**

17
18 A complementary data analysis by an independent team (DNAlytics, Belgium) on the
19
20 same 23-marker signature led to the same conclusions in terms of
21
22 performances. Another (per-marker) normalization procedure has been proposed.
23
24 This normalization is applied on the log-transformed gene expression (i.e. Δ CT
25
26 values) and consists in computing, on the training set (for example Cohort II, both HV
27
28 and CRC), the mean and standard deviation of each marker. When a prediction has
29
30 to be made on a new, potentially isolated sample, each marker measurement of this
31
32 new sample is normalized by subtracting the corresponding mean, and by dividing by
33
34 the corresponding standard deviation. A modified procedure has also been proposed
35
36 for the imputation of missing values, making it dependent on the reference cohort
37
38 only. This avoids the need for a new reference HV batch as prediction has to be
39
40 made on a new (set of) sample(s).
41
42
43
44

45
46 The first experiment consisted in cross-validating a model on Cohort II (BER: 8.4%
47
48 [3.4%;13.4%]; AUC: 0.93 [0.88;0.98]). A second experiment consisted in learning the
49
50 same type of model on Cohort II and having it make predictions on Cohort III (BER:
51
52 13.2%; AUC: 0.92). All analyses were performed in R with scripts designed by
53
54 DNAlytics, fully independent from other analyses described in this paper.
55
56
57
58
59
60

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.

1
2
3 Expression data were represented as mean \pm SEM of the indicated number of
4
5 measurements. Statistical significance of differential expression was assessed with
6
7 Prism (GraphPad, Inc.) by two-tailed unpaired t-test (for two conditions) and ANOVA
8
9 followed by Bonferroni correction (for more than two conditions) after ensuring equal
10
11 variance using F test.
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
51
52
53
54
55
56
57
58
59
60

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 context^{14 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

1
2
3 shift. However, we could not identify key regulators of known pathways, probably due
4
5 to the limited information on reciprocal effects of PBM and tumour cells¹⁶. Though of
6
7 high interest with regards to the biological function, functional biological knowledge is
8
9 dispensable to exploit the full potential of the gene signature as a diagnostic tool in
10
11 analogy to other important clinical tests, which are devoid of a biological
12
13 understanding (e.g., prostate specific antigen, PSA, and pro-calcitonin, PCT).
14
15
16
17

18 **Supplementary Note 3**

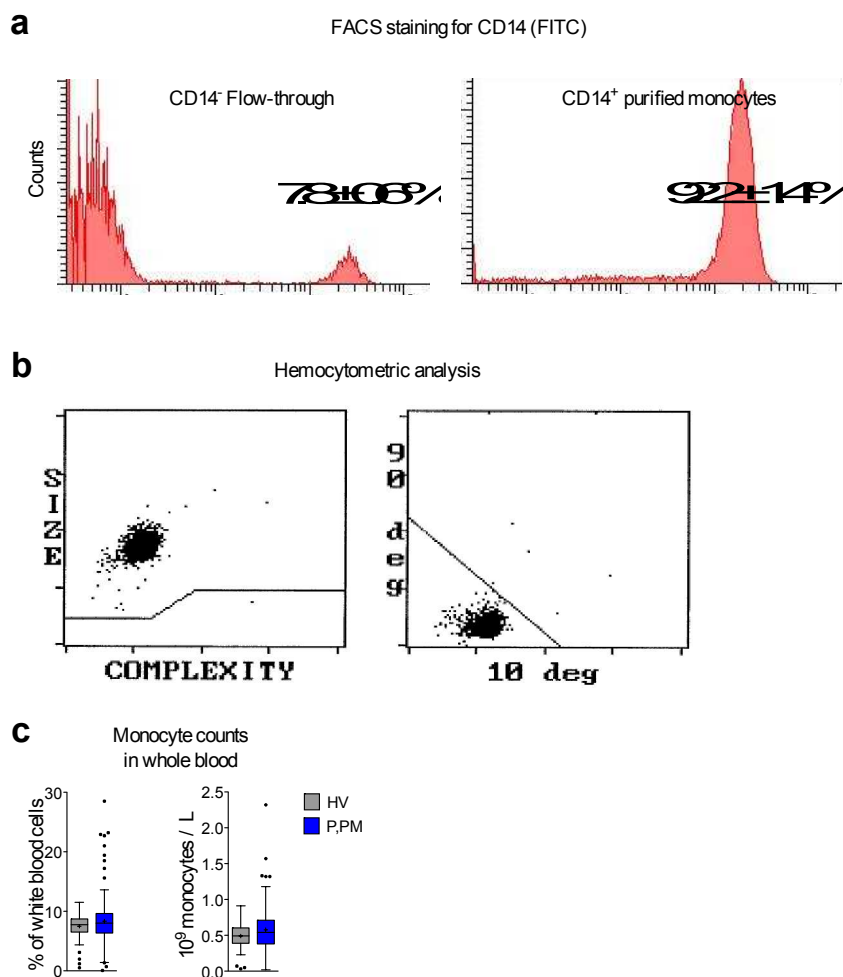
19
20 In accordance with our initial screening results, we found no differences in
21
22 expression patterns of P versus PM (data not shown). Moreover, as cumulating
23
24 evidence is suggesting subcategories of CRC according to its location¹⁷, we
25
26 investigated if the gene signature was capable of separating left versus right CRC or
27
28 colon versus rectal cancer, respectively. In line with the homogeneous clustering of
29
30 samples, we found no differences by location (AUC of 0.45 [0.20-0.73] for left versus
31
32 right CRC and AUC of 0.47 [0.28-0.70] for colon versus rectal cancer).
33
34
35
36
37

38 **Supplementary Note 4**

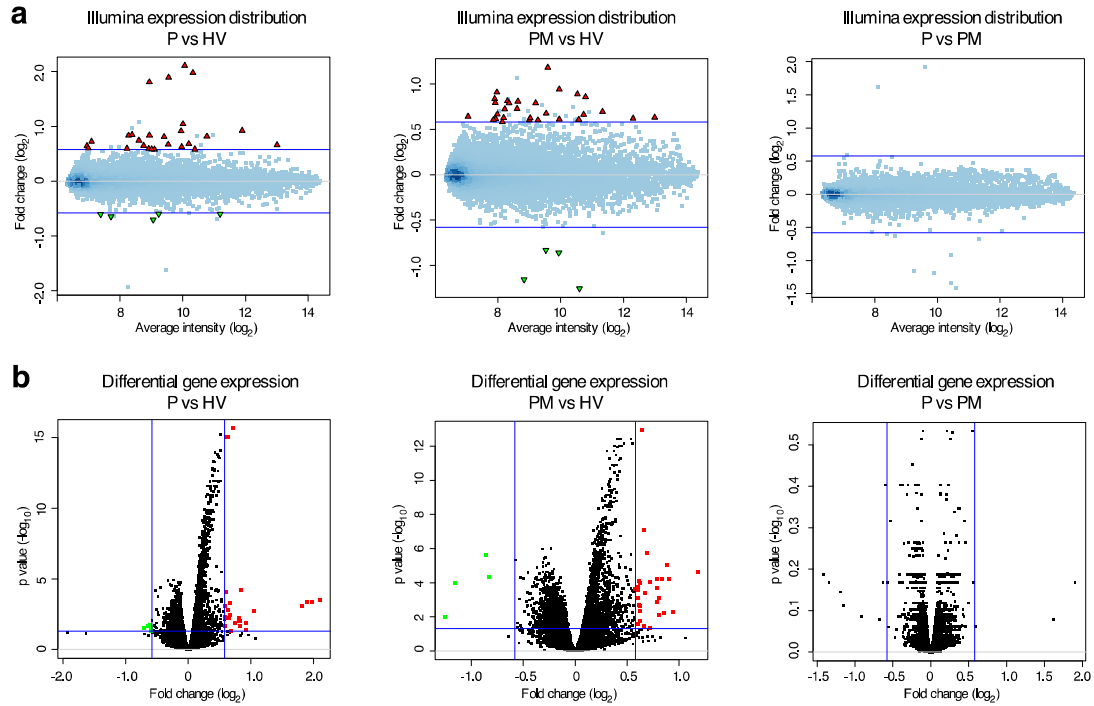
39
40 We sought to confirm our findings from the screening in independent samples by
41
42 independent techniques to rule out bias by the chosen technique and maximize
43
44 chances of extrapolation to other clinical centres. Our first step was a random re-
45
46 processing of collected samples and assessment by qPCR, which led to an initial
47
48 refinement of the gene signature, while some genes in this subset of samples
49
50 performed well even as single markers. By assessing Spearman correlation values
51
52 between expression data in the Illumina platform (used for screening) and the qPCR
53
54 technique (used for confirmation), we could rule out discrepancies in expression
55
56
57
58
59
60

1
2
3 between both analyses (Supplementary Figure 8). Consistently, a multicentric
4 validation trial revealed that the established gene signature retained the promising
5 performance observed in the screening phase, regardless of the centre and method
6 of analysis, while our multi-gene classification model allows to exploit the highest
7 informative content obtained from the expression analyses.
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
51
52
53
54
55
56
57
58
59
60

SUPPLEMENTARY FIGURES

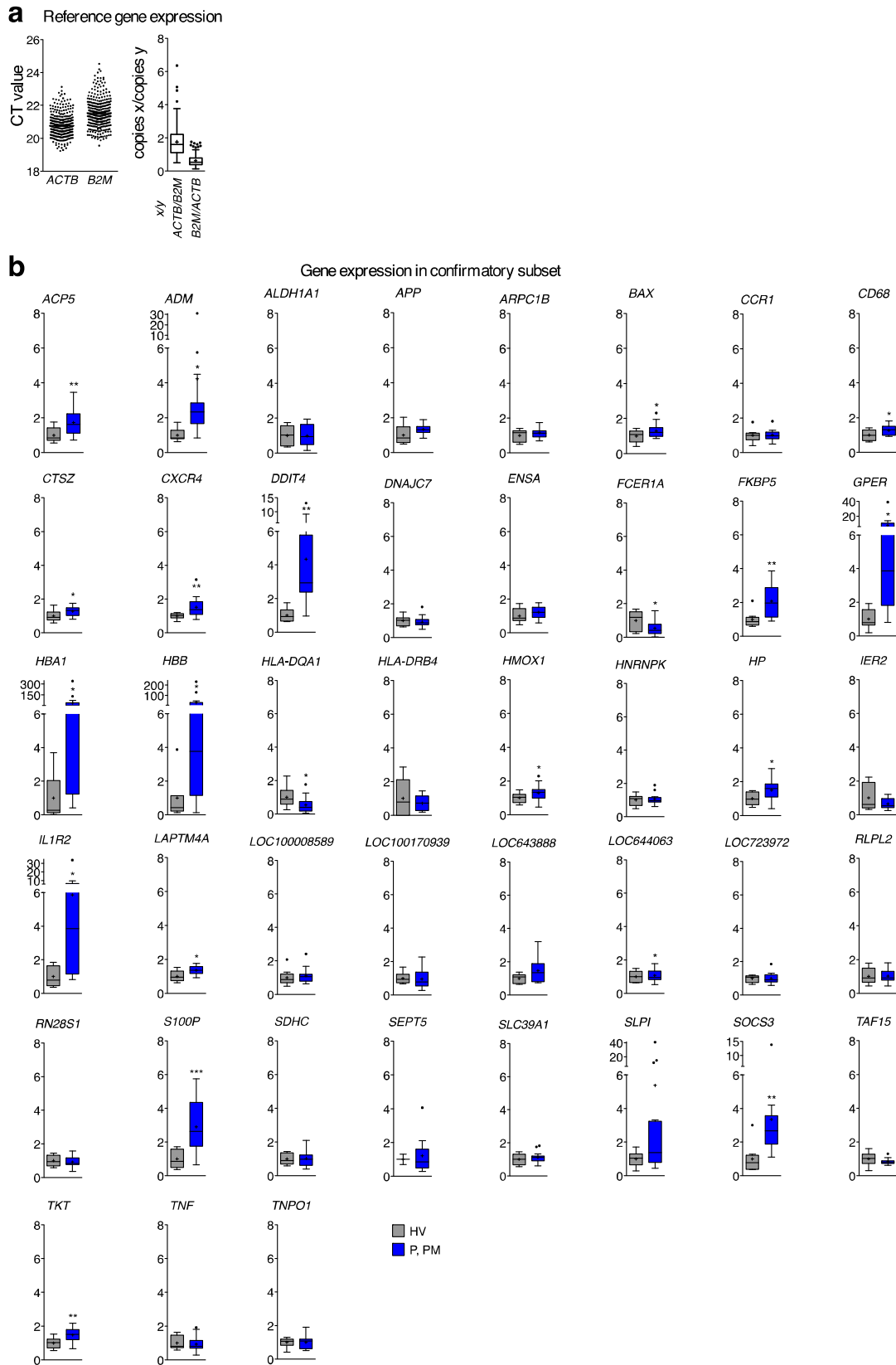
**Supplementary Figure 1: Isolation of PBM and monocyte counts**

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).



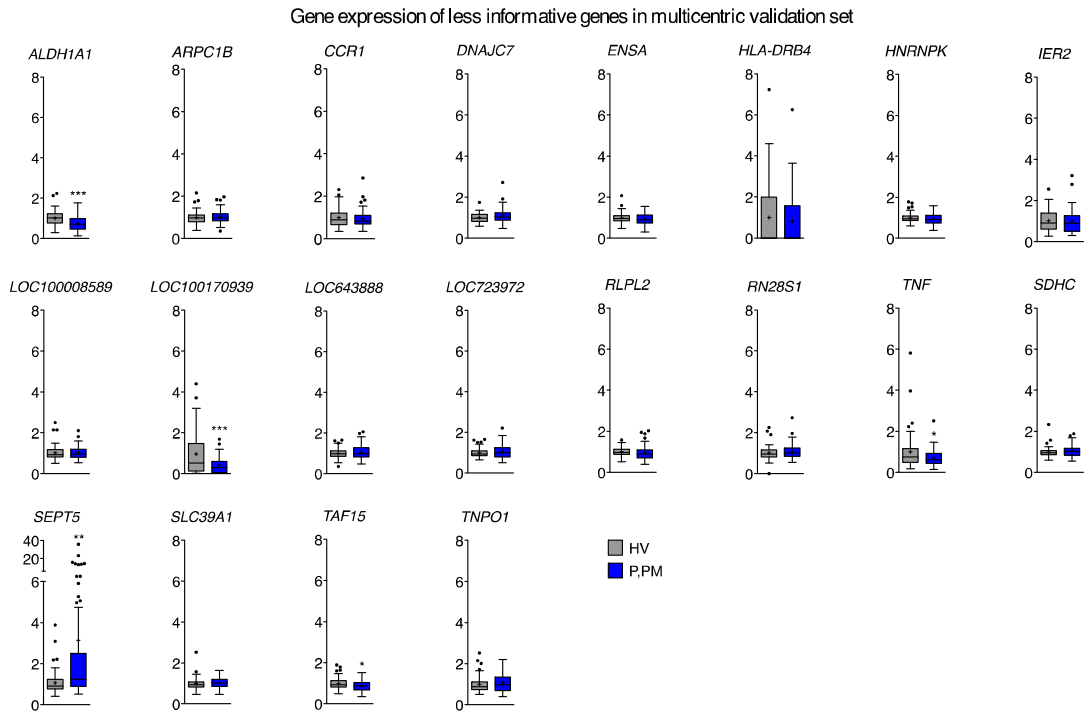
Supplementary Figure 2: Differentially expressed genes in PBM

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 $p < 0.05$.



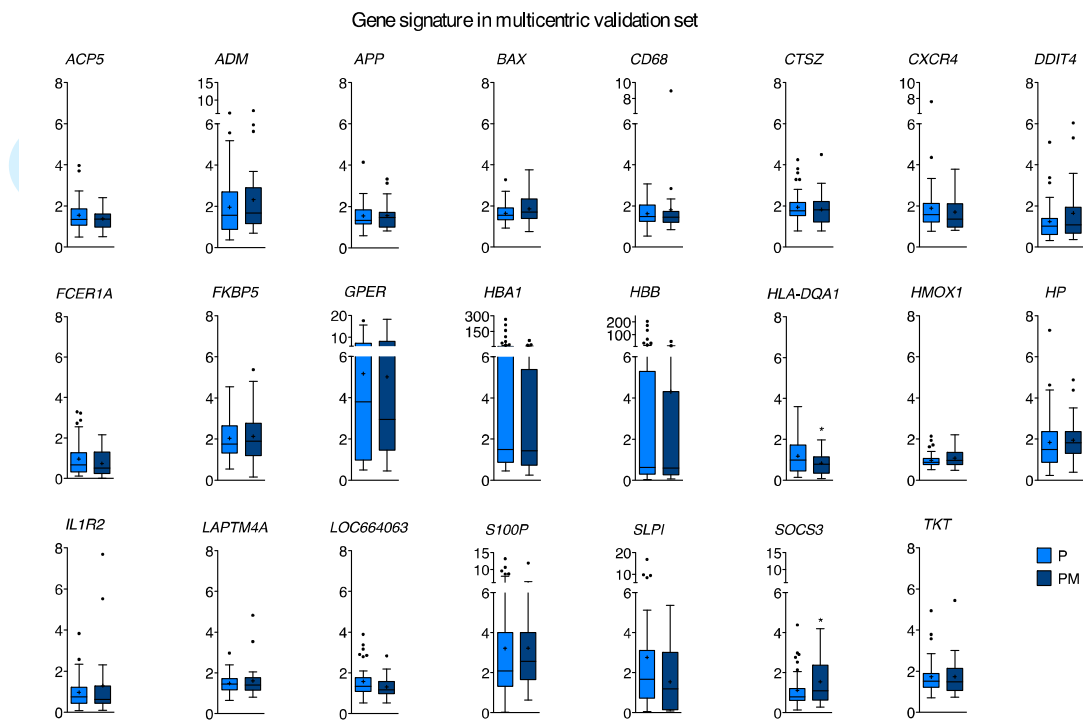
Supplementary Figure 3: Technical validation (subset of cohort I)

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; *, $p < 0.1$; **, $p < 0.01$; ***, $p < 0.001$.



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

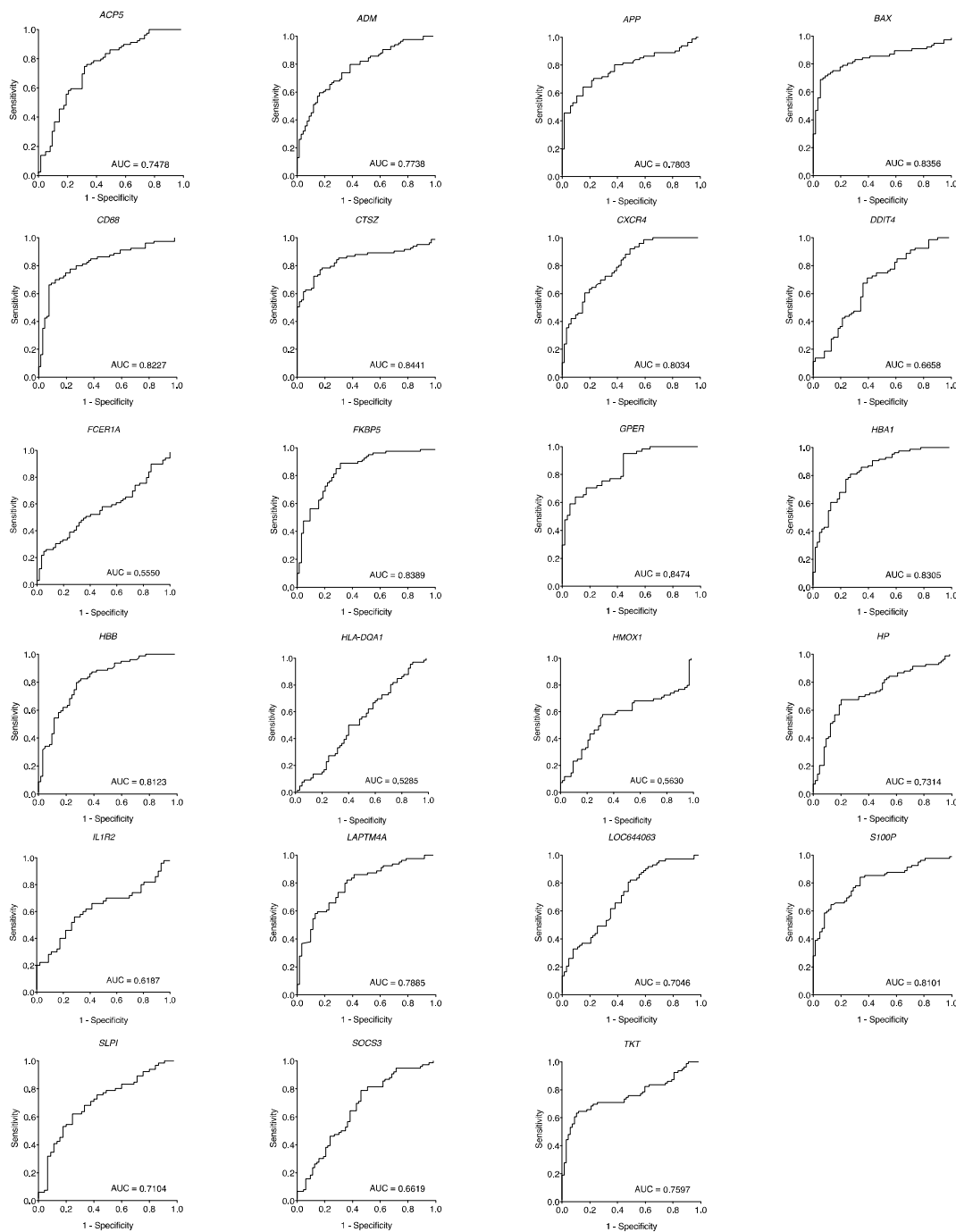
Expression levels are displayed as expression relative to the HV mean; boxes, first to third quartile; Whiskers, range; dots, values outside 1.5-times the interquartile distance; horizontal line, median; +, mean; *, p<0.05; **, p<0.01; ***, p<0.001.



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; *, $p < 0.05$; **, $p < 0.01$; ***, $p < 0.001$.

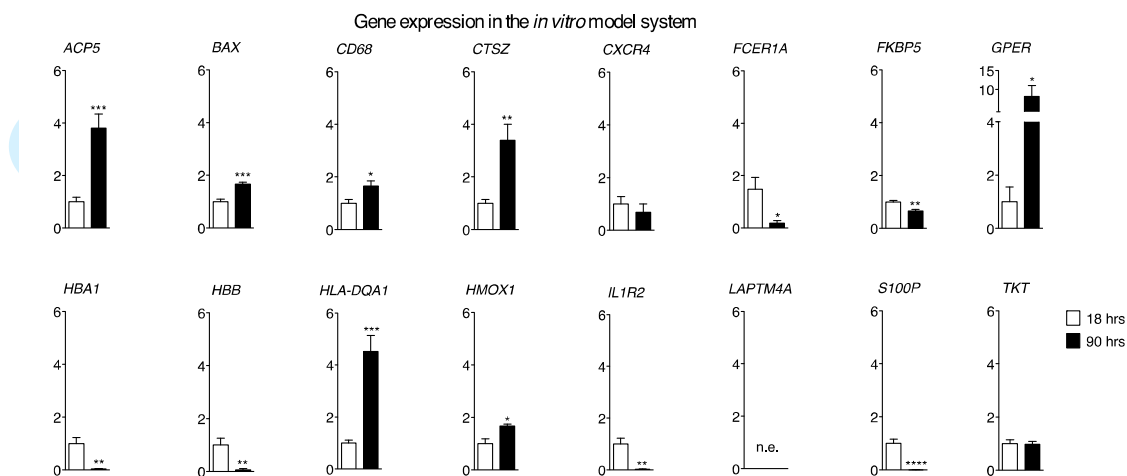
ROC analysis for individual signature genes in full multicentric validation set



Supplementary Figure 6: Single gene ROC analysis

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

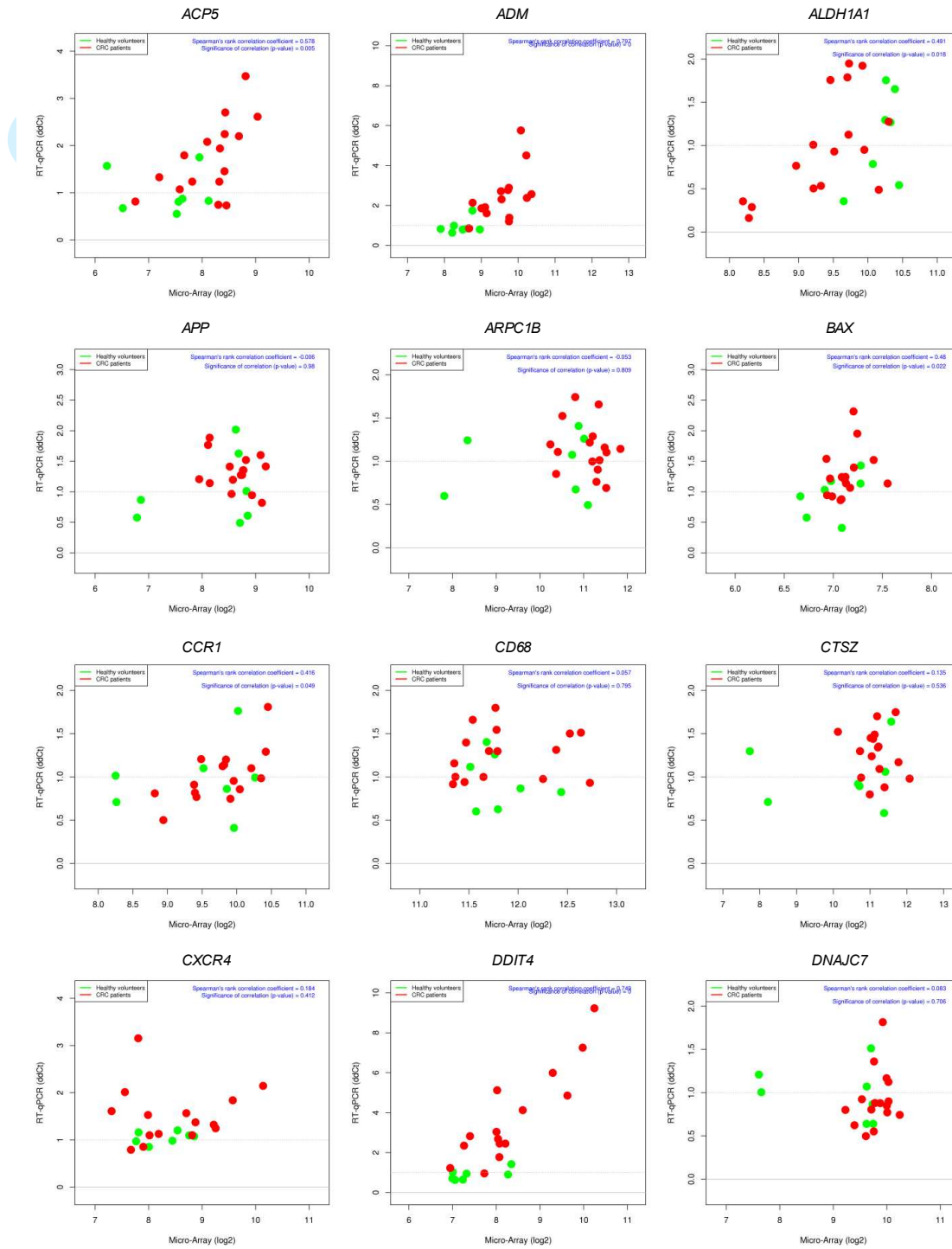




Supplementary Figure 7: Identification of putative markers in the *in vitro* model

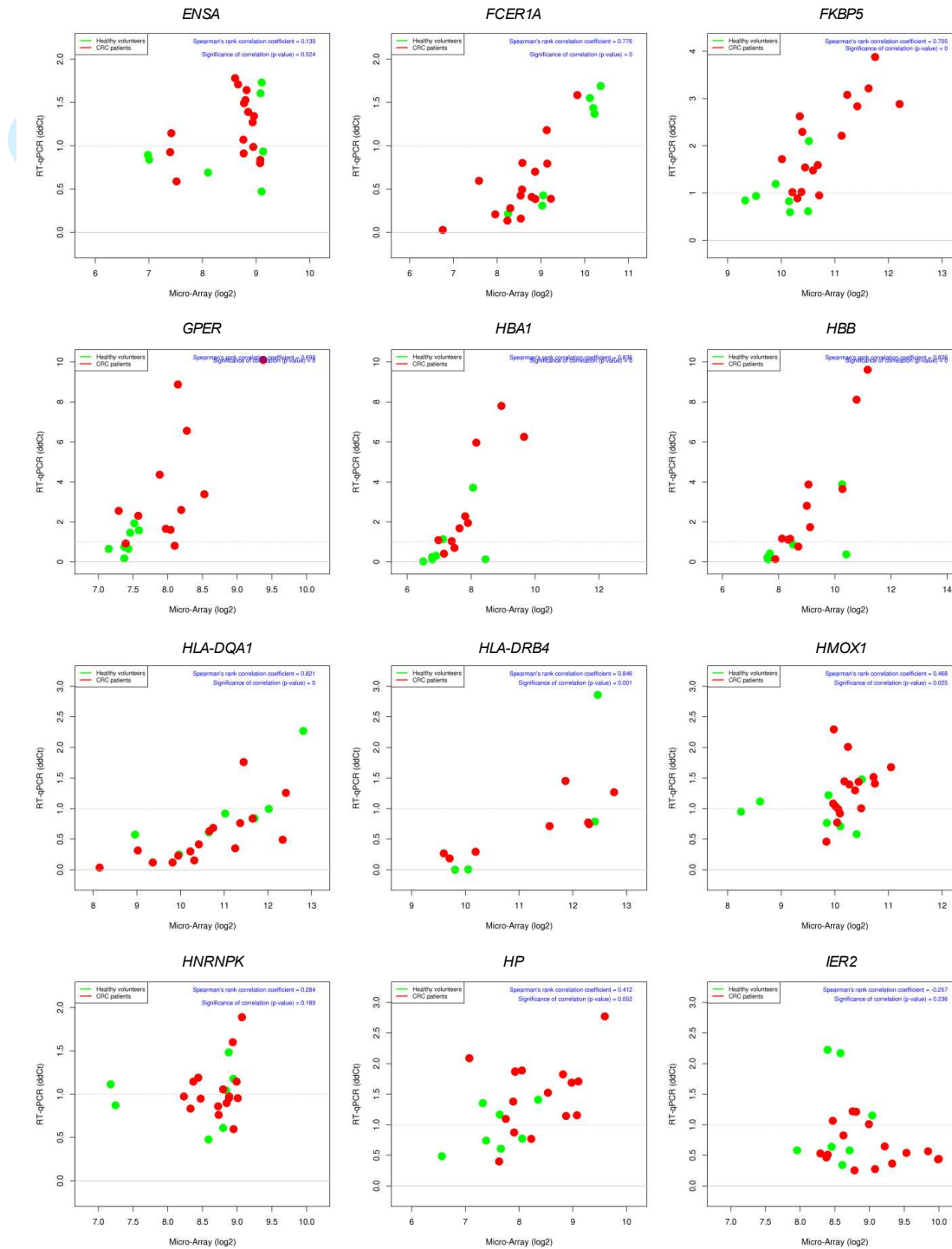
Shown are the expression levels of the 16 genes not selected out of the gene signature, which show altered expression levels in culture without any stimulus. Expression levels are shown as mean with SEM at 18 hours and 72 hours later (90 hours).

*, $p < 0.05$; **, $p < 0.01$; ***, $p < 0.001$; ****, $p < 0.0001$; n.e., not expressed *in vitro*.



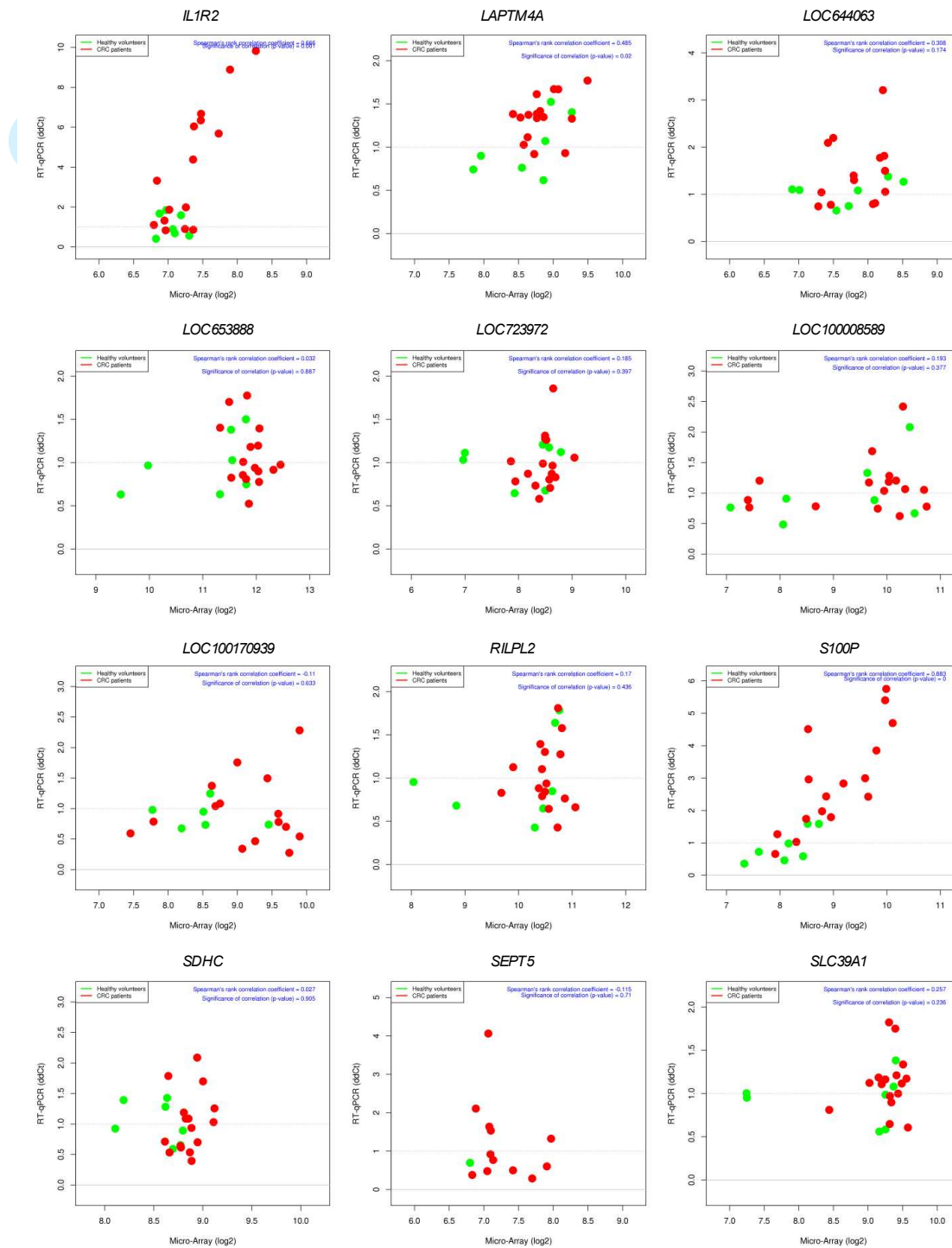
Supplementary Figure 8:

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.



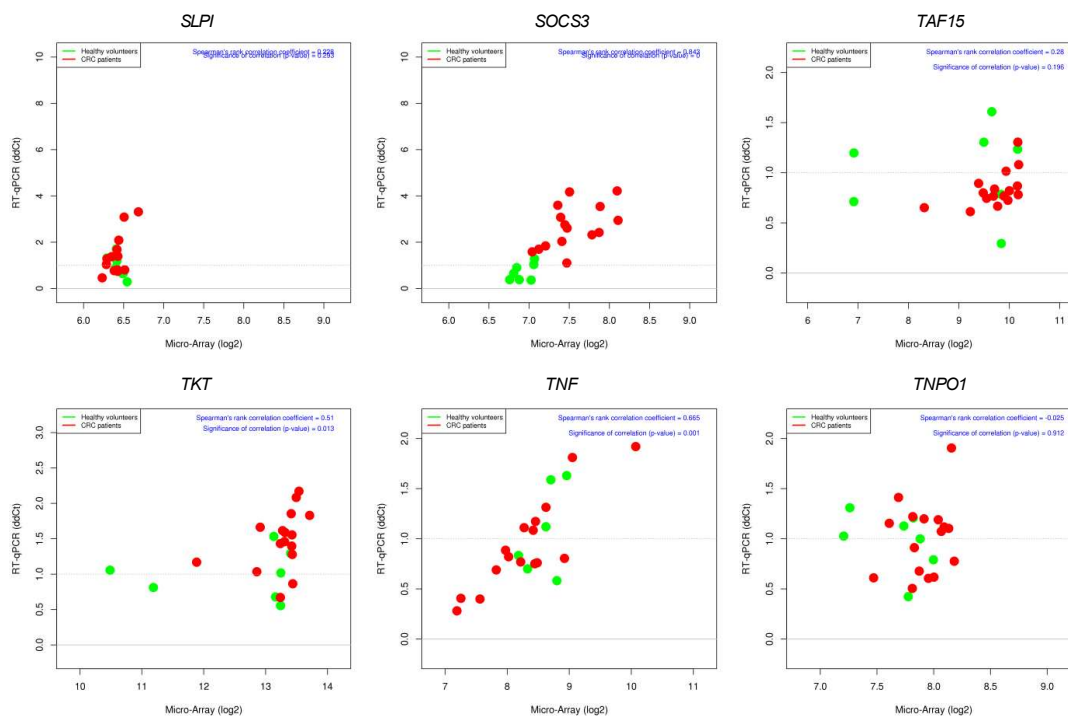
Supplementary Figure 8 – continued





Supplementary Figure 8 – continued





Supplementary Figure 8 – continued

SUPPLEMENTARY TABLES

Supplementary Table 1: IDT PrimeTime qPCR Assays

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

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

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
51
52
53
54
55
56
57
58
59
60

Supplementary Table 2: Dependence of class label on number of missing values (Fisher's exact test)

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

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

Genomewide Screening						Confirmation and Validation					
P,PM vs HV ^a			P vs HV			PM vs HV			P,PM vs. HV		
	Ratio	p		Ratio	p		Ratio	p		Ratio	p
<i>ADM</i>	2,00	<0,0001	<i>ADM</i>	1,75	0,0059	<i>ADM</i>	2,27	<0,0001	<i>ACP5^b</i>	1,61	<0,0001
<i>ALDH1A1</i>	0,66	0,0002	<i>CTSZ</i>	1,76	0,0103	<i>ALDH1A1</i>	0,56	<0,0001	<i>ADM</i>	2,16	<0,0001
<i>ARPC1B</i>	1,55	0,0209	<i>DDIT4</i>	1,78	0,0226	<i>AQP9</i>	1,62	<0,0001	<i>ALDH1A1</i>	0,88	<0,0001
<i>BAX</i>	1,50	0,0001	<i>DNAJC7</i>	1,59	0,0005	<i>BAX</i>	1,52	0,0008	<i>APP</i>	1,61	<0,0001
<i>CTSZ</i>	1,79	0,0007	<i>FCER1A</i>	0,61	0,0296	<i>CTSZ</i>	1,81	0,0056	<i>ARPC1B</i>	0,98	0,6497
<i>DDIT4</i>	1,71	0,0063	<i>HBA1</i>	3,51	0,0008	<i>DDIT4</i>	1,65	0,0477	<i>BAX</i>	1,76	<0,0001
<i>DNAJC7</i>	1,59	<0,0001	<i>HBA2</i>	4,31	0,0004	<i>DNAJC7</i>	1,60	0,0004	<i>CCR1</i>	0,90	0,3981
<i>FCER1A</i>	0,52	0,0002	<i>HBB</i>	3,95	0,0004	<i>DYSF</i>	1,52	0,0002	<i>CD68</i>	1,76	<0,0001
<i>FKBP5</i>	1,61	0,0001	<i>HMOX1</i>	1,55	0,0017	<i>FCER1A</i>	0,45	0,0001	<i>CTSZ</i>	1,96	<0,0001
<i>GPER</i>	1,58	0,0006	<i>HNRNPK</i>	1,60	0,0497	<i>FCGR1A</i>	1,52	0,0003	<i>CXCR4</i>	2,24	<0,0001
<i>HBA1</i>	2,33	0,0078	<i>HS.143909</i>	1,56	<0,0001	<i>FKBP5</i>	1,85	<0,0001	<i>DDIT4</i>	1,47	0,0025
<i>HBA2</i>	2,69	0,0051	<i>HS.581828</i>	1,52	<0,0001	<i>GPER</i>	1,78	0,0001	<i>DNAJC7</i>	1,07	0,1045
<i>HBB</i>	2,39	0,0099	<i>HS.61208</i>	1,65	<0,0001	<i>HLA-DRB6</i>	0,42	0,0102	<i>ENSA</i>	0,89	0,1122
<i>HMOX1</i>	1,54	0,0001	<i>IER3</i>	1,50	0,0009	<i>HMOX1</i>	1,53	0,0020	<i>FCER1A</i>	0,97	0,7541
<i>HNRNPK</i>	1,58	0,0125	<i>LOC100008589</i>	1,68	0,0131	<i>HP</i>	1,75	0,0080	<i>FKBP5</i>	2,45	<0,0001
<i>HP</i>	1,54	0,0131	<i>LOC100128274</i>	0,66	0,0195	<i>HS.61208</i>	1,56	<0,0001	<i>GPER</i>	5,29	<0,0001
<i>HS.143909</i>	1,51	<0,0001	<i>LOC100130707</i>	1,51	0,0232	<i>LOC100170939</i>	1,65	0,0001	<i>HBA1</i>	15,07	0,0165
<i>HS.61208</i>	1,60	<0,0001	<i>LOC100132394</i>	1,79	0,0095	<i>LOC100190986</i>	1,53	0,0001	<i>HBB</i>	11,96	0,0281
<i>IL1R2</i>	1,50	0,0482	<i>LOC100132727</i>	0,66	0,0282	<i>LOC153561</i>	1,73	0,0001	<i>HLA-DQA1</i>	1,01	0,8425
<i>LOC100008589</i>	1,55	0,0079	<i>LOC100134364</i>	1,57	0,0057	<i>LOC441087</i>	1,54	0,0177	<i>HLA-DRB4</i>	0,77	0,3931
<i>LOC100129685</i>	1,71	0,0356	<i>LOC153561</i>	1,50	0,0049	<i>RNF146</i>	1,50	0,0002	<i>HMOX1</i>	0,95	0,7338
<i>LOC100132394</i>	1,65	0,0045	<i>LOC649143</i>	1,90	0,0133	<i>S100P</i>	1,75	0,0007	<i>HNRNPK</i>	0,92	0,2280
<i>LOC100134364</i>	1,53	0,0009	<i>LOC723972</i>	1,51	0,0001	<i>SEPT5</i>	1,59	0,0347	<i>HP</i>	1,92	<0,0001
<i>LOC100170939</i>	1,54	<0,0001	<i>LOC728755</i>	0,64	0,0210	<i>SLC39A1</i>	1,54	0,0025	<i>IER2</i>	0,97	0,8782
<i>LOC153561</i>	1,61	<0,0001	<i>SLC39A1</i>	1,50	0,0058	<i>SOCS3</i>	1,73	0,0014	<i>IL1R2</i>	0,86	0,4209

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

Genomewide Screening						Confirmation and Validation					
P,PM vs HV ^a			P vs HV			PM vs HV			P,PM vs. HV		
	Ratio	p		Ratio	p		Ratio	p		Ratio	p
<i>LOC649143</i>	1,56	0,0356	<i>TAF15</i>	2,06	0,0001	<i>TAF15</i>	1,73	0,0002	<i>LAPTM4A</i>	1,59	<0,0001
<i>LOC653156</i>	1,73	0,0443	<i>TKT</i>	1,58	0,0034	<i>TKT</i>	1,55	0,0048	<i>LOC100008589</i>	0,99	0,9313
<i>LOC653737</i>	1,86	0,0472	<i>ZNF223</i>	0,66	0,0478	<i>TNPO1</i>	1,55	0,0001	<i>LOC100170939</i>	1,09	0,0004
<i>LOC728755</i>	0,66	0,0066				<i>UPP1</i>	1,58	<0,0001	<i>LOC643888</i>	1,05	0,2617
<i>S100P</i>	1,53	0,0020				<i>ZBTB16</i>	1,52	0,0252	<i>LOC644063</i>	1,54	<0,0001
<i>SEPT5</i>	1,57	0,0094							<i>LOC723972</i>	1,03	0,1904
<i>SLC39A1</i>	1,52	0,0003							<i>RLPL2</i>	0,95	0,3618
<i>SOCS3</i>	1,51	0,0043							<i>RN28S1</i>	1,03	0,3003
<i>TAF15</i>	1,76	<0,0001							<i>S100P</i>	3,35	<0,0001
<i>TKT</i>	1,56	0,0003							<i>SDHC</i>	1,04	0,3017
									<i>SEPT5</i>	3,47	0,0020
									<i>SLC39A1</i>	1,02	0,3827
									<i>SLPI</i>	15,76	0,0090
									<i>SOCS3</i>	1,60	0,0158
									<i>TAF15</i>	0,84	0,0154
									<i>TKT</i>	1,79	<0,0001
									<i>TNF</i>	0,75	0,0205
									<i>TNPO1</i>	1,02	0,3165

^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).

Supplementary Table 4: Confirmation in random subset of cohort I

	Mean expression ^a	Fold ratio ^b	p	AUC
<i>ACP5</i>	81,336	1,73	0.0081^c	0.79
<i>ADM</i>	73,107	4,23	0.0941	0.95
<i>ALDH1A1</i>	47,649	0,99	0.9624	0.51
<i>APP</i>	176,332	1,32	0.1576	0.73
<i>ARPC1B</i>	1873,873	1,15	0.3336	0.57
<i>BAX</i>	11,474	1,29	0.0978	0.67
<i>CCR1</i>	338,797	1,01	0.9292	0.52
<i>CD68</i>	1821,912	1,27	0.0640	0.76
<i>CTSZ</i>	1580,418	1,28	0.0637	0.76
<i>CXCR4</i>	508,754	1,52	0.0065	0.84
<i>DDIT4</i>	36,963	4,34	0.0010	0.96
<i>DNAJC7</i>	105,494	0,92	0.5431	0.62
<i>ENSA</i>	250,287	1,21	0.2580	0.67
<i>FCER1A</i>	410,118	0,54	0.0768	0.73
<i>FKBP5</i>	39,131	2,08	0.0013	0.89
<i>GPER</i>	1,875	7,59	0.0138	0.93
<i>HBA1</i>	5243,339	41,40	0.0861	0.86
<i>HBB</i>	440,188	31,10	0.0773	0.85
<i>HLA-DQ1</i>	1748,345	0,53	0.0918	0.77
<i>HLA-DRB4</i>	1135,072	0,71	0.6316	0.53
<i>HMOX1</i>	405,989	1,30	0.0729	0.70
<i>HNRNPK</i>	1648,472	1,05	0.7356	0.52
<i>HP</i>	129,478	1,50	0.0218	0.76
<i>IER2</i>	0,657	0,65	0.2556	0.63
<i>IL1R2</i>	4,794	5,85	0.0288	0.87
<i>LAPTM4A</i>	338,391	1,35	0.0206	0.78
<i>LOC100008589</i>	18500877,250	1,12	0.5782	0.63
<i>LOC100170939</i>	252,768	0,96	0.8506	0.56
<i>LOC643888</i>	308,104	1,46	0.6143	0.55
<i>LOC644063</i>	1504,972	1,07	0.0415	0.71
<i>LOC723972</i>	568,957	1,00	0.9645	0.56
<i>RLPL2</i>	186,476	1,02	0.9078	0.53
<i>RN28S1</i>	14924567,167	0,92	0.5908	0.58
<i>S100P</i>	8,494	2,90	0.0003	0.91
<i>SDHC</i>	313,373	1,02	0.9145	0.53
<i>SEPT5</i>	0,298	1,22	0.6497	0.50
<i>SLC39A1</i>	166,812	1,12	0.4069	0.62
<i>SLPI</i>	1,656	5,39	0.2477	0.71
<i>SOCS3</i>	128,894	3,36	0.0081	0.91
<i>TAF15</i>	327,194	0,83	0.3084	0.65
<i>TKT</i>	986,111	1,48	0.0061	0.82
<i>TNF</i>	28,056	0,94	0.7325	0.53
<i>TNPO1</i>	140,447	1,00	0.9722	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

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

Gene	Full Name	Biological Function	Potential Function in Monocytes
<i>ACP5</i>	acid phosphatase 5, tartrate resistant	iron containing glycoprotein involved in adhesion and migration	negative regulation of inflammatory response in interleukin pathways
<i>ADM</i>	adrenomedullin	vasodilation, regulation of hormone secretion, promotion of angiogenesis	antimicrobial activity, wound healing
<i>APP</i>	amyloid beta (A4) precursor protein	protein basis of amyloid plaques in Alzheimer disease	antimicrobial activity, mitotic activity
<i>BAX</i>	BCL2-associated X protein	p53-mediated activator of apoptosis	myeloid cell homeostasis
<i>CD68</i>	CD68 molecule	integral membran glycoprotein of scavenger receptor family	highly expressed on monocytes and macrophages, mediator of recruitment and activation
<i>CTSZ</i>	cathepsin Z	lysosomal cystein proteinase, involved in migration and adhesion	unknown
<i>CXCR4</i>	chemokine (C-X-C motif) receptor 4	CXC chemokine receptor specific for stromal cell-derived factor-1	mediator of recruitment, chemotaxis, and activation
<i>DDIT4</i>	DNA-damage-inducible transcript 4	negative regulation of mTOR signalling upon cellular stress	defense response to microbial signals
<i>FCER1A</i>	Fc fragment of IgE, high affinity I, receptor for; alpha polypeptide	alpha subunit of IgE-mediated allergic response	positive regulation of type-I immune response and macrophage differentiation
<i>FKBP5</i>	FK506 binding protein 5	member of immunophilin protein family, immunoregulation	receptor for FK506 and rapamycin, mediating calcineurin inhibition
<i>GPER</i>	G protein-coupled estrogen receptor 1	non-genomic signalling of estrogen stimulus	negative regulator of leukocyte activation; innate immune response
<i>HBA1</i>	hemoglobin, alpha 1	alpha chain of hemoglobin	unknown
<i>HBB</i>	hemoglobin, beta	beta chain of hemoglobin	positive regulation of nitric oxide synthesis,
<i>HLA-DQA1</i>	major histocompatibility complex, class II, DQ alpha 1	MHC class II receptor activity; peptide antigen binding	antigen processing and presentation
<i>HMOX1</i>	heme oxygenase (decycling) 1	heme catabolism	regulation of phagocytosis and migration, chemokine synthesis, wound healing, and angiogenesis
<i>HP</i>	haptoglobin	preproprotein of haptoglobin subunit	acute-phase defense response
<i>IL1R2</i>	interleukin 1 receptor, type II	cytokine receptor for IL-1	cytokine-mediated immune response
<i>LAPTM4A</i>	lysosomal protein transmembrane 4 alpha	unknown	unknown
<i>LOC644063</i>	heterogeneous nuclear ribonucleoprotein K pseudogene 4	unknown	unknown
<i>S100P</i>	S100 calcium binding protein P	cell cycle progression and differentiation	unknown
<i>SLPI</i>	secretory leukocyte peptidase inhibitor	secreted inhibitor of serin proteinases	negative regulation of endopeptidase activity
<i>SOCS3</i>	suppressor of cytokine signaling 3	negative regulator of cytokine signalling	modulator of immune response, particularly IFN-γ mediated
<i>TKT</i>	transketolase	enzyme of pentose phosphate pathway	metabolic modulator

SUPPLEMENTARY REFERENCES

1. Weitz J, Koch M, Debus J, et al. Colorectal cancer. *Lancet* 2005;**365**(9454):153-65.
2. 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.
3. 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.
6. 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. <http://bioinformatics.mdanderson.org/MicroarraySampleSize/>.
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.
11. Burges CJC. A Tutorial on Support Vector Machines for Pattern Recognition. *Data Min Knowl Discov* 1998;**2**(2):121-67.

- 1
2
3 12. Liaw A, Wiener M. Classification and Regression by randomForest. R News: The
4
5 Newsletter of the R Project 2002;**2**(3):18-22.
6
7 13. Ambroise C, McLachlan GJ. Selection bias in gene extraction on the basis of
8
9 microarray gene-expression data. Proceedings of the National Academy of
10
11 Sciences of the United States of America 2002;**99**(10):6562-6.
12
13 14. Piehler A, Grimholt R, Ovstebo R, et al. Gene expression results in
14
15 lipopolysaccharide-stimulated monocytes depend significantly on the choice of
16
17 reference genes. BMC Immunology 2010;**11**(1):21.
18
19 15. Guo C, Liu S, Wang J, et al. ACTB in cancer. Clinica chimica acta; international
20
21 journal of clinical chemistry 2013;**417**:39-44.
22
23 16. Khatri P, Sirota M, Butte AJ. Ten Years of Pathway Analysis: Current Approaches
24
25 and Outstanding Challenges. PLoS Comput Biol 2012;**8**(2):e1002375.
26
27 17. Jess P, Hansen IO, Gamborg M, et al. A nationwide Danish cohort study
28
29 challenging the categorisation into right-sided and left-sided colon cancer.
30
31 BMJ open 2013;**3**(5).
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

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

CONTENTS

Supplementary Methods	Page 3
Supplementary Notes	Page 14
Supplementary Figures	Page 17
Supplementary Tables	Page 28
Supplementary References	Page 35

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
51
52
53
54
55
56
57
58
59
60

Confidential: For Review Only

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

1
2
3 gut flora. Samples were excluded from further analysis when final histology of the
4 surgical specimen did not confirm adenocarcinoma of the large intestine (assessed
5 by board-certified pathologists within clinical routine procedures).

6
7
8
9 The following four oncological centres contributed samples to this study: Digestive
10 Oncology, University Hospitals Leuven and Department of Oncology, KU Leuven,
11 Leuven, Belgium; Department of General, Visceral, and Transplantation Surgery,
12 University of Heidelberg, Heidelberg, Germany; Department of Oncology, San Filippo
13 Neri, Rome, Italy; Medical Oncology Clinic, Institut Jules Bordet, Brussels, Belgium.
14
15 The responsible scientists in each centre (1-2 per centre) were trained in the protocol
16 for isolation of PBM to ensure uniformity of the procedure. All participants gave
17 written informed consent, and the study was approved by the respective institutional
18 review boards (Leuven: B322201215873, Brussels: CE1950, Heidelberg: 323/2004,
19 Rome: 319/51). No adverse events from blood collection or colonoscopy were
20 recorded in included participants.

31 32 **Isolation of PBM**

33
34
35
36
37
38 20ml of EDTA-anticoagulated peripheral venous blood was collected following clinical
39 routine procedure, stored at 4°C and processed within 2 hours of blood collection.
40
41 For further isolation, blood was diluted 1:2 with DPBS (free of Ca²⁺ and Mg²⁺) and
42 layered carefully on Lymphoprep (Axis-Shield) in two separate tubes. All blood
43 collection and isolation steps were performed identical for samples of all origin.
44
45 Density gradient centrifugation was performed at 1,200g for 20 minutes at low
46 acceleration and no brake. Samples with macroscopically visible hemolysis were
47 excluded from further analysis. The PBMC interphase was collected carefully and
48
49 washed twice for 12 minutes at 250g and 175g with PBS. Hemocytometric analysis
50
51
52
53
54
55
56
57
58
59
60

1
2
3 was performed to ensure purity of PBMCs, and the pellet was pooled for further
4
5 processing and washed once for 10 minutes at 300g. Cells were then incubated with
6
7 CD14 magnetically-conjugated beads (BD) for 15 minutes at 4°C, washed 10
8
9 minutes at 300g and positively separated with the MACS system (Miltenyi) following
10
11 the manufacturer's instructions. The CD14+ fraction was flushed out and washed
12
13 once 10 minutes at 300g. Purity was assessed by FACS analysis for CD14 in the
14
15 pilot phase and by hemocytometric analysis (CellDyn 3700, Abbott) in every further
16
17 sample. Only samples with purity of >90% and viability >95% (assessed by Trypan
18
19 Blue staining) were retained for further analysis. Cell pellets were lysed in Buffer RLT
20
21 (Qiagen) at 10^6 monocytes in 350 μ l of Buffer RLT and stored at -80°C. For each
22
23 respective expression study, all samples were extracted simultaneously with the
24
25 RNeasy Mini Kit (Qiagen) following the manufacturer's instructions. Quality control
26
27 was performed by checking RNA quality on the Nanodrop system, and RNA integrity
28
29 was checked for microarray samples on the Agilent Bio-Analyzer. Only samples with
30
31 an extinction fraction 260/280 > 1.8 and 260/230 > 1.5, and an RNA integrity index of
32
33 >6 were retained for further analysis.
34
35
36
37
38
39
40

41 **Genome-wide expression analysis**

42
43 For genome-wide expression analysis, RNA was amplified and biotinylated using
44
45 Illumina TotalPrep RNA Amplification Kit (Ambion) following the manufacturer's
46
47 instructions to obtain biotinylated cRNA, which was hybridized to Illumina HumanHT-
48
49 12 v4 Expression BeadChips (Illumina) with the Illumina Whole-Genome Gene
50
51 Expression Direct Hybridization Assay (Illumina) following the manufacturer's
52
53 instructions. The Illumina HumanHT-12 v4 Expression BeadChip Kit contains 47,323
54
55 probes and 887 controls. After scanning, background-corrected expression values
56
57
58
59
60

1
2
3 and detection scores were extracted with GenomeStudio GX (version 1.5.4). For
4
5 each array, we used the summarized expression level (AVG_Signal), standard error
6
7 of the bead replicates (BEAD_STERR), number of beads used (AVG_NBEADS) and
8
9 a detection score, which estimates the probability of a gene being detected above the
10
11 background. Resulting expression data was analyzed with R, using the lumi
12
13 package³. A variance stabilizing transformation⁴ was applied, followed by quantile
14
15 normalization to compensate for batch effects of the individual bead chips. For each
16
17 probe, the number of present calls over all samples was determined (the threshold
18
19 on the detection was $p < 0.01$), and probes absent in all samples were omitted in the
20
21 analysis. This omitted subset consisted of 18,396 probes. Hence, analysis was
22
23 performed for 28,927 probes. Differential expression was assessed with the limma
24
25 package of R⁵.
26
27
28
29
30
31

32 **Quantitative RT-PCR (qPCR)**

33
34 For qPCR analyses, 400ng of RNA was reverse transcribed with SuperScript III First
35
36 Strand Kit (Invitrogen) following the manufacturer's instructions, and qPCR was
37
38 performed in duplicates on a 7500Fast System (Applied Biosystems) using intron-
39
40 spanning PrimeTime qPCR Assays (Integrated DNA Technologies) listed in
41
42 Supplementary Table 2. Wherever possible, qPCR assays were selected that
43
44 covered the exon in which the Illumina Expression BeadChip probe was located.
45
46 Raw data was analyzed with SDS v1.4 (Applied Biosystems), and expression was
47
48 normalized within samples with the $\Delta\Delta$ CT method to reference gene *B2M*. Data was
49
50 expressed relative to the average expression of that gene in the healthy volunteers in
51
52 the dataset. Data points where duplicates differed by more than 1 CT were
53
54 discarded. Inter-run validity was verified by both processing and running previously
55
56
57
58
59
60

1
2
3 analyzed samples as internal controls and ensuring correct clustering within their
4
5 respective groups. Where necessary for normalization purposes, stored and
6
7 validated healthy volunteer samples were re-profiled along with samples from cohorts
8
9 IV and V.
10

11 **Identification of a gene signature**

12
13
14 For each pair-wise comparison between HV, P and PM, we evaluated all probes with
15
16 a moderated t-test, as implemented in the limma-package⁵ of R. P-values were
17
18 adjusted for multiple testing with Benjamini-Hochberg to control the false discovery
19
20 rate⁶. A probe was selected as being differentially expressed between two groups
21
22 when the adjusted p-value was smaller than 0.05 and the fold change exceeded 1.5
23
24 times up- or down-regulation ($\log_2 > 0.58$ or < -0.58 , respectively). For the
25
26 comparison between PM/P and HV, differential expression of the selected genes was
27
28 further validated with qPCR in 8 randomly selected individuals from each of the
29
30 groups in cohort I. The panel of 35 candidate genes derived from the 40 Illumina
31
32 probes differentially expressed in cohort I was augmented by 8 genes which
33
34 marginally missed the applied cutoff criteria and had been identified in unpublished *in*
35
36 *vitro* and *in vivo* screens during the pilot phase. Minimal sample size for further
37
38 cohorts was chosen to be 15 after conducting a statistical power analysis with the
39
40 data from cohort I to estimate the expected variation in gene expression. Sample size
41
42 was chosen to achieve a statistical power of 0.9 with an ordinary t-test when fold
43
44 changes of 1.5 are considered and 5% false positives are accepted. Power
45
46 calculations were done with the online tool from the Department of Bioinformatics
47
48 and Computational Biology of MD Anderson Cancer Center⁷. Differential expression
49
50 was considered to be confirmed by qPCR when the p-value after a two-tailed
51
52
53
54
55
56
57
58
59
60

1
2
3 unpaired t-test was smaller than 0.1 and/or the associated area under the ROC curve
4
5 (AUC) was larger than 0.7. as calculated with Prism (GraphPad, Inc.). We chose
6
7 deliberately for loose cut-offs on p-value and AUC for the confirmation, since less
8
9 distinctly differentially expressed genes could in theory still add value to a (later
10
11 developed) multiple-gene classification strategy.
12
13

14 15 16 **Multicentric validation study**

17
18 *Overview.* The diagnostic test consists of a gene panel assay in combination with
19
20 software for decision support. The software implements an algorithm that takes the
21
22 data from the assay as input and outputs a binary decision: whether the profiled
23
24 sample comes from a CRC patient or not. The algorithm is an ensemble method
25
26 (ENS)⁸ that consults 3 subroutines, then counts the number of votes in favor of CRC
27
28 and finally proposes the decision that is supported by at least 2 subroutines. The 3
29
30 subroutines form a heterogeneous set of alternative classification algorithms: an
31
32 easily interpretable ensemble stump classifier (SGMV – single gene majority vote), a
33
34 linear support vector machine (SVM) and a more complex random forest (RF). The
35
36 parameters of the 3 subroutines were fitted in parallel to a subset of samples from
37
38 the multi-centric cohort II. This training subset was constructed via stratified random
39
40 sampling. Performance of the algorithm was assessed through a Monte Carlo cross-
41
42 validation (MCCV) procedure on the training data and further validated on the
43
44 samples from cohort II that were excluded during training.
45
46
47
48

49
50 *Stratified random sampling.* We identified combinations of the four oncology centres
51
52 and two sample classes (i.e. HV or CRC) as 8 strata. From each stratum, we
53
54 sampled 2 times as much training samples as validation samples. The actual number
55
56 of samples per stratum was chosen so that i. there was no evidence of dependence
57
58
59
60

1
2
3 of class labeling on centre in either validation or training dataset, ii. the final datasets
4
5 were balanced (i.e. as much HV as CRC). Dependence between class labeling and
6
7 centre of origin was excluded by testing with a Fisher's exact test ($p > 0.93$). The
8
9 random split was performed prior to fitting parameters and retained for all further
10
11 analyses to obtain realistic measures of classification performance. Since our
12
13 subroutines required complete data, we imputed missing values after assembling the
14
15 training and validation datasets for each dataset separately using nearest neighbor
16
17 averaging, as implemented in the impute-package in R⁹.

20
21 *Subroutines.* The SGMV compares the expression value of each input gene first to a
22
23 gene-specific cut-off and then assigns a defined class to an unknown sample
24
25 depending on whether the cut-offs are exceeded for at least half of the genes (i.e.
26
27 majority vote). The SGMV parameters hence consist of gene-specific cut-offs. The
28
29 gene-specific cut-offs are fitted by taking that value that corresponds to the point
30
31 closest to the top-left corner of the gene-associated ROC curve, using the pROC-
32
33 package in R¹⁰. The SVM with linear kernel is similar to linear discriminant analysis,
34
35 taking as input the expression values of a set of genes and comparing a linear
36
37 combination of the input values to a threshold in order to assign a defined class to an
38
39 unknown sample, thereby giving higher weight to more informative genes. The SVM
40
41 parameters hence consist of gene-specific weights and one threshold. We fitted the
42
43 parameters with the kernlab-package in R¹¹. The RF pushes the expression values of
44
45 a set of genes through a multitude of decision trees (each looking at a random subset
46
47 of genes and built from a random subset of samples from the training data), notes
48
49 down for each class the proportion of supporting individual trees and finally assigns
50
51 the class with highest support. The RF parameters hence consist of individual
52
53 decision trees. We fitted the parameters with the randomForest-package in R¹².

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

16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32 *Performance metrics.* The classifiers were validated on the qPCR test dataset,
33 constructed from healthy volunteers and patients of multi-centric cohort II who were
34 not included during development of the models (see above). To verify the similarity of
35 the test set to the training set, a Spearman-correlation between all assays was
36 performed, ensuring that test assays did not cluster separately from training assays.
37 A separate clustering would have been an indication that the training dataset was not
38 representative for the test samples. Two types of performance were finally reported:
39 ranking performance and classification performance. Ranking performance is the
40 capability of an algorithm to give a higher score to an individual from class CRC than
41 to an individual from class HV. We measured ranking performance by the area under
42 the ROC curve (AUC). For all 4 routines (SGMV, SVM, RF and ENS), we provided
43 the AUC as well as the lower bound and upper bound of its 95% confidence interval,
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 as computed after 2,000 bootstraps with the pROC-package in R¹⁰. Classification
4
5 performance measures the capability of an algorithm to assign an individual to the
6
7 correct class. We reported for all routines the balanced error rate (BER), sensitivity
8
9 (Se) and specificity (Sp). For Se and Sp, we also computed the lower bound and
10
11 upper bound of the 95% confidence interval after 2,000 bootstraps.
12
13

14 15 16 **Complementary data analysis**

17
18 A complementary data analysis by an independent team (DNAlytics, Belgium) on the
19
20 same 23-marker signature led to the same conclusions in terms of
21
22 performances. Another (per-marker) normalization procedure has been proposed.
23
24 This normalization is applied on the log-transformed gene expression (i.e. Δ CT
25
26 values) and consists in computing, on the training set (for example Cohort II, both HV
27
28 and CRC), the mean and standard deviation of each marker. When a prediction has
29
30 to be made on a new, potentially isolated sample, each marker measurement of this
31
32 new sample is normalized by subtracting the corresponding mean, and by dividing by
33
34 the corresponding standard deviation. A modified procedure has also been proposed
35
36 for the imputation of missing values, making it dependent on the reference cohort
37
38 only. This avoids the need for a new reference HV batch as prediction has to be
39
40 made on a new (set of) sample(s).
41
42
43
44

45
46 The first experiment consisted in cross-validating a model on Cohort II (BER: 8.4%
47
48 [3.4%;13.4%]; AUC: 0.93 [0.88;0.98]). A second experiment consisted in learning the
49
50 same type of model on Cohort II and having it make predictions on Cohort III (BER:
51
52 13.2%; AUC: 0.92). All analyses were performed in R with scripts designed by
53
54 DNAlytics, fully independent from other analyses described in this paper.
55
56
57
58
59
60

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.

1
2
3 Expression data were represented as mean \pm SEM of the indicated number of
4
5 measurements. Statistical significance of differential expression was assessed with
6
7 Prism (GraphPad, Inc.) by two-tailed unpaired t-test (for two conditions) and ANOVA
8
9 followed by Bonferroni correction (for more than two conditions) after ensuring equal
10
11 variance using F test.
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
51
52
53
54
55
56
57
58
59
60

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 context^{14 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

1
2
3 shift. However, we could not identify key regulators of known pathways, probably due
4 to the limited information on reciprocal effects of PBM and tumour cells¹⁶. Though of
5 high interest with regards to the biological function, functional biological knowledge is
6 dispensable to exploit the full potential of the gene signature as a diagnostic tool in
7 analogy to other important clinical tests, which are devoid of a biological
8 understanding (e.g., prostate specific antigen, PSA, and pro-calcitonin, PCT).
9
10
11
12
13
14
15
16
17

18 **Supplementary Note 3**

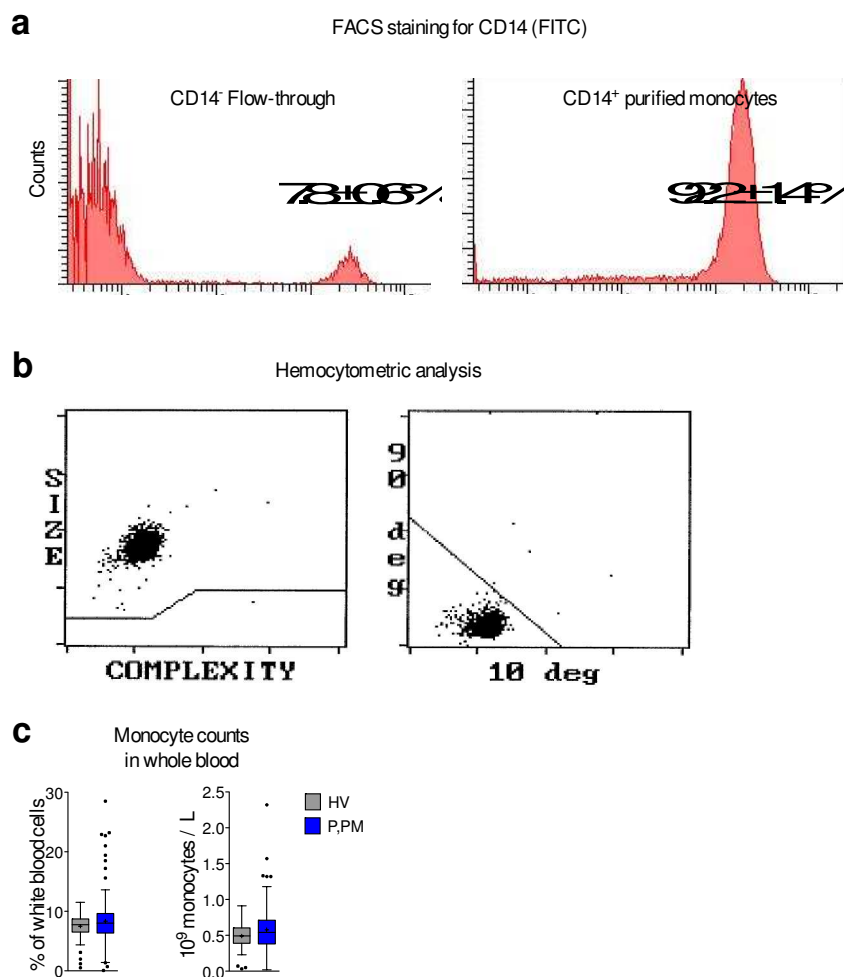
19
20 In accordance with our initial screening results, we found no differences in
21 expression patterns of P versus PM (data not shown). Moreover, as cumulating
22 evidence is suggesting subcategories of CRC according to its location¹⁷, we
23 investigated if the gene signature was capable of separating left versus right CRC or
24 colon versus rectal cancer, respectively. In line with the homogeneous clustering of
25 samples, we found no differences by location (AUC of 0.45 [0.20-0.73] for left versus
26 right CRC and AUC of 0.47 [0.28-0.70] for colon versus rectal cancer).
27
28
29
30
31
32
33
34
35
36
37

38 **Supplementary Note 4**

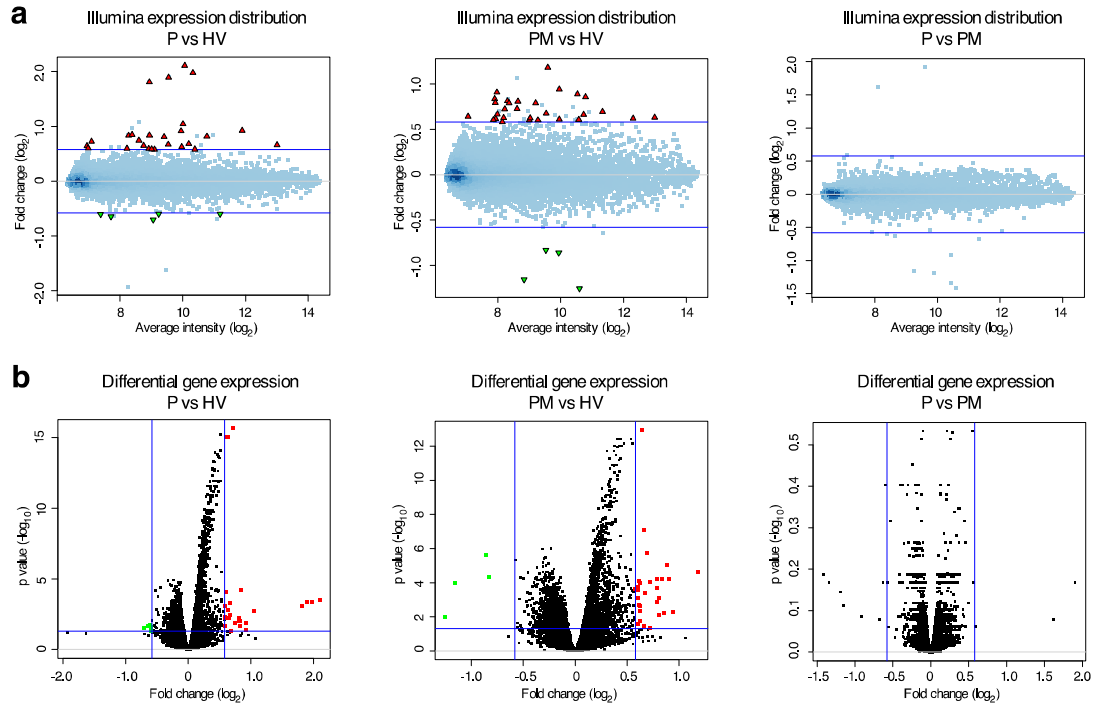
39
40 We sought to confirm our findings from the screening in independent samples by
41 independent techniques to rule out bias by the chosen technique and maximize
42 chances of extrapolation to other clinical centres. Our first step was a random re-
43 processing of collected samples and assessment by qPCR, which led to an initial
44 refinement of the gene signature, while some genes in this subset of samples
45 performed well even as single markers. By assessing Spearman correlation values
46 between expression data in the Illumina platform (used for screening) and the qPCR
47 technique (used for confirmation), we could rule out discrepancies in expression
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 between both analyses (Supplementary Figure 8). Consistently, a multicentric
4 validation trial revealed that the established gene signature retained the promising
5 performance observed in the screening phase, regardless of the centre and method
6 of analysis, while our multi-gene classification model allows to exploit the highest
7 informative content obtained from the expression analyses.
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
51
52
53
54
55
56
57
58
59
60

SUPPLEMENTARY FIGURES

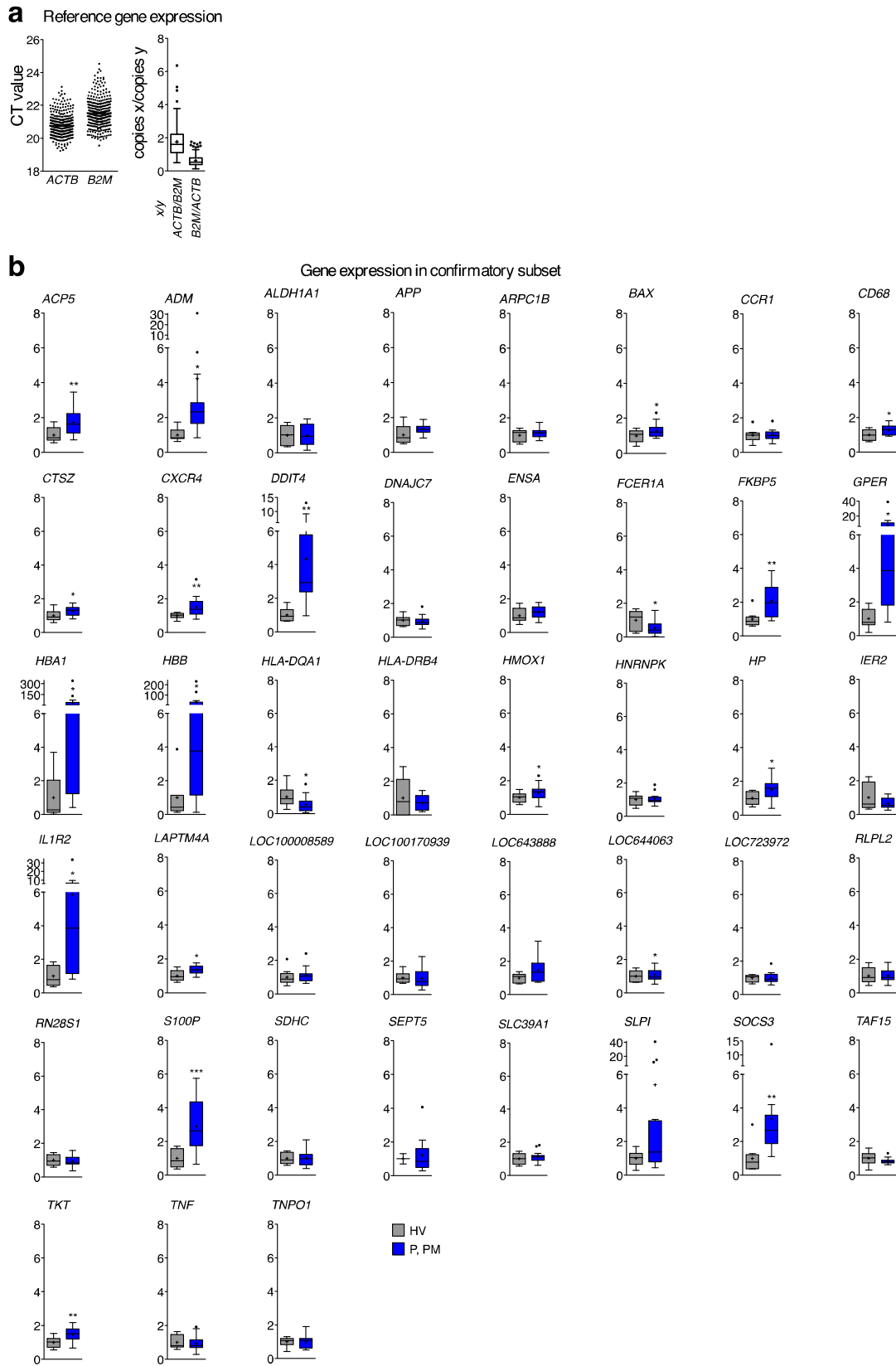
**Supplementary Figure 1: Isolation of PBM and monocyte counts**

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).



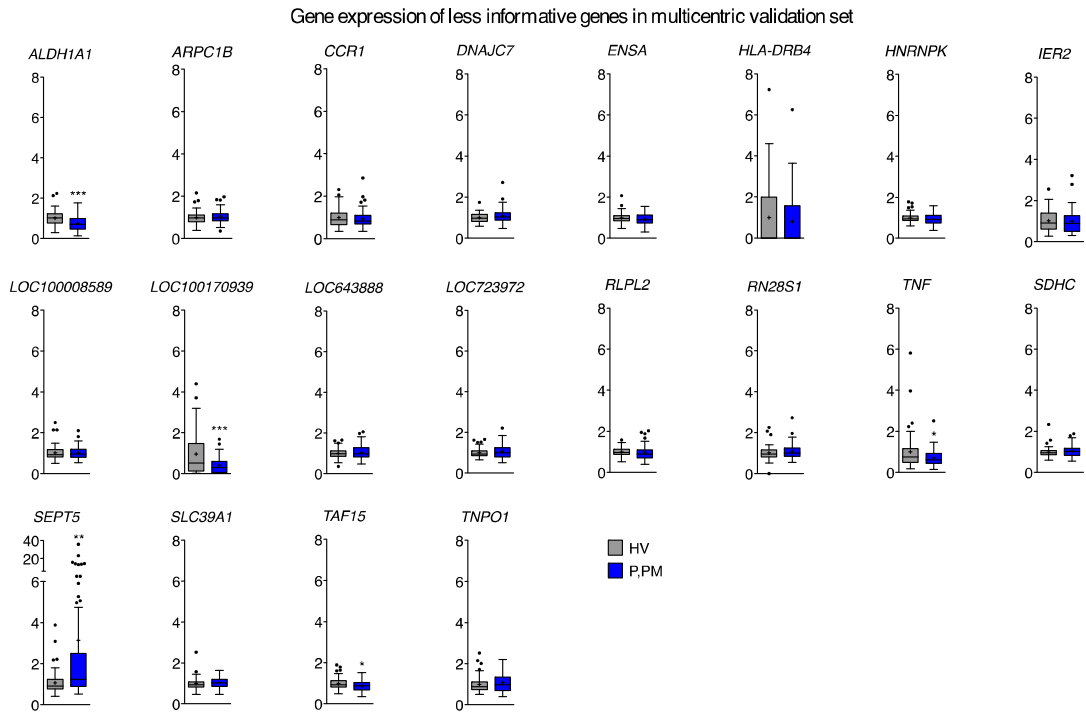
Supplementary Figure 2: Differentially expressed genes in PBM

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 $p < 0.05$.



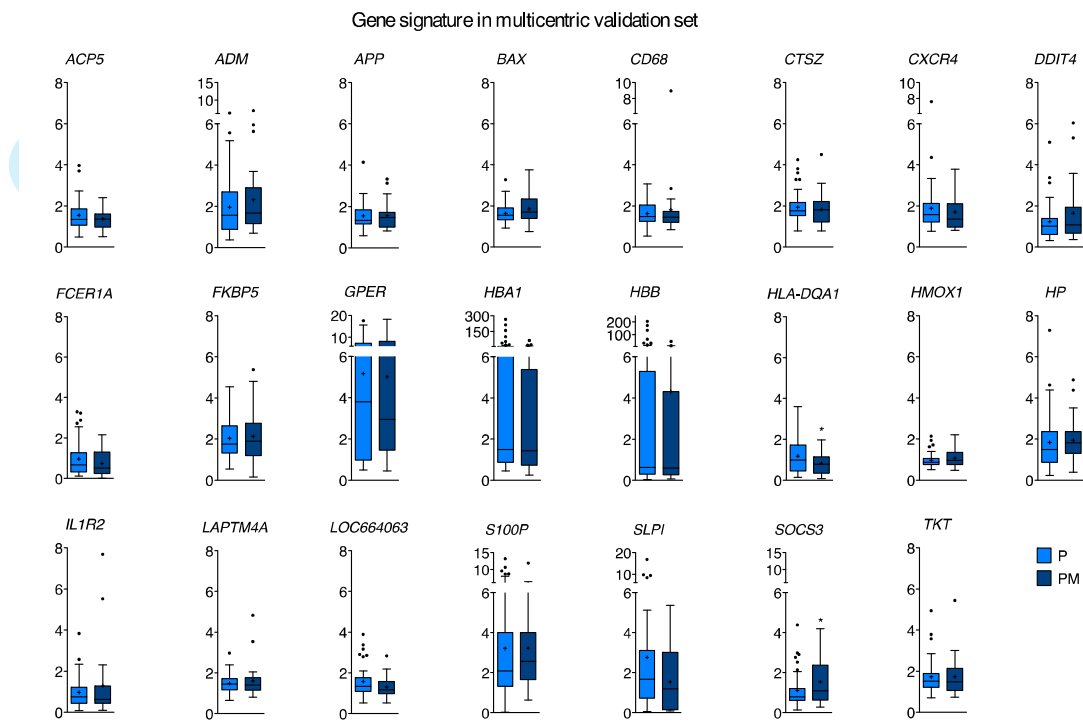
Supplementary Figure 3: Technical validation (subset of cohort I)

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 < 0.001.



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

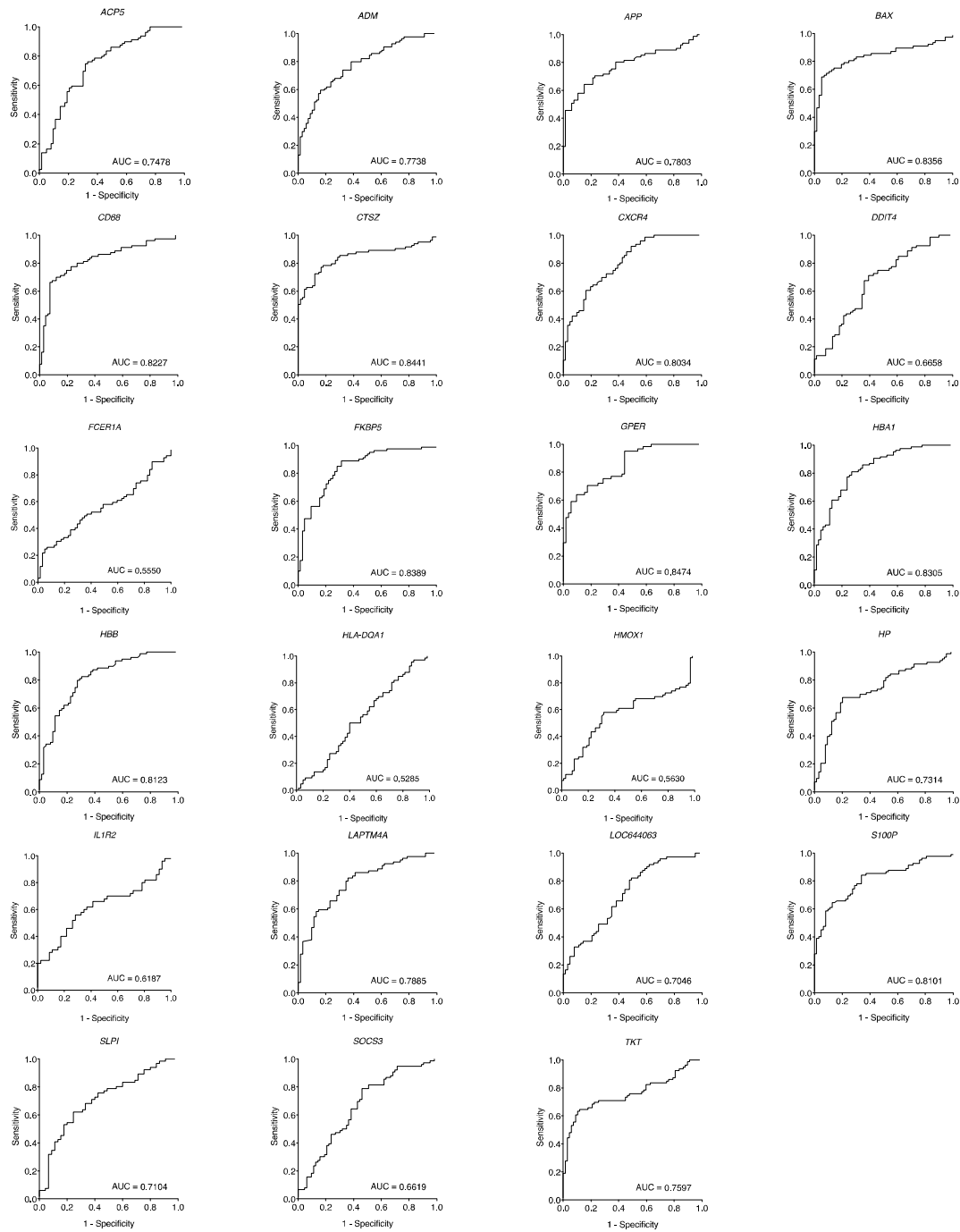
Expression levels are displayed as expression relative to the HV mean; boxes, first to third quartile; Whiskers, range; dots, values outside 1.5-times the interquartile distance; horizontal line, median; +, mean; *, $p < 0.05$; **, $p < 0.01$; ***, $p < 0.001$.



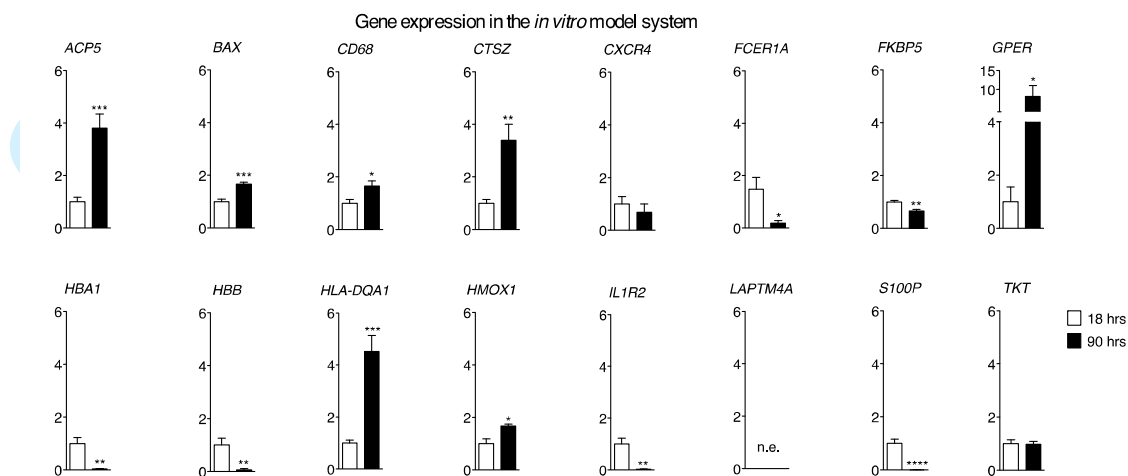
Supplementary Figure 5: The gene signature stays robust over disease progression (cohort II)

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; *, $p < 0.05$; **, $p < 0.01$; ***, $p < 0.001$.

ROC analysis for individual signature genes in full multicentric validation set

**Supplementary Figure 6: Single gene ROC analysis**

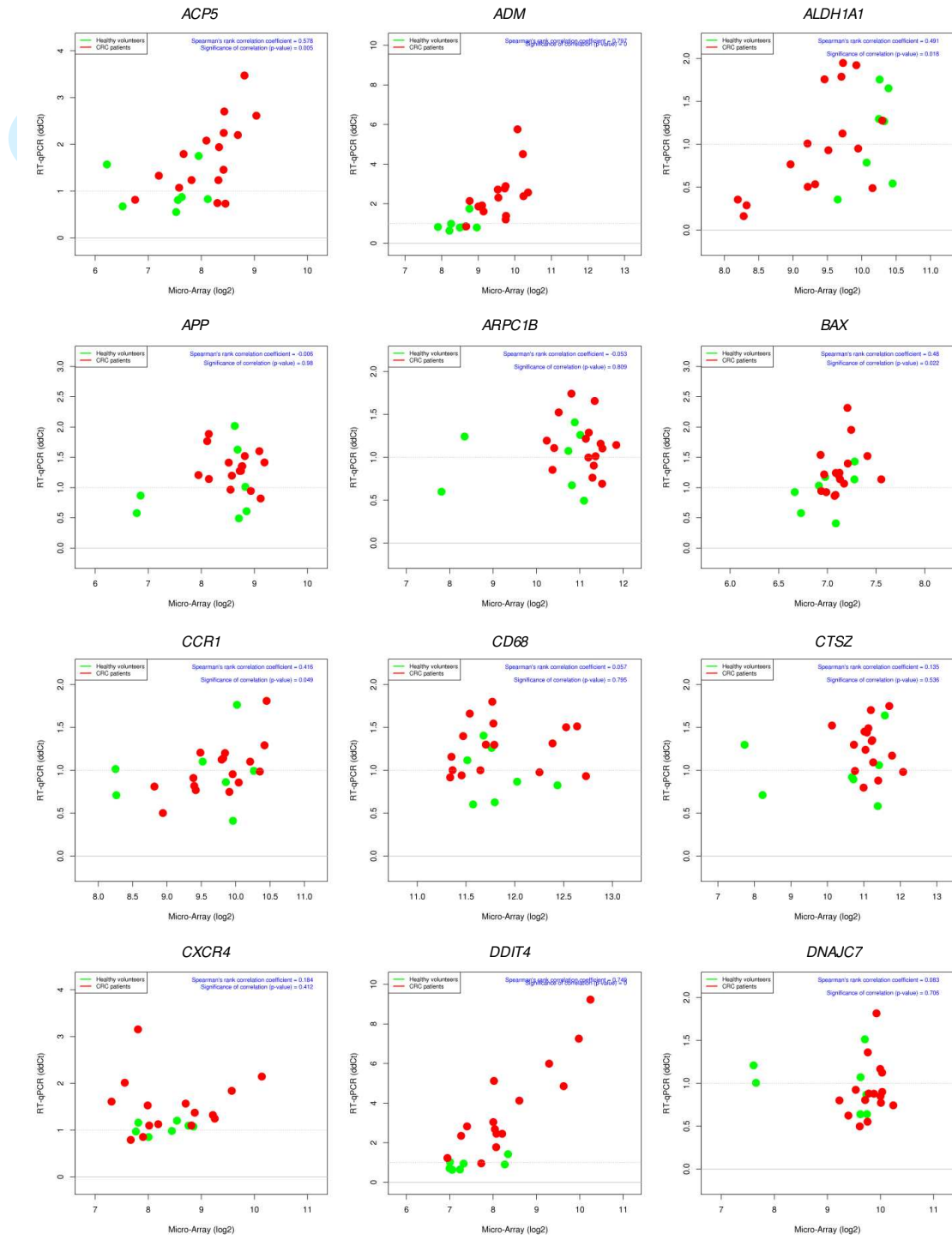
ROC analyses for each individual in cohort II. AUC, area under the curve.



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

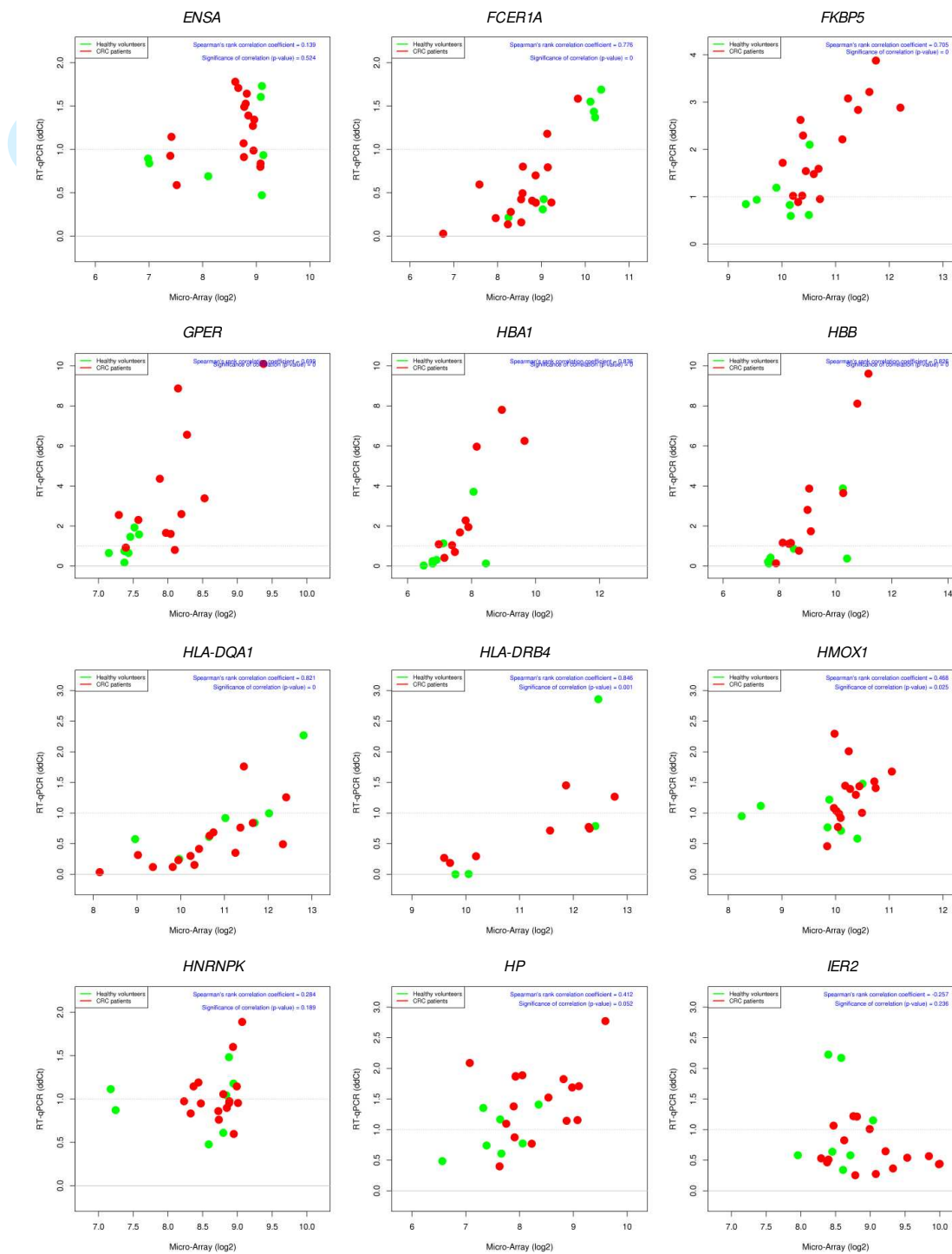
Shown are the expression levels of the 16 genes not selected out of the gene signature, which show altered expression levels in culture without any stimulus. Expression levels are shown as mean with SEM at 18 hours and 72 hours later (90 hours).

*, $p < 0.05$; **, $p < 0.01$; ***, $p < 0.001$; ****, $p < 0.0001$; n.e., not expressed *in vitro*.



Supplementary Figure 8:

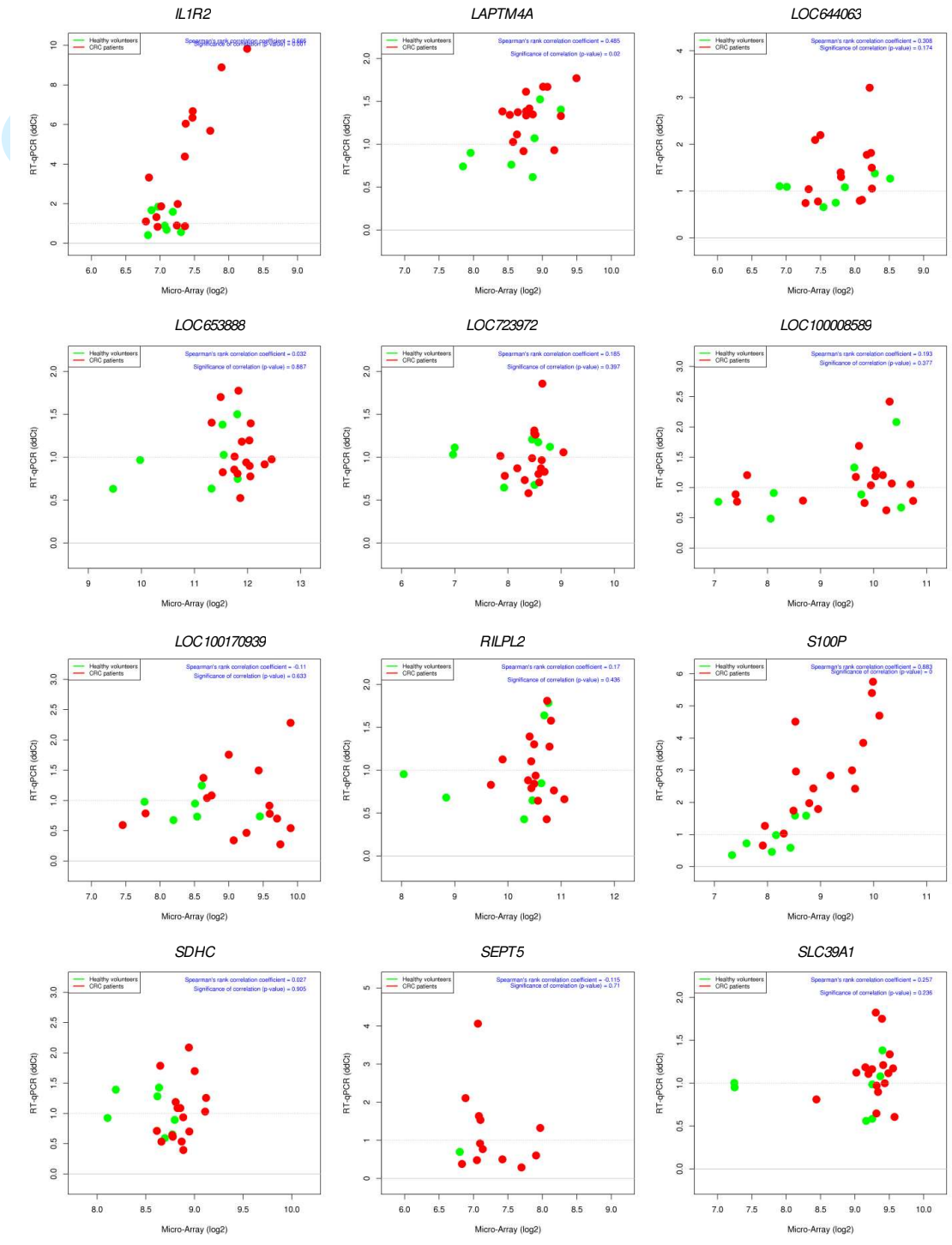
Scatter plots of Cohort I displaying correlation between Illumina microarray (x axis) and qPCR data (y axis). Spearman correlation values and p values are noted in the figures.



Supplementary Figure 8 – continued

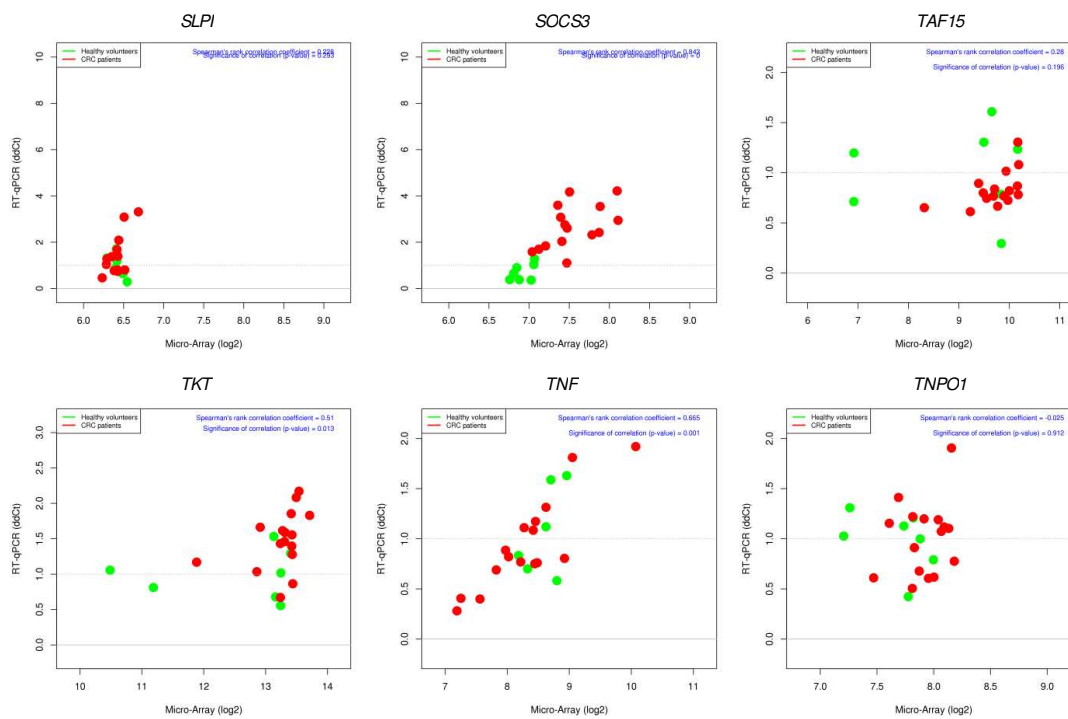


1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



Supplementary Figure 8 – continued





Supplementary Figure 8 – continued

SUPPLEMENTARY TABLES

Supplementary Table 1: IDT PrimeTime qPCR Assays

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

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

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
51
52
53
54
55
56
57
58
59
60

Supplementary Table 2: Dependence of class label on number of missing values (Fisher's exact test)

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

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

Genomewide Screening						Confirmation and Validation					
P,PM vs HV ^a			P vs HV			PM vs HV			P,PM vs. HV		
	Ratio	p		Ratio	p		Ratio	p		Ratio	p
<i>ADM</i>	2,00	<0,0001	<i>ADM</i>	1,75	0,0059	<i>ADM</i>	2,27	<0,0001	<i>ACP5^b</i>	1,61	<0,0001
<i>ALDH1A1</i>	0,66	0,0002	<i>CTSZ</i>	1,76	0,0103	<i>ALDH1A1</i>	0,56	<0,0001	<i>ADM</i>	2,16	<0,0001
<i>ARPC1B</i>	1,55	0,0209	<i>DDIT4</i>	1,78	0,0226	<i>AQP9</i>	1,62	<0,0001	<i>ALDH1A1</i>	0,88	<0,0001
<i>BAX</i>	1,50	0,0001	<i>DNAJC7</i>	1,59	0,0005	<i>BAX</i>	1,52	0,0008	<i>APP</i>	1,61	<0,0001
<i>CTSZ</i>	1,79	0,0007	<i>FCER1A</i>	0,61	0,0296	<i>CTSZ</i>	1,81	0,0056	<i>ARPC1B</i>	0,98	0,6497
<i>DDIT4</i>	1,71	0,0063	<i>HBA1</i>	3,51	0,0008	<i>DDIT4</i>	1,65	0,0477	<i>BAX</i>	1,76	<0,0001
<i>DNAJC7</i>	1,59	<0,0001	<i>HBA2</i>	4,31	0,0004	<i>DNAJC7</i>	1,60	0,0004	<i>CCR1</i>	0,90	0,3981
<i>FCER1A</i>	0,52	0,0002	<i>HBB</i>	3,95	0,0004	<i>DYSF</i>	1,52	0,0002	<i>CD68</i>	1,76	<0,0001
<i>FKBP5</i>	1,61	0,0001	<i>HMOX1</i>	1,55	0,0017	<i>FCER1A</i>	0,45	0,0001	<i>CTSZ</i>	1,96	<0,0001
<i>GPER</i>	1,58	0,0006	<i>HNRNPK</i>	1,60	0,0497	<i>FCGR1A</i>	1,52	0,0003	<i>CXCR4</i>	2,24	<0,0001
<i>HBA1</i>	2,33	0,0078	<i>HS.143909</i>	1,56	<0,0001	<i>FKBP5</i>	1,85	<0,0001	<i>DDIT4</i>	1,47	0,0025
<i>HBA2</i>	2,69	0,0051	<i>HS.581828</i>	1,52	<0,0001	<i>GPER</i>	1,78	0,0001	<i>DNAJC7</i>	1,07	0,1045
<i>HBB</i>	2,39	0,0099	<i>HS.61208</i>	1,65	<0,0001	<i>HLA-DRB6</i>	0,42	0,0102	<i>ENSA</i>	0,89	0,1122
<i>HMOX1</i>	1,54	0,0001	<i>IER3</i>	1,50	0,0009	<i>HMOX1</i>	1,53	0,0020	<i>FCER1A</i>	0,97	0,7541
<i>HNRNPK</i>	1,58	0,0125	<i>LOC100008589</i>	1,68	0,0131	<i>HP</i>	1,75	0,0080	<i>FKBP5</i>	2,45	<0,0001
<i>HP</i>	1,54	0,0131	<i>LOC100128274</i>	0,66	0,0195	<i>HS.61208</i>	1,56	<0,0001	<i>GPER</i>	5,29	<0,0001
<i>HS.143909</i>	1,51	<0,0001	<i>LOC100130707</i>	1,51	0,0232	<i>LOC100170939</i>	1,65	0,0001	<i>HBA1</i>	15,07	0,0165
<i>HS.61208</i>	1,60	<0,0001	<i>LOC100132394</i>	1,79	0,0095	<i>LOC100190986</i>	1,53	0,0001	<i>HBB</i>	11,96	0,0281
<i>IL1R2</i>	1,50	0,0482	<i>LOC100132727</i>	0,66	0,0282	<i>LOC153561</i>	1,73	0,0001	<i>HLA-DQA1</i>	1,01	0,8425
<i>LOC100008589</i>	1,55	0,0079	<i>LOC100134364</i>	1,57	0,0057	<i>LOC441087</i>	1,54	0,0177	<i>HLA-DRB4</i>	0,77	0,3931
<i>LOC100129685</i>	1,71	0,0356	<i>LOC153561</i>	1,50	0,0049	<i>RNF146</i>	1,50	0,0002	<i>HMOX1</i>	0,95	0,7338
<i>LOC100132394</i>	1,65	0,0045	<i>LOC649143</i>	1,90	0,0133	<i>S100P</i>	1,75	0,0007	<i>HNRNPK</i>	0,92	0,2280
<i>LOC100134364</i>	1,53	0,0009	<i>LOC723972</i>	1,51	0,0001	<i>SEPT5</i>	1,59	0,0347	<i>HP</i>	1,92	<0,0001
<i>LOC100170939</i>	1,54	<0,0001	<i>LOC728755</i>	0,64	0,0210	<i>SLC39A1</i>	1,54	0,0025	<i>IER2</i>	0,97	0,8782
<i>LOC153561</i>	1,61	<0,0001	<i>SLC39A1</i>	1,50	0,0058	<i>SOCS3</i>	1,73	0,0014	<i>IL1R2</i>	0,86	0,4209

Genomewide Screening						Confirmation and Validation					
P,PM vs HV ^a			P vs HV			PM vs HV			P,PM vs. HV		
	Ratio	p		Ratio	p		Ratio	p		Ratio	p
<i>LOC649143</i>	1,56	0,0356	<i>TAF15</i>	2,06	0,0001	<i>TAF15</i>	1,73	0,0002	<i>LAPTM4A</i>	1,59	<0,0001
<i>LOC653156</i>	1,73	0,0443	<i>TKT</i>	1,58	0,0034	<i>TKT</i>	1,55	0,0048	<i>LOC100008589</i>	0,99	0,9313
<i>LOC653737</i>	1,86	0,0472	<i>ZNF223</i>	0,66	0,0478	<i>TNPO1</i>	1,55	0,0001	<i>LOC100170939</i>	1,09	0,0004
<i>LOC728755</i>	0,66	0,0066				<i>UPP1</i>	1,58	<0,0001	<i>LOC643888</i>	1,05	0,2617
<i>S100P</i>	1,53	0,0020				<i>ZBTB16</i>	1,52	0,0252	<i>LOC644063</i>	1,54	<0,0001
<i>SEPT5</i>	1,57	0,0094						<i>LOC723972</i>	1,03	0,1904	
<i>SLC39A1</i>	1,52	0,0003						<i>RLPL2</i>	0,95	0,3618	
<i>SOCS3</i>	1,51	0,0043						<i>RN28S1</i>	1,03	0,3003	
<i>TAF15</i>	1,76	<0,0001						<i>S100P</i>	3,35	<0,0001	
<i>TKT</i>	1,56	0,0003						<i>SDHC</i>	1,04	0,3017	
								<i>SEPT5</i>	3,47	0,0020	
								<i>SLC39A1</i>	1,02	0,3827	
								<i>SLPI</i>	15,76	0,0090	
								<i>SOCS3</i>	1,60	0,0158	
								<i>TAF15</i>	0,84	0,0154	
								<i>TKT</i>	1,79	<0,0001	
								<i>TNF</i>	0,75	0,0205	
								<i>TNPO1</i>	1,02	0,3165	

^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).

Supplementary Table 4: Confirmation in random subset of cohort I

	Mean expression ^a	Fold ratio ^b	p	AUC
<i>ACP5</i>	81,336	1,73	0.0081^c	0.79
<i>ADM</i>	73,107	4,23	0.0941	0.95
<i>ALDH1A1</i>	47,649	0,99	0.9624	0.51
<i>APP</i>	176,332	1,32	0.1576	0.73
<i>ARPC1B</i>	1873,873	1,15	0.3336	0.57
<i>BAX</i>	11,474	1,29	0.0978	0.67
<i>CCR1</i>	338,797	1,01	0.9292	0.52
<i>CD68</i>	1821,912	1,27	0.0640	0.76
<i>CTSZ</i>	1580,418	1,28	0.0637	0.76
<i>CXCR4</i>	508,754	1,52	0.0065	0.84
<i>DDIT4</i>	36,963	4,34	0.0010	0.96
<i>DNAJC7</i>	105,494	0,92	0.5431	0.62
<i>ENSA</i>	250,287	1,21	0.2580	0.67
<i>FCER1A</i>	410,118	0,54	0.0768	0.73
<i>FKBP5</i>	39,131	2,08	0.0013	0.89
<i>GPER</i>	1,875	7,59	0.0138	0.93
<i>HBA1</i>	5243,339	41,40	0.0861	0.86
<i>HBB</i>	440,188	31,10	0.0773	0.85
<i>HLA-DQ1</i>	1748,345	0,53	0.0918	0.77
<i>HLA-DRB4</i>	1135,072	0,71	0.6316	0.53
<i>HMOX1</i>	405,989	1,30	0.0729	0.70
<i>HNRNPK</i>	1648,472	1,05	0.7356	0.52
<i>HP</i>	129,478	1,50	0.0218	0.76
<i>IER2</i>	0,657	0,65	0.2556	0.63
<i>IL1R2</i>	4,794	5,85	0.0288	0.87
<i>LAPTM4A</i>	338,391	1,35	0.0206	0.78
<i>LOC100008589</i>	18500877,250	1,12	0.5782	0.63
<i>LOC100170939</i>	252,768	0,96	0.8506	0.56
<i>LOC643888</i>	308,104	1,46	0.6143	0.55
<i>LOC644063</i>	1504,972	1,07	0.0415	0.71
<i>LOC723972</i>	568,957	1,00	0.9645	0.56
<i>RLPL2</i>	186,476	1,02	0.9078	0.53
<i>RN28S1</i>	14924567,167	0,92	0.5908	0.58
<i>S100P</i>	8,494	2,90	0.0003	0.91
<i>SDHC</i>	313,373	1,02	0.9145	0.53
<i>SEPT5</i>	0,298	1,22	0.6497	0.50
<i>SLC39A1</i>	166,812	1,12	0.4069	0.62
<i>SLPI</i>	1,656	5,39	0.2477	0.71
<i>SOCS3</i>	128,894	3,36	0.0081	0.91
<i>TAF15</i>	327,194	0,83	0.3084	0.65
<i>TKT</i>	986,111	1,48	0.0061	0.82
<i>TNF</i>	28,056	0,94	0.7325	0.53
<i>TNPO1</i>	140,447	1,00	0.9722	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

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

<u>Gene</u>	<u>Full Name</u>	<u>Biological Function</u>	<u>Potential Function in Monocytes</u>
<u>ACP5</u>	<u>acid phosphatase 5, tartrate resistant</u>	<u>iron containing glycoprotein involved in adhesion and migration</u>	<u>negative regulation of inflammatory response in interleukin pathways</u>
<u>ADM</u>	<u>adrenomedullin</u>	<u>vasodilation, regulation of hormone secretion, promotion of angiogenesis</u>	<u>antimicrobial activity, wound healing</u>
<u>APP</u>	<u>amyloid beta (A4) precursor protein</u>	<u>protein basis of amyloid plaques in Alzheimer disease</u>	<u>antimicrobial activity, mitotic activity</u>
<u>BAX</u>	<u>BCL2-associated X protein</u>	<u>p53-mediated activator of apoptosis</u>	<u>myeloid cell homeostasis</u>
<u>CD68</u>	<u>CD68 molecule</u>	<u>integral membran glycoprotein of scavenger receptor family</u>	<u>highly expressed on monocytes and macrophages, mediator of recruitment and activation</u>
<u>CTSZ</u>	<u>cathepsin Z</u>	<u>lysosomal cystein proteinase, involved in migration and adhesion</u>	<u>unknown</u>
<u>CXCR4</u>	<u>chemokine (C-X-C motif) receptor 4</u>	<u>CXC chemokine receptor specific for stromal cell-derived factor-1</u>	<u>mediator of recruitment, chemotaxis, and activation</u>
<u>DDIT4</u>	<u>DNA-damage-inducible transcript 4</u>	<u>negative regulation of mTOR signalling upon cellular stress</u>	<u>defense response to microbial signals</u>
<u>FCER1A</u>	<u>Fc fragment of IgE, high affinity I, receptor for; alpha polypeptide</u>	<u>alpha subunit of IgE-mediated allergic response</u>	<u>positive regulation of type-I immune response and macrophage differentiation</u>
<u>FKBP5</u>	<u>FK506 binding protein 5</u>	<u>member of immunophilin protein family, immunoregulation</u>	<u>receptor for FK506 and rapamycin, mediating calcineurin inhibition</u>
<u>GPER</u>	<u>G protein-coupled estrogen receptor 1</u>	<u>non-genomic signalling of estrogen stimulus</u>	<u>negative regulator of leukocyte activation; innate immune response</u>
<u>HBA1</u>	<u>hemoglobin, alpha 1</u>	<u>alpha chain of hemoglobin</u>	<u>unknown</u>
<u>HBB</u>	<u>hemoglobin, beta</u>	<u>beta chain of hemoglobin</u>	<u>positive regulation of nitric oxide synthesis,</u>
<u>HLA-DQA1</u>	<u>major histocompatibility complex, class II, DQ alpha 1</u>	<u>MHC class II receptor activity; peptide antigen binding</u>	<u>antigen processing and presentation</u>
<u>HMOX1</u>	<u>heme oxygenase (decycling) 1</u>	<u>heme catabolism</u>	<u>regulation of phagocytosis and migration, chemokine synthesis, wound healing, and angiogenesis</u>
<u>HP</u>	<u>haptoglobin</u>	<u>preproprotein of haptoglobin subunit</u>	<u>acute-phase defense response</u>
<u>IL1R2</u>	<u>interleukin 1 receptor, type II</u>	<u>cytokine receptor for IL-1</u>	<u>cytokine-mediated immune response</u>
<u>LAPTM4A</u>	<u>lysosomal protein transmembrane 4 alpha</u>	<u>unknown</u>	<u>unknown</u>
<u>LOC644063</u>	<u>heterogeneous nuclear ribonucleoprotein K pseudogene 4</u>	<u>unknown</u>	<u>unknown</u>
<u>S100P</u>	<u>S100 calcium binding protein P</u>	<u>cell cycle progression and differentiation</u>	<u>unknown</u>
<u>SLPI</u>	<u>secretory leukocyte peptidase inhibitor</u>	<u>secreted inhibitor of serin proteinases</u>	<u>negative regulation of endopeptidase activity</u>
<u>SOCS3</u>	<u>suppressor of cytokine signaling 3</u>	<u>negative regulator of cytokine signalling</u>	<u>modulator of immune response, particularly IFN-γ mediated</u>
<u>TKT</u>	<u>transketolase</u>	<u>enzyme of pentose phosphate pathway</u>	<u>metabolic modulator</u>

SUPPLEMENTARY REFERENCES

1. Weitz J, Koch M, Debus J, et al. Colorectal cancer. *Lancet* 2005;**365**(9454):153-65.
2. 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.
3. 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.
6. 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. <http://bioinformatics.mdanderson.org/MicroarraySampleSize/>.
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.
11. Burges CJC. A Tutorial on Support Vector Machines for Pattern Recognition. *Data Min Knowl Discov* 1998;**2**(2):121-67.

- 1
2
3 12. Liaw A, Wiener M. Classification and Regression by randomForest. R News: The
4
5 Newsletter of the R Project 2002;**2**(3):18-22.
6
- 7 13. Ambroise C, McLachlan GJ. Selection bias in gene extraction on the basis of
8
9 microarray gene-expression data. Proceedings of the National Academy of
10
11 Sciences of the United States of America 2002;**99**(10):6562-6.
12
- 13 14. Piehler A, Grimholt R, Ovstebo R, et al. Gene expression results in
14
15 lipopolysaccharide-stimulated monocytes depend significantly on the choice of
16
17 reference genes. BMC Immunology 2010;**11**(1):21.
18
- 19 15. Guo C, Liu S, Wang J, et al. ACTB in cancer. Clinica chimica acta; international
20
21 journal of clinical chemistry 2013;**417**:39-44.
22
- 23 16. Khatri P, Sirota M, Butte AJ. Ten Years of Pathway Analysis: Current Approaches
24
25 and Outstanding Challenges. PLoS Comput Biol 2012;**8**(2):e1002375.
26
- 27 17. Jess P, Hansen IO, Gamborg M, et al. A nationwide Danish cohort study
28
29 challenging the categorisation into right-sided and left-sided colon cancer.
30
31 BMJ open 2013;**3**(5).
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60