**Alternative splicing downstream of EMT enhances phenotypic plasticity and malignant behaviour in colon cancer.**

Tong Xu[1], Mathijs Verhagen[1], Rosalie Joosten[1], Wenjie Sun[2], Andrea Sacchetti[1], Leonel Munoz Sagredo[3,4], Véronique Orian-Rousseau[4], and Riccardo Fodde[1*]

[1]Department of Pathology, Erasmus University Medical Center, Rotterdam, The Netherlands. [2]Institut Curie, Laboratory of Genetics and Developmental Biology, Paris, France. [3]Faculty of Medicine, University of Valparaiso, Chile; [4]Institute of Biological and Chemical Systems - Functional Molecular Systems (IBCS-FMS), Karlsruhe Institute of Technology (KIT), Germany.

*to whom correspondence should be addressed at:

Erasmus MC, Department of Pathology, PO Box 2040, 3000 CA Rotterdam, The Netherlands

E-mail: r.fodde@erasmusmc.nl  Tel: +31 10 7043896

## Abstract

Phenotypic plasticity allows carcinoma cells to transiently acquire the quasi-mesenchymal features necessary to detach from the primary mass and proceed along the invasion-metastasis cascade. A broad spectrum of epigenetic mechanisms is likely to cause the epithelial–to-mesenchymal (EMT) and mesenchymal-to-epithelial (MET) transitions necessary to allow local dissemination and distant metastasis. Here, we report on the role played by alternative splicing (AS) in eliciting phenotypic plasticity in colon cancer.

By taking advantage of the identification of a subpopulation of quasi-mesenchymal and highly metastatic EpCAM$^{lo}$ colon cancer cells, we here show that the differential expression of *ESRP1* and other RNA-binding proteins (RBPs) downstream of the EMT master regulator *ZEB1*, alters the AS pattern of a broad spectrum of targets including *CD44* and *NUMB*, thus resulting in the generation of specific isoforms functionally associated with increased invasion and metastasis. Additional functional and validation studies indicate that both the newly identified RBPs and the CD44s and NUMB2/4 splicing isoforms promote local invasion and distant metastasis and are associated with poor survival.

The systematic elucidation of the spectrum of EMT-related RBPs and AS targets in colon cancer and in other malignancies, apart from the insights in the mechanisms underlying phenotypic plasticity, will lead to the identification of novel and tumor-specific therapeutic targets.

1 **Introduction**

2      Colon cancer still represents one of the major causes of cancer-related morbidity and mortality
3      worldwide. Apart from its high incidence, the adenoma-carcinoma sequence along which colon cancer
4      progresses has served as a classic model to elucidate the underlying genetic alterations representative of
5      virtually all of the hallmarks of cancers[1], possibly with the only exception of "activating invasion and
6      metastasis". As also reported in other epithelial cancers, the several steps of the invasion-metastasis
7      cascade are not caused by genetic alterations but rather by transient morphological and gene expression
8      changes of epigenetic nature[2,3]. In this context, epithelial–mesenchymal transition (EMT), and its reverse
9      MET, likely to represent the main mechanisms underlying local dissemination and distant metastasis[4,5].
10     EMT is triggered at the invasive front of the primary colon carcinoma in cells earmarked by nuclear β-
11     catenin and enhanced Wnt signaling, as the result of their physical and paracrine interactions with the
12     microenvironment[6]. The acquisition of quasi-mesenchymal features allows local invasion and
13     dissemination through the surrounding stromal compartment. Of note, EMT/MET should not be regarded
14     as binary processes in view of previous reports highlighting the existence of metastable hybrid E/M states
15     (partial EMT or pEMT) endowed with phenotypic plasticity and likely to underlie the reversible
16     morphological and functional transitions necessary to successfully complete the invasion-metastasis
17     cascade[7].
18     The molecular basis of the epigenetic changes underlying EMT and MET is likely to encompass a broad
19     spectrum of mechanisms ranging from chromatin remodeling and histone modifications to promoter DNA
20     methylation, non-coding RNAs (e.g. micro RNAs), and alternative splicing (AS). The inclusion/exclusion of
21     specific exons in mature mRNAs results in different protein isoforms with distinct biological functions. AS
22     occurs in 92–94% of human genes leading to enriched protein density[8,9]. Several sequence-specific RNA-
23     binding proteins (RBPs) have been identified which bind pre-mRNAs to control AS in context-dependent
24     fashion[10]. Multiple cancer-specific AS variants have been found to underlie progression and metastasis[11].
25     Likewise, alternative splicing has been suggested to play key roles in EMT/MET[12,13] and phenotypic
26     plasticity[14] in cancer by expression changes in RBP-encoding genes and their consequences for the
27     modulation of downstream AS targets.
28     The *ESRP1* (epithelial splicing regulatory protein 1) gene encodes for an epithelial-specific RBP and splicing
29     regulator shown to play a central role in EMT by modulating AS of EMT-associated genes including *FGFR2*,
30     Mena, *CD44* and p120-catenin[4]. Relevant to the present study, ESRP1 was reported to regulate the EMT
31     transition from CD44v (variable) to CD44s (standard) isoforms in breast and lung cancer progression[15,16].

3

1    As for colon cancer, whether ESRP1 regulates alternative splicing of CD44 and other target genes

2    downstream of EMT/MET activation during invasion and metastasis, is yet poorly understood.

3        Recently, we identified and thoroughly characterized subpopulations of CD44$^{hi}$/EpCAM$^{lo}$ colon cancer

4    cells (here referred to as EpCAM$^{lo}$) that coexist with their epithelial counterparts (CD44$^{hi}$/EpCAM$^{hi}$; for

5    brevity EpCAM$^{hi}$) through stochastic state transitions governed by phenotypic plasticity and pEMT[17].

6    Accordingly, EpCAM$^{lo}$ cells feature highly invasive and metastatic capacities. Here, we took advantage of

7    these *in vitro* models of phenotypic plasticity to study differential gene expression changes at upstream

8    RBPs and AS variations at target genes. Among the top AS targets, CD44 and NUMB were selected for

9    validation and functional studies in view of their known splicing isoforms and roles in stemness and cancer.

10   Moreover, we provide an extensive list of additional EMT- and colon cancer-related RBPs and AS targets

11   and show that the same CD44s and NUMB2/4 isoforms are conserved in ovarian and cervical cancer, i.e.

12   independently of the distinct modalities through which these malignant cells metastasize.

13

1  **Results**

2  ***Differential expression of RNA-binding proteins in quasi-mesenchymal and highly metastatic colon***

3  ***cancer cells (EpCAM<sup>lo</sup>) affects alternative splicing of a broad spectrum of downstream target genes.***

4  As previously reported, the EpCAM[lo] subpopulation of colon cancer cells is earmarked by increased

5  expression of the *ZEB1* transcription factor, responsible for EMT activation and for their quasi-

6  mesenchymal and highly metastatic phenotype[17]. It has been established that in breast and pancreatic

7  cancer *ZEB1*-driven EMT downregulates the expression of the RNA-binding protein and splicing regulator

8  *ESRP1* as part of a self-enforcing feedback loop[18]. Accordingly, among the top differentially expressed

9  genes between EpCAM[lo] and EpCAM[hi] in SW480 and HCT116 colon cancer cells, *ESRP1* was found to be

10  downregulated both at the RNA and protein level in the quasi-mesenchymal subpopulation where *ZEB1*

11  expression is upregulated (Figure 1A-C). Gain- and loss-of-function analyses of both genes confirmed the

12  inter-dependence of their expression levels in both cell lines (Figure 1D-E). Of note, *ESRP1*-overexpression

13  in the HCT116 and SW480 cell lines resulted in the dramatic reduction of their EpCAM[lo] subpopulations

14  and the expansion of the epithelial bulk, as shown by FACS analysis (Figure 1F). However, *ESRP1*

15  knockdown (KD) only marginally affected CD44 though not EpCAM expression levels (Suppl. Figure 1A).

16  The latter suggests that additional RNA binding proteins are likely to be involved in the alternative splicing

17  regulation of the EpCAM[lo] colon cancer subpopulation. Indeed, by taking advantage of the RBPDB

18  database[19], we found that, apart from *ESRP1*, consistent differential expression in the quasi-mesenchymal

19  subpopulation of both cell lines was observed for *ESRP2*, *RBM47*, *MBNL3* (down-regulated) and *NOVA2*,

20  *MBNL2*, (up-regulated). Other RBPs were found to be differentially expressed though in only one of the

21  two cell lines (Suppl. Figure 1B). In validation of the clinical relevance of the RBPs found to be differential

22  expressed between the EpCAM[hi/lo] subpopulations derived from the SW480 and HCT116 cell lines, the

23  RBP-coding genes *QKI*, *RBM24*, and *MBNL2* (up in EpCAM[lo]), and *ESRP1/2* and *RBM47* (down in EpCAM[lo])

24  were found to be respectively up- and down-regulated in the consensus molecular subtype 4 (CMS4) of

25  colon cancers, responsible for ~25% of the cases and earmarked by poor prognosis and a pronounced

26  mesenchymal component (Suppl. Figure 1C)[20].

27  Differentially spliced target genes between EpCAM[lo] and EpCAM[hi] colon cancer cells from the SW480

28  and HCT116 cell lines were selected based on exon skip splicing events with ΔPSI (differential Percentage

29  Spliced In) values > 10%. The PSI value ranges from 0 to 1 and is a measurement of the percentage of

30  isoform with an alternative exon included[21]. This resulted in a large and rather heterogeneous group of

31  alternative spliced targets (n=1495; Suppl. Table 1a) with no clear enrichment in any specific gene

1   ontology class (data not shown). In order to identify differentially spliced target genes in RBP-specific

2   fashion, we took advantage of RNAseq data sets from previous *ESRP1-*, *ESRP2-*, *RBM47-*, and *QKI-*

3   knockdown studies in different cancer cell lines and compared them with our own AS data relative to the

4   EpCAM[hi/lo] colon cancer subpopulations[17] (Suppl. Figure 2A). A total of 32 common skipped exons events

5   in 20 genes were identified between EpCAM[lo] colon (both cell lines) and *ESRP1* KD H358 lung cancer cells[22]

6   (Figure 2A). More extensive lists of common *ESRP1* AS events and target genes were obtained when the

7   SW480 and HCT116 cell lines were individually compared with the lung cancer study (Suppl. Table 1b-c).

8   As for the alternative splicing targets of RBPs other than *ESRP1*, based on the available RNAseq data from

9   knockdown studies of *ESRP2* (in the LNCaP cell line[23]), *RBM47* (H358[22]), and *QKI* (CAL27; GEO Accession:

10  GSM4677985), several common and unique genes were found (Suppl. Figure 2 and Suppl. Table 2).

11  Notably, four EMT-related genes (*CTNND1*[24], *LSR*[25], *SLK*[26], and *TCF7L2*[27]) were common to all RBP KD

12  studies analyzed (Suppl. Figure 2).

13

14  ***The CD44s and NUMB2/4 ESRP1-specific AS isoforms are preferentially expressed in EpCAM[lo] colon***

15  ***cancer cells.***

16      From the newly generated lists of RBP-specific alternative splicing targets, we selected *CD44* and

17  *NUMB* for further analysis, based both on their *ESRP1*-specific AS patterns and on their well-established

18  roles in EMT, stemness/differentiation, and cancer progression.

19  CD44, a transmembrane cell surface glycoprotein, has been show to play key roles in inflammatory

20  responses and in cancer metastasis[28]. The *CD44* gene encompasses 20 exons of which 1-5 and 16-20 are

21  constant and exist in all isoforms. In contrast, exons 6-14, also referred to as variants exons v2-v10, are

22  alternatively spliced and often deregulated in cancer[28]. The *NUMB* gene and its protein product have been

23  involved in a broad spectrum of cellular phenotypes including cell fate decisions, maintenance of stem

24  cell niches, asymmetric cell division, cell polarity, adhesion, and migration. In cancer, NUMB is a tumor

25  suppressor that regulates, among others,  Notch and Hedgehog signaling[29]. The mammalian *NUMB* gene

26  encodes for 4 isoforms, ranging from 65 to 72 KD, differentially encompassing two key functional domains,

27  i.e. the amino-terminal phosphotyrosine-binding (PTB) domain, and a C-terminal proline-rich region (PRR)

28  domain[29].

29  Based on the above ΔPSI-based AS analysis, decreased expression of CD44v (variable) isoforms was

30  observed in EpCAM[lo] and *ESRP1*-KD cells, accompanied by increased CD44s (standard) isoform expression

31  (Figure 2B). Likewise, the NUMB2/4 isoforms appear to be preferentially expressed in EpCAM[lo] and *ESRP1*-

32  KD, accompanied by decreased NUMB1/3 expression (Figure 2B). RTqPCR and western analyses validated

1    these *in silico* data: CD44s and NUMB2/4 isoforms were preferentially expressed in EpCAM$^{lo}$ colon cancer

2    cells, in contrast with the increased CD44v and NUMB1/3 levels in EpCAM$^{hi}$ cells (Figure 2C-D). In view of

3    its previously suggested role in invasion and metastasis[30], we focused on the CD44v6 isoform.

4    As reported above, AS events at the *NUMB* and *CD44* genes correlate with decreased ESRP1 expression.

5    To confirm this observation, we up- and down-regulated *ESRP1* in the SW480 and HCT116 cell lines. The

6    dox-inducible shRNA vector used for the KD studies reduces ESRP1 expression by 5-10 fold (Figure 1D-E)

7    and resulted in the upregulation of the CD44s and NUMB2/4 isoforms at the mRNA and protein level in

8    both cell lines (Figure 3A-B and Suppl. Figure 3). Likewise, *ESRP1* overexpression led to an increase in the

9    CD44v6 and NUMB1/3 isoforms, found in association with the bulk of epithelial colon cancer cells (Figure

10   3C-D and Suppl. Figure 3).

11

12   ***Transcriptional and functional consequences of the CD44s and NUMB2/4 isoforms on colon cancer***

13   ***invasion and metastasis.***

14   In order to elucidate the functional contribution exerted by the newly identified CD44s and NUMB2/4

15   isoforms on the overall invasive and metastatic capacities of colon cancer cells, we first ectopically

16   expressed each of them (individually and in combination for NUMB1/3 and 2/4) in the HCT116 and SW480

17   cell lines (Suppl. Figure 3E-H), and analyzed their consequences *in vitro* by cell proliferation, transwell

18   migration assay, RTqPCR, western, FACS, and RNAseq, and *in vivo* by spleen transplantation. A significant

19   increase in migratory capacity, comparable to that of EpCAM$^{lo}$ cells sorted from the parental lines, was

20   observed in SW480 and HCT116 upon overexpression of the CD44s and NUMB2/4 isoforms. Likewise,

21   ectopic expression of the single NUMB2 or 4 isoforms resulted in increased migration rates when

22   compared with NUMB1 and 3. In contrast, overexpression of CD44v6 and NUMB1/3, normally prevalent

23   in the epithelial bulk (EpCAM$^{hi}$) of both cell lines, did not affect their migratory properties (Suppl. Figure

24   4A-B).

25   In agreement with the migration assays, overexpression of CD44s and NUMB2/4 results in the significant

26   upregulation of the EMT transcription factors (EMT-TFs) *ZEB1*, accompanied by the up- and

27   downregulation regulation of mesenchymal and epithelial markers such as *VIM* (vimentin), *CDH1* (E-

28   cadherin), and *EpCAM*, respectively (Suppl. Figure 4C). Of note, expression of *ESRP1*, the main upstream

29   splicing regulator of both CD44 and NUMB, was also decreased in CD44s- and NUMB2/4-OE cells, in

30   confirmation of the self-enforcing feedback loop that characterize its interaction with ZEB1 and EMT

31   activation[18]. In agreement with the well-established regulation of Notch signaling by NUMB isoforms[29],

7

1    established Notch target genes and were accordingly up- (*HES1*, *HEY1*) and down-regulated (*ID2*) upon

2    overexpression of NUMB2/4 (Suppl. Figure 4D).

3    FACS analysis was then employed to evaluate the overall effect of the ectopic expression of the specific

4    CD44 and NUMB isoforms on the relative percentages of the EpCAM$^{hi/lo}$ subpopulations in the HCT116 and

5    SW480 cell lines. As shown in Figure 4A, CD44s overexpression led to a dramatic increase of the EpCAM$^{lo}$

6    subpopulation at the expenses of EpCAM$^{hi}$ cells. The opposite effect was observed with CD44v6, i.e. the

7    enlargement of the EpCAM$^{hi}$ gate and the corresponding decrease of EpCAM$^{lo}$ cells. As for NUMB, ectopic

8    expression of NUMB2/4 significantly increased the relative proportion of EpCAM$^{lo}$ cells while reducing the

9    size of the EpCAM$^{hi}$ subpopulation, while the opposite was observed with NUMB1/3 (Figure 4B-C). Of note,

10    the single NUMB2 and NUMB4 isoforms appear dominant in their capacity to enlarge the HCT116 and

11    SW480 EpCAM$^{lo}$ subpopulations, respectively. The same was true for NUMB1 and NUMB3 in the

12    consequences of their ectopic expression in reducing the size of the HCT116 and SW480 EpCAM$^{lo}$ fractions,

13    respectively (Figure 4B-C). In agreement with the RTqPCR analysis of EMT markers, CD44s overexpression

14    negatively affected overall proliferation rates in both cell lines, whereas the opposite was observed upon

15    CD44v6 expression (Suppl. Figure 5A-B). Likewise, NUMB1/3 expression positively affected proliferation

16    rates in HCT116 and SW480, whereas the NUMB2/4 isoforms exert the opposite effects. In both cases,

17    synergistic effects were observed upon co-expression of NUMB1/3 and 2/4, when compared to the

18    individual isoforms (Suppl. Figure 5C-D).

19    In order to assess the *in vivo* the consequences of the ectopic expression of the CD44 and NUMB isoforms

20    on the capacity of colon cancer cells to form metastatic lesions in the liver, parental HCT116 and SW480

21    cells and their CD44s-, CD44v6-, NUMB1/3-, and NUMB1/4-overexpressing counterparts were injected in

22    the spleen of immune-incompetent recipient mice. In agreement with the *in vitro* results, overexpression

23    of both NUMB2/4 and CD44s isoforms significantly increased the multiplicity of liver metastases, whereas

24    CD44v6 and NUMB1/3 did not differ from the parental controls (Figure 4D-E).

25    Next, in order to elucidate the signaling pathways and molecular and cellular mechanisms triggered by

26    the CD44 isoforms, we analyzed by RNAseq HCT116 and SW480 cells ectopically expressing CD44s and

27    CD44v6. After dimension reduction with principal component analysis (PCA), the samples separated by

28    group (i.e. CD44s-OE, CD44v6-OE, and controls) (Figure 5A). Notably, the CD44s-OE samples showed most

29    distinct expression in both cell lines when compared to the parental and CD44v6-OE cell lines. In HCT116,

30    the CD44v6 samples shared most similarity with the CD44s samples, while in SW480, the CD44v6 samples

31    were most similar to the parental cell line. Thus, we observed both an isoform independent effect,

32    presumably as the result of the ectopic CD44 expression (and most dominantly visible in HCT116), and an

1    isoform dependent effect as depicted by the separation of CD44s and CD44v6 samples (Figure 5A). As

2    expected, differential expression analysis of the CD44s and v6 isoforms overexpressing samples compared

3    with the parental cell lines revealed an overall upregulation of gene expression (Suppl. Figure 6A). Next,

4    in order to identify which genes are specifically upregulated by the different CD44 isoforms, we performed

5    differential expression analysis between the CD44s samples and the CD44v6 samples. To this aim, we

6    employed k-means clustering on the scaled expression values to separate genes specific for the CD44s

7    isoform (e.g. *SPARC, ZEB1, VIM*), the CD44v6 isoform (e.g. *IL32, TACSTD2, CSF2*), and genes that were

8    indiscriminative for the CD44v6 isoform or the parental cell lines (e.g. *MAL2, ESRP1, CDH1*) (Figure 5B).

9    Finally, to identify the most distinct differences in signaling pathways and GO functional categories, we

10   performed a gene set enrichment analysis (GSEA) by comparing the CD44s- with the CD44v6-

11   overexpressing samples in the individual cell lines. Among the significantly altered pathways (normalized

12   enrichment score > 1, pval < 0.01), epithelial mesenchymal transition (EMT) was the only one upregulated

13   in CD44s vs. CD44v6 in both cell lines (Figure 5C). Additional pathways and GO categories activated by

14   CD44s appeared to be cell line specific, e.g. Wnt beta catenin signaling (HCT116) and oxidative

15   phosphorylation (SW480). Of note, the detailed GSEA analysis evidenced how several inflammatory

16   (TNF/NF$\kappa$B; IL6/JAK/STAT3; IF$\alpha$/$\gamma$; ILK2/STAT5) and signaling (KRAS, MYC, E2F) pathways were common

17   to both CD44s and v6, presumably as the result of the ectopic CD44 expression, regardless of the isoform

18   (Suppl. Figure 6B).

19

20   ***Increased ZEB1 and decreased ESRP1 expression correlate with the NUMB2/4 and CD44s isoforms and***

21   ***with poor overall survival***

22   In order to assess the clinical relevance of the results obtained with the SW480 and HCT116 cell lines,

23   we analyzed RNAseq data from patient-derived colon cancers available from the public domain and the

24   scientific literature. To this aim, the TCGA Splicing Variants Database (TSVdb; www.tsvdb.com) was

25   employed to integrate clinical follow-up data with RBP and AS expression profiles obtained from The

26   Cancer Genome Atlas project (TCGA) and from the Guinney et al study[20] on the classification of human

27   colon cancers into four consensus molecular subtypes (CMS1–4). The main limitation of this approach is

28   the low representation of quasi-mesenchymal (EpCAM$^{lo}$-like) subpopulations in bulk RNAseq preparations

29   and the masking effect that the majority of epithelial (EpCAM$^{hi}$-like) cancer cells are likely to cause. To

30   identify tumors enriched in EpCAM$^{lo}$-like cells, we first stratified them based on *ZEB1* expression (*ZEB1*>8.6:

31   ZEB1$^{hi}$; ZEB1<8.3: ZEB1$^{lo}$; 8.2<ZEB1<8.6: Intermediate). Subsequently, we used *ESPR1* expression levels to

32   further define the tumors into *ZEB1*$^{hi}$*ESRP1*$^{lo}$ (*ESRP1*<11.8; hereafter referred to as *ZEB1*$^{hi}$), *ZEB1*$^{lo}$*ESRP1*$^{hi}$

9

1   (*ESRP1*>11.6; hereafter referred to as *ZEB1*[lo]). Tumors with intermediate *ZEB1* expression levels and

2   tumors with *ESRP1* expression levels outside these thresholds were defined as intermediate (Figure 6A).

3   Kaplan-Meier analysis showed that *ZEB1*[hi] tumors have an overall decreased survival probability (p = 0.045)

4   (Figure 6B). Next, we compared the expression of CD44 and NUMB isoforms across the *ZEB1*[hi/lo] tumors.

5   Notably, while no significant differences were observed based on the expression level of the whole CD44

6   and NUMB genes, significant differences were found for their specific isoforms (Figure 6C). Analysis of the

7   specific isoforms expression across the different consensus molecular subtypes[20] revealed elevated CD44s

8   and NUMB2/4 expression in the CMS4 subtype, known to be enriched in mesenchymal lineages in tumor

9   and TME cells, and strongly associated with poor survival and the greatest propensity to form distant

10  metastases (Figure 6D). Likewise, the majority of the *ZEB1*[hi] group was composed of the CMS4 subtype

11  (72%), while the *ZEB1*[lo] group was mainly contributed by CMS2 (49%) and CMS3 tumors (31%), with few

12  CMS4 tumors (1%) (Figure 6E).

13  Next, we correlated the expression of CD44s/v6 isoforms in patient-derived colon tumors with the

14  differentially expressed genes (DEGs) identified in the isoform-overexpressing cell lines (Figure 7A). While

15  overall *CD44* expression correlated with both isoforms, the DEGs from the CD44s-OE samples showed

16  specific correlation with CD44s expression in patient-derived tumors (e.g. *SPARC, ZEB1*), the DEGs from

17  the CD44v6 samples correlated with CD44v6 but not with CD44s (e.g. *KDF1, ESRP1*).

18  Last, we correlated the CD44 and NUMB isoforms expression in patient-derived colon cancers with

19  functional signatures obtained by averaging the scaled expression levels for each of the hallmark sets[31].

20  The CD44s and NUMB2/4 isoforms showed overall similar correlating hallmarks and pathways. However,

21  the same was not true when compared to the CD44v6- and NUMB1/3-associated functional signatures.

22  Here, most invasion/metastasis-relevant hallmarks (e.g. EMT, angiogenesis, apical junctions) showed a

23  positive correlation with CD44s and NUMB2/4, though not with CD44v6 and NUMB1/3 (Figure 7B).

24      In sum, we confirmed a switch in isoform expression (CD44v6 vs. CD44s and NUMB1/3 vs. NUMB2/4)

25  as a function of *ESRP1* and *ZEB1* expression in colon cancer. Expression of the EpCAM[lo]–specific isoforms

26  (CD44s and NUMB2/4) is elevated in CMS4 tumors overall survival.

27

28  ***Upregulation of the NUMB2/4 and CD44s isoforms is common to quasi-mesenchymal cells from cancers***

29  ***other than colon.***

30      In order to assess whether the preferential expression of the NUMB2/4 and CD44s isoforms is specific

31  to the modalities of local invasion and distant metastasis characteristic of colon cancer, we interrogated

32  expression profiling data previously obtained by comparing epithelial and quasi-mesenchymal

10

1    subpopulations from ovarian (OV90) and cervical (SKOV6) cancer cell lines (*manuscript in preparation*).

2    Ovarian cancer, because of the distinct anatomical localization of the primary lesion, metastasizes the

3    abdominal cavity with very different modalities than colon cancer, namely by peritoneal dissemination

4    rather than local dissemination into the stroma microenvironment followed by intra- and extravasation

5    of the portal blood stream[32,33]. On the other hand, metastasis in carcinoma of the cervix occurs both by

6    lymphatic or hematogenous spread to the lung, liver, and bones. We asked whether, notwithstanding the

7    distinctive patterns of metastatic spread, the CD44s and NUMB2/4 isoforms were preferentially expressed

8    in the corresponding EpCAM[lo] RNAseq profiles. To this aim, EpCAM[hi/lo] subpopulations from OV90 and

9    SKOV6 were sorted and analyzed by RNAseq and RTqPCR, similar to our previous study on colon cancer[17].

10   As shown in Suppl. Figure 7, both NUMB2/4 and CD44s isoforms appear to be upregulated in the OV90

11   and SKOV6 cell lines, as also validated by RTqPCR.

**Discussion**

The capacity to invade the tumor microenvironment and to form distant metastases undoubtedly represents the most clinically relevant hallmark of epithelial cancer cells. However, the complexity and diversity of the obstacles that carcinoma cells encounter along the invasion-metastasis cascade require transient and reversible changes that cannot be explained by the *de novo* acquisition of genetic alterations. Instead, epigenetic modifications underlie phenotypic plasticity, i.e. the capacity of cancer cells with a given genotype to acquire more than one phenotype in a context-dependent fashion[34]. Epithelial-to-mesenchymal and mesenchymal-to-epithelial transitions (EMT/MET) are central to the phenotypic plasticity characteristic of metastasizing carcinoma cells and are prompted by a broad spectrum of epigenetic mechanisms ranging from chromatin remodeling by histone modifications, DNA promoter methylation, non-coding RNAs, and alternative splicing (AS)[35]. Here, we have taken advantage of our previous identification of phenotypic plastic and highly metastatic EpCAM$^{lo}$ colon cancer cells[17] to characterize the genome-wide AS events that accompany EMT/MET state transitions between the epithelial bulk (EpCAM$^{hi}$) and the quasi-mesenchymal subpopulation.

In view of the central role played by RNA-binding proteins in AS, we first identified RBP-coding genes differentially expressed between the EpCAM$^{lo}$ and EpCAM$^{hi}$ fractions of two commonly employed colon cancer cell lines , representative of the chromosomal- and microsatellite-instable subtypes (SW480, CIN; HCT116, MIN)[36]. The Epithelial Splicing Regulatory Protein 1 and 2 genes (*ESRP1/2*)[37], the "*splicing masterminds*" of EMT[38,39], were found among the top downregulated RBP coding genes in EpCAM$^{lo}$ colon cancer cells, as part of a self-enforcing feedback loop with the EMT-TF *ZEB1*[18]. Accordingly, *ZEB1* upregulation in EpCAM$^{lo}$ colon cancer cells is invariably accompanied by *ESRP1/2* downregulation, and *ZEB1*$^{hi}$/*ESRP1*$^{lo}$ colon cancers, predominantly belonging to the mesenchymal CMS4 subgroup, have a significantly worse survival outcome when compared with *ZEB1*$^{lo}$/*ESRP1*$^{hi}$ patients.

Apart from *ESRP1*, several other RBP-coding genes were found to be differentially expressed between epithelial and quasi-mesenchymal colon cancer cells. Whereas the majority of RBP-coding DEGs, like *ESRP1*, appear to be downregulated upon EMT induction (*ESRP1/2*, *RBM14/19/47*, *MBNL3*, *HNRPAB/PF*, *USAF2*), others were activated in the quasi-mesenchymal EpCAM$^{lo}$ fraction (*NOVA2*, *MBNL2*, *QKI*, *SRSF5*, *HNRNPH*, *RBM24/43*). Accordingly, in patient-derived colon cancers stratified according to their consensus molecular signature, the same *QKI*, *RBM24*, and *MBNL2* genes were found to have increased expression in CMS4 tumors, known for their pronounced mesenchymal composition and poor prognosis[20]. Of note, the mesenchymal nature of CMS4 tumors has previously been questioned as these lesions often feature pronounced infiltration from the surrounding microenvironment, the extent of which might cover

12

1    their true cellular identity other than   representing a mere contamination from the tumor

2    microenvironment[40,41]. As shown in our previous study[17], the EpCAM[lo] cells do represent *bona fide* quasi-

3    mesenchymal colon cancer cells, enriched among CMS4 cases, and likely responsible for their poor

4    prognosis. The observed upregulation of RBPs such as quaking (*QKI*) is caused by the presence in its 3'UTR

5    of target sequences of the miR-200 family of microRNAs[42,43]. The latter is analogous to the regulation of

6    the expression of the EMT-TF *ZEB1* gene, whose activation during EMT is regulated by the same microRNA

7    family[44]. Accordingly, the significantly reduced levels of all five miR-200 members in EpCAM[lo] cells[17]

8    underlies the coordinated upregulation of both *ZEB1* and *QKI*.

9    The increased expression of other RBP-coding genes such as *RBM24* and *MBNL2* (muscleblind-like 2) in

10    CMS4 tumors and in EpCAM[lo] cells is also in sharp contradiction with their alleged tumor suppressing roles

11    in colon and other cancers[45,46]. Of note, MBNL2 regulates cancer migration and invasion through

12    PI3K/AKT-mediated EMT[46] and its overexpression in breast and cancer cell lines inhibits their metastatic

13    potential[47]. In contrast to *MBNL2*, *MBNL3*, a distinct member of the muscleblind family, is downregulated

14    in EpCAM[lo] colon cancer cells, similar to what reported in prostate cancer by Lu and colleagues[48]. *NOVA2*,

15    a member of the Nova family of neuron-specific RNA-binding proteins, was also upregulated in the quasi-

16    mesenchymal cells from both cell lines, possibly as the result of the differential expression of miR-7-5p[49],

17    as previously shown in non-small cell lung[49] and prostate[48] cancer. The identification the AS targets

18    downstream of specific RBPs in quasi-mesenchymal cancer cells from different malignancies will likely

19    clarify these apparent contradictions and shed light the functional roles of distinct members of the splicing

20    machinery in EMT and metastasis.

21    The spectrum of AS target genes downstream of the RBPs differentially expressed in EpCAM[lo] colon

22    cancer cells appears extremely broad when it comes to specific cellular processes or signaling pathways.

23    Nonetheless, comparison of our RNAseq data with KD studies of specific RBPs from the public domain

24    (*ESRP1/2*[23], *RBM47*[22], and *QKI* [GEO Accession: GSM4677985]) allowed us to identify common and unique

25    AS target genes associated with specific downstream effectors. By following this admittedly imperfect

26    approach, the top 4 AS targets common to all of the above-mentioned RBPs notwithstanding their up- or

27    downregulation in EpCAM[lo] colon cancer cells, i.e. *CTNND1* ($\delta$- or p120-catenin), *LSR* (Lipolysis Stimulated

28    Lipoprotein Receptor), *SLK* (STE20 Like Kinase), and *TCF7L2* (Transcription Factor 7-Like 2, or TCF4) are

29    known regulators and effectors of epithelial-to-mesenchymal transition[24-27], thus pointing to the central

30    role played by alternative splicing in the regulation of EMT in the malignant evolution of colon cancer.

31    Here, we have focused on CD44 and NUMB as two ESRP1-specific AS target genes with well-established

32    functional roles in EMT and in cancer invasion and metastasis. The CD44s and NUMB2/4 isoforms appear

13

1    to be specifically expressed in quasi-mesenchymal colon cancer cells both from the immortalized cell lines

2    and from patient-derived tumors, with a striking enrichment in the CMS4 subgroup of colon cancer

3    patients. In contrast, the CD44v6 and NUMB1/3 isoforms are preferentially expressed in the epithelial

4    bulk of the tumor. The latter, as far as CD44v6 is concerned, sharply contrasts what previously reported

5    by Todaro et al.[30] where this specific isoform was found to earmark the colon cancer stem cells (CSCs)

6    which underlie metastasis. CD44v6 and other 'variable' CD44 isoforms (CD44v4-10) earmark *Lgr5*[+]

7    intestinal stem cells (ISCs), i.e. the stem cells of origin of intestinal tumors, and accordingly promote

8    adenoma formation *in vivo*[50-52]. A plausible explanation for the discordant results lies in the epithelial

9    nature of the models employed in the above study and in the requirement of both EMT and MET for the

10   completion of the invasion-metastasis cascade[5]. By employing tumor spheres and freshly sorted CD133[+]

11   tumor cells, Todaro et al. focused on epithelial CSCs where, as observed in normal ISCs, the CD44v6

12   isoform is predominantly expressed, and is necessary for EMT to occur upon interaction with c-MET[30]. The

13   CD44v6 isoform is strictly required for c-MET activation by hepatocyte growth factor (HGF, or scatter

14   factor)[53] and as such plays an essential role in triggering EMT at the invasive front where tumor cells are

15   exposed to these TME-secreted factors. Our own immunoprecipitation studies confirmed that CD44v6 but

16   not CD44s binds to cMET in response to HGF stimulation (*data not shown*). Therefore, HGF/SF stimulation

17   of colon cancer cells along the invasive front will trigger the acquisition of quasi-mesenchymal

18   characteristics and the AS-driven switch from CD44v6 to CD44s, the latter unable to bind HGF and as such

19   controlling the extension of EMT activation. The reverse switch will take place upon the activation of the

20   mesenchymal-to-epithelial transitions necessary for the colonization of  the distal metastatic site. From

21   this perspective, both CD44 isoforms are essential for the completion of the invasion-metastasis cascade.

22   The functional relevance of the CD44s isoforms has been highlighted in malignancies other than colon

23   cancer, namely in prostate[48] and breast cancer where it activates, among others, PDGFRβ/Stat3 and Akt

24   signaling to promote EMT and CSC traits[15,54]. GO analysis of the RNAseq profiles from colon cancer cells

25   ectopically expressing CD44s highlighted a broader spectrum of signaling pathways likely to underlie EMT.

26   Accordingly, analysis of RNAseq data from primary colon cancers stratified for their CD44s expression

27   revealed an equally broad spectrum of downstream EMT-related biological processes. Of note, among the

28   DEGs identified upon CD44s ectopic expression which correlate with *ZEB1*[hi]/*ESRP1*[lo] (and CMS4) colon

29   cancers, the *SPARC* gene, a partial EMT marker in the EpCAM[hi/lo] state transitions[17], was found.

30   Expression of NUMB2/4 isoforms both in cells lines and in patient-derived colon tumors is associated with

31   signaling pathways and GO categories largely overlapping with those linked to CD44s (and CD44v6 with

32   NUMB1/3), possibly suggesting synergism between AS at these genes. Accordingly, NUMB is involved in a

1    broad spectrum of cellular phenotypes in homeostasis and in cancer where it mainly function as a tumor

2    suppressor[29]. NUMB inhibits EMT by suppressing the Notch signaling pathway. As such, downregulation

3    of NUMB can induce an EMT phenotype in isoform-specific fashion. Analysis of colon cancer cells

4    individually overexpressing each of the four isoforms revealed an increased basal Notch signaling in

5    NUMB2 and 4, as shown by the expression of the 'universal' targets *HES1* and *HEY1*. Instead, ectopic

6    expression of NUMB1/3 resulted in increased transcriptional levels of the more atypical Notch signaling

7    target *ID2*. Although the functional consequences of the NUMB2/4 (and 1/3) isoforms on Notch regulation

8    of EMT is yet unclear, it seems plausible that the complex network of AS targets activated downstream

9    the RBP-coding DEGs, including CD44, NUMB and many others as shown here, will eventually lead to the

10   'just-right' level of plasticity needed to allow both the 'mesenchymalization' during local invasion and

11   systemic dissemination, and the reacquisition of epithelial features at the distant site of metastasis.

12   Overall, it appears that alternative splicing substantially contribute to the epigenetic mechanisms that

13   underlie EMT/MET in colon cancer metastasis. The systematic elucidation of the RBPs and AS targets will

14   not only elucidate the cellular and molecular mechanisms underlying phenotypic plasticity as the most

15   clinically relevant hallmark of cancer, but it will also offer novel tumor-specific targets for therapeutic

16   intervention based on small molecule inhibitors and even RNA vaccination.

17

**18   Authors' contribution**

19   T.X. performed most of the experiments and wrote the first manuscript draft. M.V. contributed the *in*

20   *silico* analysis and validation of the AS results in patient-derived RNAseq data. R.J. and A.S. contributed to

21   the implementation of PCR, mouse, and FACS experiments. W.S. analysed AS splicing in the RNAseq data.

22   L.M.S. and V.O-R. contributed to the CD44 AS analysis and critically revised the manuscript. RF conceived

23   the experimental strategy and wrote the manuscript.

24

28

1 **Materials and Methods**

2 *Cell Cultures*

3 The human colon cancer cell lines HCT116 and SW480, obtained from the European Collection of

4 Authenticated Cell Culture (ECACC), were cultured in DMEM (11965092, Thermo Fisher Scientific) with 10%

5 FBS (Thermo Fisher Scientific), 1% penicillin/streptomycin (Thermo Fisher Scientific, 15140122), and 1%

6 glutamine (Gibco, 25030024), in humidified atmosphere at 37°C with 5% $CO_2$.

7 *Plasmid transfection and lentiviral transduction*

8 Stable transfection of the *ESRP1* (Sino Biological plasmid # HG13708-UT), *CD44s*, *CD44v6*, and NUMB1-4

9 (from V.O.R.) expression plasmids was performed using FuGENE HD transfection reagent (Promega, E2311)

10 according to the manufacturer's protocol and selected with Geneticin (Gibco, 10131035). As for the

11 knockdown constructs, the *ESRP1*-shRNA plasmid (Horizon, V3THS_335722) was packaged by pPAX2

12 (Addgene # 12260) and pMD2.G (Addgene # 12259) into HEK293T. The virus-containing supernatant was

13 collected 24 hrs. after transfection, filtered, and employed to infect the HCT116 and SW480 cell line.

14 Selection was applied with 750 ng/ml puromycin (Invivogen, San Diego, USA) or 800 μg/ml of Geneticin

15 selection for 1-2 weeks. The efficiency of overexpression and knockdown was assessed by qPCR and

16 western blot 48-72 h after transfection.

17 *qRT-PCR and PCR analyses*

18 Total RNA was isolated using TRIzol reagent (Thermo Fisher Scientific, 15596018) and was reverse-

19 transcribed using high-capacity cDNA reverse transcription kit (Life Technologies, 4368814), according to

20 the manufacturer's instructions. qRT-PCR was performed using the Fast SYBR Green Master Mix (Thermo

21 Fisher Scientific) on an Applied Biosystems StepOne Plus Real-Time Thermal Cycling Research with three

22 replicates per group. Relative gene expression was determined by normalizing the expression of each

23 target gene to GAPDH. Results were analyzed using the 2-(ΔΔCt) method. To validate isoform switches by

24 RT-PCR, hCD44 and mCD44 primers were used (Table S3). qRT-PCR and PCR primers are listed in

25 Supplementary Table S3.

26 *Western analysis*

27 Cells were lysed in 2X Laemmli buffer containing 4% SDS, 48% Tris 0.5M pH6.8, 20% glycerol, 18% $H_2O$,

28 bromophenol blue and 10% 1M DTT, and subjected to sodium dodecyl sulfate (SDS)- polyacrylamide gel

29 electrophoresis (PAGE), followed by transfer onto polyvinylidene fluoride (PVDF) membranes (Bio-Rad).

30 After blocking with 5% milk in TBS-Tween, the membranes were incubated with primary antibodies

1  against ESRP1 (1:1000, Invitrogen), CD44 (1:100, Invitrogen), CD44v6 (1:1000, Abcam), NUMB (1:1000,

2  Cell Signaling) and β-actin (1:2000, Cell Signaling), followed by polyclonal goat anti-mouse/ rabbit

3  immunoglobulins horseradish peroxidase (HRP)-conjugated secondary antibody (Dako) at appropriate

4  dilutions. The signals were detected with Pierce ECT western blotting subtrade (Thermo) using Amersham

5  AI600 (GE Healthcare, USA).

6  *Flow cytometry analysis and sorting*

7  Single-cell suspensions generated in PBS supplemented with 1% FBS were incubated with anti-EpCAM-

8  FITC (1:20, Genetex), and anti-CD44-APC (1:20, BD Pharmingen) antibodies for 30 min on ice and analyzed

9  on a FACSAria III Cell Sorter (BD Biosciences). CD44$^{hi}$EpCAM$^{hi}$and CD44$^{hi}$EpCAM$^{lo}$ HCT116 and SW480 cells

10 were sorted and cultured in humidified atmosphere at 37°C with 5% CO$_2$ for 3-5 days before collecting

11 RNA or protein, as previously described[17].

12 *MTT assay*

13 For MTT assay, 2×10$^3$ HCT116, SW480 parental, CD44v6, CD44s, and NUMB1-4 OE cells were plated in 96

14 well plates and incubated at 37°C, 5% CO$_2$. 24 hours later, in the culture medium was supplemented with

15 100μl 0.45 mg/mL MTT (3-(4,5-dimethylthiazol-2-yl)-2,5-diphenyltetrazolium bromide; Sigma-Aldrich)

16 and again incubated for 3 hrs.. The 96-well plates were then centrifuged at 1,000 rpm for 5 min and the

17 culture medium removed. MTT formazan precipitates were solubilized with DMSO. OD reading was

18 performed at 595 nm with microplate reader (Model 550, Bio-Rad). Background measurements were

19 subtracted from each data point. Experiments were performed in duplicate for each individual cell line

20 and drug. Cell numbers were calculated every 24 hrs. for a 6 days period for proliferation analysis.

21 *Cell migration assay*

22 Migration assay were conducted with 8-μm pore PET transwell inserts (BD Falcon™) and TC-treated multi-

23 well cell culture plate (BD Falcon™). 5×10$^4$ cells were seeded in 100 μl of serum-free growth medium in

24 the top chamber. Growth medium containing 10% FBS was used as a chemoattractant in the lower

25 chamber. After 24 hrs., cells migrated to the lower chamber were fixed with 4% PFA, stained with 0.1%

26 trypan blue solution, and counted under the microscope.

27 *Mouse spleen transplantation*

28 All mice experiment were implemented according to the Code of Practice for Animal Experiment in Cancer

29 Research from the Netherlands Inspectorate for Health Protections, Commodities and Veterinary Public

30 Health. Mice were fed in the Erasmus MC animal facility (EDC). NOD.Cg-Prkdc$^{scid}$ Il2rg$^{tm1Wjl}$/SzJ (NSG) mice

17

1    from 8 to12 week-old were used for spleen transplantation. Anesthetics Ketamine (Ketalin®, 0.12 mg/ml)

2    and xylazine (Rompun®, 0.61mg/ml) were given intraperitoneally, while the analgesic Carpofen (Rimadyl®,

3    5 mg/ml) was injected subcutaneously. $5 \times 10^4$ HCT116 and SW480 cells resuspended in 50 µl PBS were

4    injected into the exposed spleen with an insulin syringe and left for 15 minutes before splenectomy.

5    Transplanted mice were sacrificed after 4 and 8 weeks and analyzed for the presence of liver metastases.

6    *Alternative splicing analysis*

7    The following public available RNASeq (SRA database) data relative to RBP (RNA binding protein) knock-

8    down (KD) studies were used: ESRP1-KD and RMB47-KD in the human non-small cell lung cancer cell line

9    H358[22] with accession ID SRP066789 and SRP066793; ESRP2-KD in the human prostate adenocarcinoma

10   cancer cell line LNCaP[23] with accession ID SRP191570; the QKI-KD in the oral squamous cell carcinoma cell

11   line CAL27 datasets with accession number SRX8772405. Together with our own EpCAM[hi/lo] RNASeq data

12   obtained from the colon cancer cell lines[17], the sequencing reads were mapped to GRCh37.p13.genome

13   by STAR[55] (https://www.gencodegenes.org/human/release_19.html). MISO[56] was used to quantify AS

14   events with annotation from https://miso.readthedocs.io/en/fastmiso/index.html#iso-centric. The MISO

15   uses the alternative exon reads and adjacent conservative reads to measure the percentage of transcript

16   isoform with specific exon included, termed Percentage Spliced In (PSI or Ψ). The PSI ranges from 0 (i.e.

17   no isoform includes a specific alternative exon) to 1 (i.e. all of the isoforms detected comprise the

18   alternative exon).

19   We removed alternative events with low expression of related transcript isoforms if less than 3 samples

20   in a dataset had more than 10 informative reads to calculate the PSI. Next, we compared the PSI between

21   RBPs KD and wild type in each cell line, as well as the PSI between EpCAM[hi] and EpCAM[lo] groups in the

22   SW480 and HCT116 colon cancer cell lines. AS events were defined as differentially spliced events when

23   the difference of mean PSI between two groups (Δpsi; differential Percentage Spliced In) was >10%.

24

25   *RNAseq analysis*

26   RNA quality was first evaluated by NanoDrop and further purified by DNAse treatment followed by the

27   TURBO DNA-free Kit protocol (Invitrogen). Samples were sequenced with the DNA nanoball (DNB) seq

28   protocol (BGI) to a depth of 50 million reads per sample. Adapter sequences and low-quality sequences

29   were filtered from the data using SOAPnuke software (BGI). Reads were aligned to the human reference

30   genome build hg19 with the RNAseq aligner STAR (v2.7.9a) and the Homo sapiens GENCODE v35

31   annotation. Duplicates were marked with Sambamba (0.8.0) and raw counts were summed using

32   FeatureCounts (subread 2.0.3). Downstream analysis was performed in R using the DESeq2 package

1   (v1.30.1). After variance stabilizing transformation, principal component analysis was performed on each

2   cell line separately. Differentially expressed genes were identified by comparing the different groups of

3   ectopically expressing CD44 samples with a Wald test, and by selecting the genes with absolute log fold

4   change above 1.5 and padj < 0.1. Gene set enrichment analysis was performed with the Fsgsea package

5   using the HallMark geneset from the molecular signature database, and by selecting significant pathways

6   based on normalized enrichment score (NES) > 1 and pvalue < 0.05.

7   *RNAseq data from primary (patient-derived) colon cancers*

8   Patient data from The Cancer Genome Atlas (TCGA), with annotation of the consensus molecular subtypes

9   (CMS) as described in Guinney et al.[20] were integrated with splicing data from the TCGA splicing variant

10  database (TSVdb, www.tsvdb.com). For splicing analysis, RNA-seq by expectation maximization (RSEM)

11  values were log transformed and expression levels of each isoform (CD44std: isoform_uc001mvx, CD44v6:

12  exon_chr11.35226059.35226187, NUMB1: isoform_uc001xny, NUMB2: isoform_uc001xoa, NUMB3:

13  isoform_uc001xnz, NUMB4: isoform_uc001xob) were annotated to the patients. Isoform expression was

14  compared in groups based on the CMS groups and tumor expression levels (*ZEB1*, *ESRP1*). Tumors were

15  stratified on *ZEB1* expression levels using a log rank test top optimize overall survival differences

16  (thresholds: 8.3, 8.6). Next, ESRP1 expression was used to purify the groups into $ZEB1^{hi}ESRP1^{lo}$ and

17  $ZEB1^{lo}ESRP1^{hi}$ (thresholds: 11.6, 11.8). Survival analysis was done using the Kaplan-Meier method with the

18  survival and survminer packages in R. Correlation analysis was done by computing the Pearson Correlation

19  between the isoforms and whole gene expression levels as processed in Guinney et al.[20]. Likewise,

20  association between isoform expression and pathway activity was evaluated by computing the Pearson

21  Correlation between the isoforms and the average scaled expression values of the pathways, as defined

22  in the HallMark gene set from the molecular signature database[31].

23  *Data accessibility*

24  The RNA-sequencing data from this study have been submitted to the Gene Expression Omnibus (GEO)

25  database under the accession number GSE192877. Other data referenced in this study are publicly

26  available and can be accessed from the GEO using GSE154927[17] , GSE154730 and Synapse using identifier

27  syn2623706[20] .

28

29

### *Figure Legends*

### Figure 1

**a.** Gene rank plot showing differentially expressed genes between EpCAM$^{hi}$ and EpCAM$^{lo}$ with combined analysis of HCT116 and SW480.

**b.** RT-qPCR *ESRP1* expression analysis of HCT116 and SW480 EpCAM$^{hi}$, EpCAM$^{lo}$, and bulk subpopulations. *GAPDH* expression was used as control normalized with the bulk subpopulation in each sample (Means±SEM, n=3). ** = p<0.01.

**c.** ESRP1 western analysis in HCT116 and SW480 EpCAM$^{hi}$, EpCAM$^{lo}$, and bulk fractions. β-actin was used as loading control.

**d.** RT-qPCR *ZEB1* and *ESRP1* expression analysis in *ZEB1*-OE and -KD HCT116 and SW480 cells. Expression values were normalized in each sample with those from the parental HCT116 and SW480 cell lines. HCT116 and SW480 cells transduced with the sh*ZEB1* lentivirus were induced by 1 µg/mL doxycycline for 72 hrs.. Expression values were normalized with those from non-induced cells; *GAPDH* expression was employed as control (Means±SEM, n=3).

**e.** RT-qPCR *ZEB1* and *ESRP1* expression analysis in *ESRP1*-OE and -KD HCT116 and SW480 cells. Two independent *ESRP1*-OE clones were selected for each cell line. Expression values were normalized in each sample with those from the parental HCT116 and SW480 cell lines. HCT116 and SW480 cells transduced with the sh*ESRP1* lentivirus were induced by 1 µg/mL doxycycline for 72 hrs. Two independent clones were selected for each cell line. Expression values were normalized with those from non-induced cells; *GAPDH* expression was employed as control (Means±SEM, n=3).

**f.** CD44/EpCAM FACS analysis of HCT116 and SW480 EpCAM$^{lo}$ and EpCAM$^{hi}$ subpopulations in ESRP1-OE cells. Two independent clones are showed for each cell lines.

### Figure 2

**a.** Heatmap of common AS events between RNAseq data from *ESRP1*-KD in human non-small cell lung cancer cells (H358) and our HCT116 and SW480 EpCAM$^{hi}$ and EpCAM$^{lo}$ RNAseq data[17].

**b.** *CD44* and *NUMB* exon chromosome sites information from AS analysis. Exon peak plot depicts the expression of different exons in the three groups; peak height is indicative of the expression level of specific exons. CD44v: CD44 exons v2 to v10. CD44v and CD44s, and NUMB1/3 and NUMB2/4 (PRR region) exons are highlighted by gray rectangles.

20

**c.** RT-qPCR expression analysis of *CD44*s, *CD44*v6*, NUMB1/3* and *NUMB2/4* isoforms in HCT116 and SW480 EpCAM^hi, EpCAM^lo, and bulk subpopulations. *GAPDH* expression was employed as control (Means±SEM, n=3). ** = p<0.01.

**d.** Western analysis of CD44s, CD44v6 and NUMB isoforms in HCT116 and SW480 EpCAM^hi, EpCAM^lo, and bulk subpopulations. β-actin was used as loading control.


**Figure 3**

**a.** RT-qPCR and western analysis of CD44s, CD44v6 and NUMB isoforms expression in *ESRP1*-KD (sh*ESRP1* transduced) HCT116 cells. Two independent HCT116 *ESRP1*-KD clones were employed. Cells were induced with 1 µg/mL doxycycline for 72 hrs. before analysis. *GAPDH* expression was employed as qRT-PCR control (Means±SEM, n=3). ** = p<0.01. β-actin was used as loading control for western blots.

**b.** RT-qPCR and western analysis of CD44s, CD44v6 and NUMB isoforms expression in *ESRP1*-KD (sh*ESRP1* transduced) SW480 cells. Two independent SW480 *ESRP1*-KD clones were employed. Cells were induced with 1 µg/mL doxycycline for 72 hrs. before analysis. *GAPDH* expression was employed as qRT-PCR control (Means±SEM, n=3). ** = p<0.01. β-actin was used as loading control for western blots.

**c.** RT-qPCR and western analysis of CD44s, CD44v6, and NUMB isoforms expression in *ESRP1*-OE HCT116 cells. Two independent HCT116 *ESRP1*-OE clones were employed. *GAPDH* expression was employed as qRT-PCR control (Means±SEM, n=3). ** = p<0.01. β-actin was used as loading control for western blots.

**d.** RT-qPCR and western analysis of CD44s, CD44v6, and NUMB isoforms expression in *ESRP1*-OE SW480 cells. *GAPDH* expression was employed as qRT-PCR control (Means±SEM, n=3). ** = p<0.01. β-actin was used as loading control for western blots.


**Figure 4**

**a.** CD44/EpCAM FACS analysis of EpCAM^lo and EpCAM^hi subpopulations in CD44s-OE (left) and CD44v6-OE HCT116 and SW480 cell lines. The bar charts on the right depict the percentages of EpCAM^lo and EpCAM^hi cells.

**b.** and **c.** CD44/EpCAM FACS analysis of EpCAM^lo and EpCAM^hi subpopulations in NUMB1 to 4-OE HCT116 and SW480 cells. The bar charts on the right depict the percentages of EpCAM^lo and EpCAM^hi cells.

**d.** Macroscopic images of livers from mice spleen-injected with CD44s-, CD44v6-, NUMB2/4-, and NUMB1/3-OE HCT116 cells. HCT116 EpCAM^lo and bulk cells were used as positive control. Scale bar: 5 mm.

**e.** Liver metastasis multiplicity after intrasplenic injection of CD44s-, CD44v6-, NUMB2/4-, and NUMB1/3-OE HCT116 cells. For each transplantation experiment, $5 \times 10^4$ cells were injected in the spleen of recipient

NSG mouse. Six weeks after injection, mice were sacrificed and individual tumors counted. * = p<0.05; ** = p<0.01.

**Figure 5**

**a.** Principal Component Analysis (PCA) of RNAseq profiles from CD44s- and CD44v6-OE HCT116 and SW480 cell lines.

**b.** Heatmap of differentially expressed gene among HCT116 and SW480 CD44s-OE, CD44v6-OE, and parental cells.

**c.** Gene Set Enrichment Analysis (GSEA) of HCT116 and SW480 expression profiles in CD44s-OE compared with CD44v6-OE cells. Plots show only significantly altered pathways, with normalized enrichment score (NES) > 1, and pval < 0.01.

**Figure 6**

**a.** RNAseq data from the Cancer Genome Atlas (TCGA) were subdivided in 3 groups based on *ZEB1* and *ESRP1* expression level: $ZEB1^{hi}ESRP1^{lo}$ (*ZEB1*$^{hi}$, red dots), $ZEB1^{lo}ESRP1^{hi}$ (*ZEB1*$^{lo}$, blue dots), and intermediate (grey dots).

**b.** Kaplan Meier analysis of overall survival in the $ZEB1^{hi}ESRP1^{hi}$ and $ZEB1^{lo}ESRP1^{lo}$ patient groups.

**c.** Box plots showing CD44 and NUMB gene and isoforms expression across the $ZEB1^{hi}ESRP1^{lo}$, $ZEB1^{lo}ESRP1^{hi}$, and intermediate patient groups.

**d.** Dot plot analysis of the z-score scaled expression values of CD44s, CD44v6, NUMB1-4 isoforms across the 4 colon cancer consensus molecular subtypes (CMS).

**e**. Stacked bar plot showing the composition of the CMS subtypes across the $ZEB1^{hi/lo}$ and intermediate patient groups.

**Figure 7**

**a.** Gene correlation analysis showing the correlation of gene expression with CD44s and CD44v6 isoform expression in the TCGA patient cohort. Differentially expressed genes from CD44s- (red) and CD44v6-OE (blue) RNAseq data are highlighted.

**b.** Pathway correlation analysis showing the correlation of pathway activity CD44 and NUMB isoform expression in the TCGA patient cohort.

## *Supplementary Figure Legends*

**Figure S1**

**a.** CD44/EpCAM FACS analysis of EpCAM$^{lo}$ and EpCAM$^{hi}$ subpopulations in ESRP1-KD (*shESRP1*-transduced) HCT116 and SW480 cells. Cells were induced with 1 µg/mL doxycycline for 72 hrs. before analysis.

**b.** List of RBPs differentially expressed between EpCAM$^{lo}$ and EpCAM$^{hi}$ subpopulation in SW480 and HCT116. The RBPs list was from reference (10).

**c.** Dot plot analysis of the z-score scaled expression values of RBPs' expression across the 4 colon cancer consensus molecular subtypes.

**Figure S2**

Heatmap of alternative splicing analysis obtained by comparison of RNAseq data from RBP-KD studies (ESRP1-KD in H358[22], ESRP2-KD in LNCaP[23], RBM47-KD in H358[22], and QKI-KD in CAL27 [GEO Accession: GSM4677985]) with the EpCAM$^{hi/lo}$ RNAseq data[17]. Shared AS targets between RBPs KD cells, and HCT116, SW480 EpCAM$^{hi/lo}$ subpopulations are shown.

**Figure S3**

RT-PCR analysis of CD44 and NUMB isoforms expression in HCT116 (**a**) and SW480 (**b**) ESRP1-KD (*shESRP1*-transduced) cells, and in HCT116 (**c**) and SW480 (**d**) *ESRP1*-OE cells. Cells were induced with 1 µg/mL doxycycline for 72 hrs. before RNA isolation. *GAPDH* was used as control.

**e**. RT-PCR analysis of CD44s and CD44v6 expression in HCT116 and SW480 CD44s- (left), and CD44v6-OE (right) cells. *GAPDH* was used as control, normalized with the HCT116 and SW480 parental in each sample. (Means±SEM, n=3). ** = p<0.01.

**f.** Western analysis of CD44s and CD44v6 expression in HCT116 and SW480 CD44s- (left), CD44v6-OE (right) cells. β-actin was used as loading control for western blots.

**g**. RT-PCR analysis of NUMB1-4 isoforms expression in HCT116 and SW480 NUMB1-4 OE cells. *GAPDH* was used as control, normalized with the HCT116 and SW480 parental in each sample. (Means±SEM, n=3). ** = p<0.01.

**h**. Western analysis of NUMB1-4 isoforms expression in HCT116 and SW480 NUMB1-4 OE cells. β-actin was used as loading control for western blots.

**Figure S4**

**a.** Migration assay analysis of HCT116 CD44s-, CD44v6-, and NUMB1/4-OE cells. EpCAM$^{lo}$ and EpCAM$^{hi}$ cells were used as controls. Each bar represents the mean ± SD of cells migrated to the bottom of the transwell from two independent experiments. * = p<0.05, ** = p<0.01.

**b.** Migration assay analysis of SW480 CD44s-, CD44v6-, and NUMB1/4-OE cells. EpCAM$^{lo}$ and EpCAM$^{hi}$ cells were used as controls. Each bar represents the mean ± SD of cells migrated to the bottom of the transwell from two independent experiments. * = p<0.05, ** = p<0.01.

**c.** RT-qPCR analysis of EMT-TFs in HCT116 and SW480 CD44s-, CD44v6-, and NUMB1/4-OE cells. *GAPDH* expression was used as control, normalized with the HCT116 or SW480 parental in each sample (Means±SEM, n=3). Increased gene expression is depicted by red bars, whereas downregulation – when compared with parental cells- is shown by blue bar. * = p<0.05, ** = p<0.01.

**d.** RT-qPCR analysis of the Notch signaling pathway markers *HES1*, *HEY1*, and *ID2* in HCT116 and SW480 NUMB1/4-OE cells. *GAPDH* expression was used as control, normalized with the HCT116 or SW480 parental in each sample (Means±SEM, n=3). * = p<0.05, ** = p<0.01.

**Figure S5.**

**a.** Proliferation assays of HCT116 CD44s- and CD44v6-OE cells. Both OD values and cell multiplicities are shown from day 1 to 6 (Means±SEM, n=3). * = p<0.05; ** = p<0.01.

**b.** Proliferation assays of SW480 CD44s- and CD44v6-OE cells. Both OD values and cell multiplicities are shown from day 1 to 6 (Means±SEM, n=3). * = p<0.05; ** = p<0.01.

**c.** Proliferation assays of HCT116 NUMB1/4-OE cells. Both OD values and cell multiplicities are shown from day 1 to 6 (Means±SEM, n=3). ** = p<0.01.

**d.** Proliferation assays of SW480 NUMB1/4-OE cells. Both OD values and cell multiplicities are shown from day 1 to 6 (Means±SEM, n=3). ** = p<0.01.

**Figure S6**

**a.** Vulcano plots showing differentially expressed genes (absLFC > 2, pval < 0.01, red) by comparing parental cell lines to the CD44s- and CD44v6-OE samples in both cell lines.

**b.** Gene Set Enrichment Analysis of CD44s- and CD44v6-OE compared with their HCT116 and SW480 parental cells, respectively. Only significantly altered pathways (NES > 1, and pval < 0.01) are shown.

**Figure S7**

**a.** *CD44* and *NUMB* exon chromosome sites information from AS analysis in the ovarian and cervical cancer cell lines OV90 and SKOV6. Exon peak plot depicts the expression of different exons in the three groups; peak height is indicative of the expression level of specific exons. CD44v: CD44 exons v2 to v10. CD44v and CD44s, and NUMB1/3 and NUMB2/4 (PRR region) exons are highlighted by gray rectangles.

**b.** RT-qPCR expression analysis of *ESRP1*, *CD44s*, *CD44v6*, *NUMB1/3*, and *NUMB2/4* isoforms in EpCAM^hi, EpCAM^lo, and bulk subpopulations in OV90 and SKOV6 ovarian cancer cell lines. *GAPDH* expression was used as control (Means±SEM, n=3). ** = p<0.01.


***Supplementary Table Legends***

**Table S1.** List of alternative splicing targets in ESRP1 knocking down H358 line, HCT116 and SW480 EpCAM^lo and EpCAM^hi subpopulation, filtered by ΔPSI > 0.1

**Table S2.** List of alternative splicing targets in *ESRP1*-KD in the H358 cell line, *ESRP2*-KD in LNCaP, *RBM*47-KD in H358 line, QKI-KD in CAL27, and HCT116 and SW480 EpCAM^lo and EpCAM^hi subpopulation, filtered by ΔPSI > 0.1.

**Table S3.** Lists of primer sequences used for RT-PCR analysis.

**Table S4.** Differential expressed gene lists from the RNAseq analysis HCT116 CD44s- and CD44v6-OE cells.

**Table S5.** Differential expressed gene lists from the RNAseq analysis SW480 CD44s- and CD44v6-OE cells.
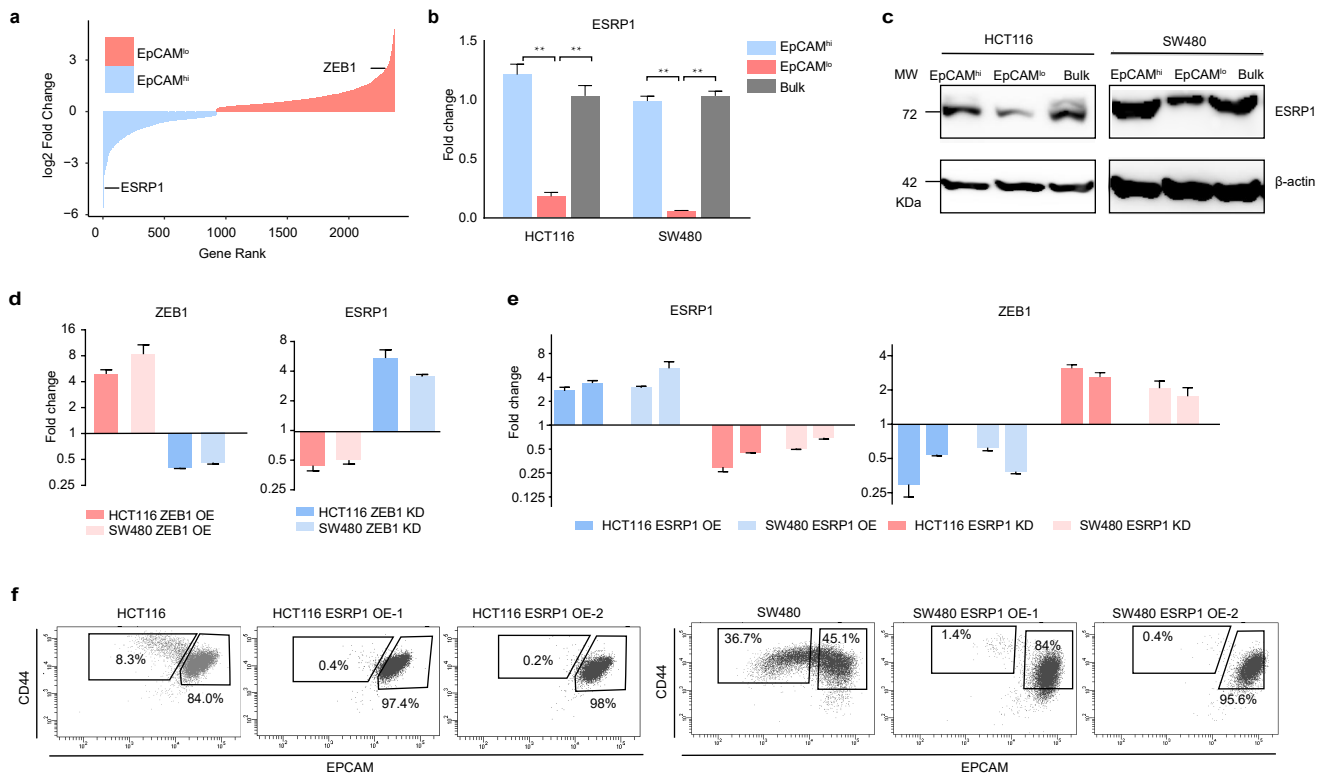
## References

1       Hanahan, D. & Weinberg, R. A. Hallmarks of cancer: the next generation. *Cell* **144**, 646-674, doi:10.1016/j.cell.2011.02.013 (2011).

2       Bernards, R. & Weinberg, R. A. A progression puzzle. *Nature* **418**, 823, doi:10.1038/418823a (2002).

3       Reiter, J. G. *et al.* Minimal functional driver gene heterogeneity among untreated metastases. *Science* **361**, 1033-1037, doi:10.1126/science.aat7171 (2018).

4       Thiery, J. P., Acloque, H., Huang, R. Y. & Nieto, M. A. Epithelial-mesenchymal transitions in development and disease. *Cell* **139**, 871-890, doi:10.1016/j.cell.2009.11.007 (2009).

5       Brabletz, T., Jung, A., Spaderna, S., Hlubek, F. & Kirchner, T. Opinion: migrating cancer stem cells - an integrated concept of malignant tumour progression. *Nat Rev Cancer* **5**, 744-749, doi:10.1038/nrc1694 (2005).

6       Fodde, R. & Brabletz, T. Wnt/beta-catenin signaling in cancer stemness and malignant behavior. *Curr Opin Cell Biol* **19**, 150-158, doi:10.1016/j.ceb.2007.02.007 (2007).

7       Teeuwssen, M. & Fodde, R. Cell Heterogeneity and Phenotypic Plasticity in Metastasis Formation: The Case of Colon Cancer. *Cancers (Basel)* **11**, doi:10.3390/cancers11091368 (2019).

8       Wang, E. T. *et al.* Alternative isoform regulation in human tissue transcriptomes. *Nature* **456**, 470-476, doi:10.1038/nature07509 (2008).

9       Blencowe, B. J. Alternative splicing: new insights from global analyses. *Cell* **126**, 37-47, doi:10.1016/j.cell.2006.06.023 (2006).

10      Fu, X. D. & Ares, M., Jr. Context-dependent control of alternative splicing by RNA-binding proteins. *Nature reviews. Genetics* **15**, 689-701, doi:10.1038/nrg3778 (2014).

11      Kahles, A. *et al.* Comprehensive Analysis of Alternative Splicing Across Tumors from 8,705 Patients. *Cancer Cell* **34**, 211-224 e216, doi:10.1016/j.ccell.2018.07.001 (2018).

12      Roy Burman, D., Das, S., Das, C. & Bhattacharya, R. Alternative splicing modulates cancer aggressiveness: role in EMT/metastasis and chemoresistance. *Mol Biol Rep* **48**, 897-914, doi:10.1007/s11033-020-06094-y (2021).

13      Oltean, S. & Bates, D. O. Hallmarks of alternative splicing in cancer. *Oncogene* **33**, 5311-5318, doi:10.1038/onc.2013.533 (2014).

14      Biamonti, G., Infantino, L., Gaglio, D. & Amato, A. An Intricate Connection between Alternative Splicing and Phenotypic Plasticity in Development and Cancer. *Cells* **9**, doi:10.3390/cells9010034 (2019).

15      Brown, R. L. *et al.* CD44 splice isoform switching in human and mouse epithelium is essential for epithelial-mesenchymal transition and breast cancer progression. *The Journal of clinical investigation* **121**, 1064-1074, doi:10.1172/JCI44540 (2011).

16      Yae, T. *et al.* Alternative splicing of CD44 mRNA by ESRP1 enhances lung colonization of metastatic cancer cell. *Nature communications* **3**, 883, doi:10.1038/ncomms1892 (2012).

17      Sacchetti, A. *et al.* Phenotypic plasticity underlies local invasion and distant metastasis in colon cancer. *Elife* **10**, doi:10.7554/eLife.61461 (2021).

18      Preca, B. T. *et al.* A self-enforcing CD44s/ZEB1 feedback loop maintains EMT and stemness properties in cancer cells. *Int J Cancer* **137**, 2566-2577, doi:10.1002/ijc.29642 (2015).

19      Cook, K. B., Kazan, H., Zuberi, K., Morris, Q. & Hughes, T. R. RBPDB: a database of RNA-binding specificities. *Nucleic Acids Res* **39**, D301-308, doi:10.1093/nar/gkq1069 (2011).

20      Guinney, J. *et al.* The consensus molecular subtypes of colorectal cancer. *Nature medicine* **21**, 1350-1356, doi:10.1038/nm.3967 (2015).

21      Schafer, S. *et al.* Alternative Splicing Signatures in RNA-seq Data: Percent Spliced in (PSI). *Curr Protoc Hum Genet* **87**, 11 16 11-11 16 14, doi:10.1002/0471142905.hg1116s87 (2015).

22      Yang, Y. *et al.* Determination of a Comprehensive Alternative Splicing Regulatory Network and Combinatorial Regulation by Key Factors during the Epithelial-to-Mesenchymal Transition. *Molecular and cellular biology* **36**, 1704-1719, doi:10.1128/MCB.00019-16 (2016).

23      <Androgen-regulated transcription of ESRP2 drives alternative splicing patterns in prostate cancer.pdf>. doi:10.7554/eLife.47678.001 10.7554/eLife.47678.002.

24      Hernandez-Martinez, R., Ramkumar, N. & Anderson, K. V. p120-catenin regulates WNT signaling and EMT in the mouse embryo. *Proc Natl Acad Sci U S A* **116**, 16872-16881, doi:10.1073/pnas.1902843116 (2019).

25      Shimada, H. *et al.* The Roles of Tricellular Tight Junction Protein Angulin-1/Lipolysis-Stimulated Lipoprotein Receptor (LSR) in Endometriosis and Endometrioid-Endometrial Carcinoma. *Cancers (Basel)* **13**, doi:10.3390/cancers13246341 (2021).

26      Conway, J., Al-Zahrani, K. N., Pryce, B. R., Abou-Hamad, J. & Sabourin, L. A. Transforming growth factor beta-induced epithelial to mesenchymal transition requires the Ste20-like kinase SLK independently of its catalytic activity. *Oncotarget* **8**, 98745-98756, doi:10.18632/oncotarget.21928 (2017).

27      Karve, K., Netherton, S., Deng, L., Bonni, A. & Bonni, S. Regulation of epithelial-mesenchymal transition and organoid morphogenesis by a novel TGFbeta-TCF7L2 isoform-specific signaling pathway. *Cell Death Dis* **11**, 704, doi:10.1038/s41419-020-02905-z (2020).

28      Orian-Rousseau, V. CD44 Acts as a Signaling Platform Controlling Tumor Progression and Metastasis. *Frontiers in immunology* **6**, 154, doi:10.3389/fimmu.2015.00154 (2015).

29      Pece, S., Confalonieri, S., P, R. R. & Di Fiore, P. P. NUMB-ing down cancer by more than just a NOTCH. *Biochimica et biophysica acta* **1815**, 26-43, doi:10.1016/j.bbcan.2010.10.001 (2011).

30      Todaro, M. *et al.* CD44v6 is a marker of constitutive and reprogrammed cancer stem cells driving colon cancer metastasis. *Cell Stem Cell* **14**, 342-356, doi:10.1016/j.stem.2014.01.009 (2014).

31      Liberzon, A. *et al.* The Molecular Signatures Database (MSigDB) hallmark gene set collection. *Cell Syst* **1**, 417-425, doi:10.1016/j.cels.2015.12.004 S2405-4712(15)00218-5 [pii] (2015).

32      Adam, R. A. & Adam, Y. G. Malignant ascites: past, present, and future. *Journal of the American College of Surgeons* **198**, 999-1011, doi:10.1016/j.jamcollsurg.2004.01.035 (2004).

33      <Numb negatively regulates the epithelial-to-mesenchymal transition in colorectal cancer through the Wnt signalling pathway..pdf>.

34      Varga, J. & Greten, F. R. Cell plasticity in epithelial homeostasis and tumorigenesis. *Nature cell biology* **19**, 1133-1141, doi:10.1038/ncb3611 (2017).

35      Dixit, A. *et al.* Perturb-Seq: Dissecting Molecular Circuits with Scalable Single-Cell RNA Profiling of Pooled Genetic Screens. *Cell* **167**, 1853-1866 e1817, doi:10.1016/j.cell.2016.11.038 (2016).

36      Lengauer, C., Kinzler, K. W. & Vogelstein, B. Genetic instability in colorectal cancers. *Nature* **386**, 623-627, doi:10.1038/386623a0 (1997).

37      Warzecha, C. C., Sato, T. K., Nabet, B., Hogenesch, J. B. & Carstens, R. P. ESRP1 and ESRP2 are epithelial cell-type-specific regulators of FGFR2 splicing. *Mol Cell* **33**, 591-601, doi:10.1016/j.molcel.2009.01.025 (2009).

38      Tavanez, J. P. & Valcarcel, J. A splicing mastermind for EMT. *EMBO J* **29**, 3217-3218, doi:10.1038/emboj.2010.234 (2010).

39      Warzecha, C. C. *et al.* An ESRP-regulated splicing programme is abrogated during the epithelial-mesenchymal transition. *EMBO J* **29**, 3286-3300, doi:10.1038/emboj.2010.195 (2010).

40      Calon, A. *et al.* Stromal gene expression defines poor-prognosis subtypes in colorectal cancer. *Nat Genet* **47**, 320-329, doi:10.1038/ng.3225 (2015).

41      Isella, C. *et al.* Stromal contribution to the colorectal cancer transcriptome. *Nat Genet* **47**, 312-319, doi:10.1038/ng.3224 (2015).

42    Pillman, K. A. *et al.* miR-200/375 control epithelial plasticity-associated alternative splicing by repressing the RNA-binding protein Quaking. *EMBO J* **37**, doi:10.15252/embj.201899016 (2018).

43    Kim, E. J. *et al.* QKI, a miR-200 target gene, suppresses epithelial-to-mesenchymal transition and tumor growth. *Int J Cancer* **145**, 1585-1595, doi:10.1002/ijc.32372 (2019).

44    Brabletz, S. & Brabletz, T. The ZEB/miR-200 feedback loop--a motor of cellular plasticity in development and cancer? *EMBO Rep* **11**, 670-677, doi:10.1038/embor.2010.117 (2010).

45    Xia, R. M., Liu, T., Li, W. G. & Xu, X. Q. RNA-binding protein RBM24 represses colorectal tumourigenesis by stabilising PTEN mRNA. *Clin Transl Med* **11**, e383, doi:10.1002/ctm2.383 (2021).

46    Lin, G. *et al.* RNA-binding Protein MBNL2 regulates Cancer Cell Metastasis through MiR-182-MBNL2-AKT Pathway. *J Cancer* **12**, 6715-6726, doi:10.7150/jca.62816 (2021).

47    Zhang, J. *et al.* The natural compound neobractatin inhibits tumor metastasis by upregulating the RNA-binding-protein MBNL2. *Cell Death Dis* **10**, 554, doi:10.1038/s41419-019-1789-5 (2019).

48    Lu, Z. X. *et al.* Transcriptome-wide landscape of pre-mRNA alternative splicing associated with metastatic colonization. *Mol Cancer Res* **13**, 305-318, doi:10.1158/1541-7786.MCR-14-0366 (2015).

49    Xiao, H. MiR-7-5p suppresses tumor metastasis of non-small cell lung cancer by targeting NOVA2. *Cell Mol Biol Lett* **24**, 60, doi:10.1186/s11658-019-0188-3 (2019).

50    Zeilstra, J. *et al.* Deletion of the WNT target and cancer stem cell marker CD44 in Apc(Min/+) mice attenuates intestinal tumorigenesis. *Cancer Res* **68**, 3655-3661, doi:10.1158/0008-5472.CAN-07-2940 (2008).

51    Zeilstra, J. *et al.* Stem cell CD44v isoforms promote intestinal cancer formation in Apc(min) mice downstream of Wnt signaling. *Oncogene* **33**, 665-670, doi:10.1038/onc.2012.611 (2014).

52    Misra, S., Hascall, V. C., De Giovanni, C., Markwald, R. R. & Ghatak, S. Delivery of CD44 shRNA/nanoparticles within cancer cells: perturbation of hyaluronan/CD44v6 interactions and reduction in adenoma growth in Apc Min/+ MICE. *J Biol Chem* **284**, 12432-12446, doi:10.1074/jbc.M806772200 (2009).

53    Orian-Rousseau, V., Chen, L., Sleeman, J. P., Herrlich, P. & Ponta, H. CD44 is required for two consecutive steps in HGF/c-Met signaling. *Genes Dev* **16**, 3074-3086, doi:10.1101/gad.242602 (2002).

54    Zhang, H. *et al.* CD44 splice isoform switching determines breast cancer stem cell state. *Genes Dev* **33**, 166-179, doi:10.1101/gad.319889.118 (2019).

55    Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15-21, doi:10.1093/bioinformatics/bts635 (2013).

56    Katz, Y., Wang, E. T., Airoldi, E. M. & Burge, C. B. Analysis and design of RNA sequencing experiments for identifying isoform regulation. *Nature Methods* **7**, 1009-1015, doi:10.1038/nmeth.1528 (2010).

Fig 1



**a**

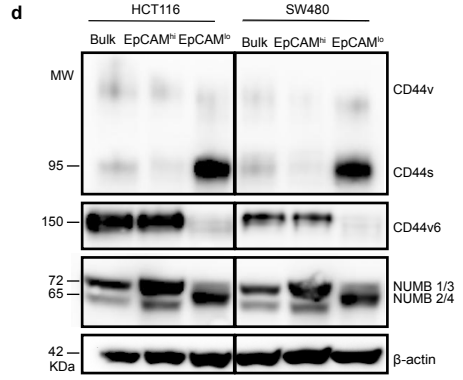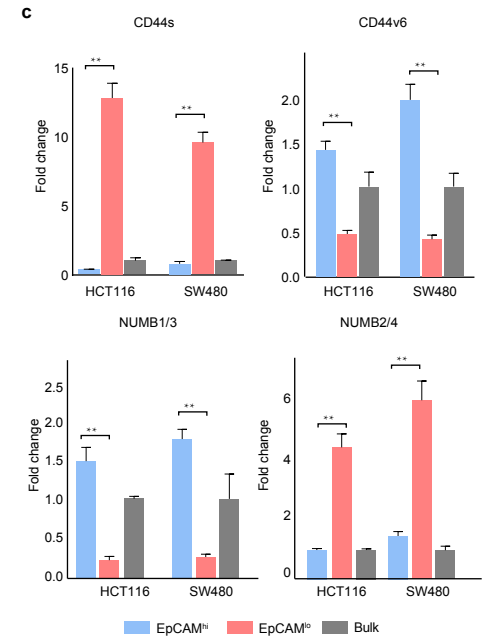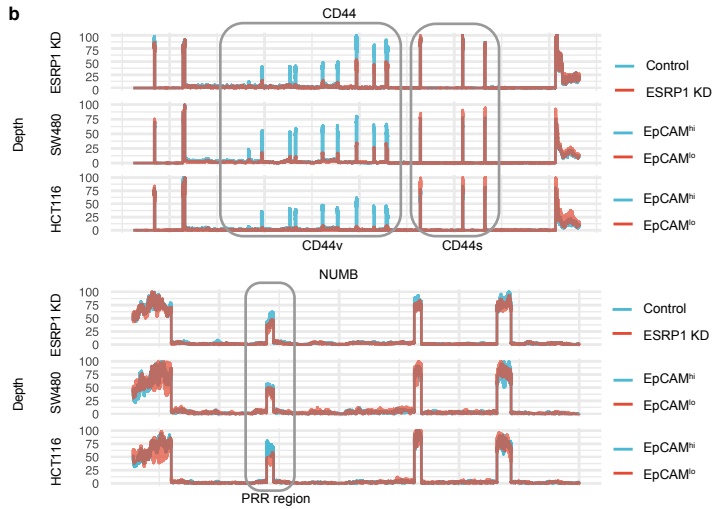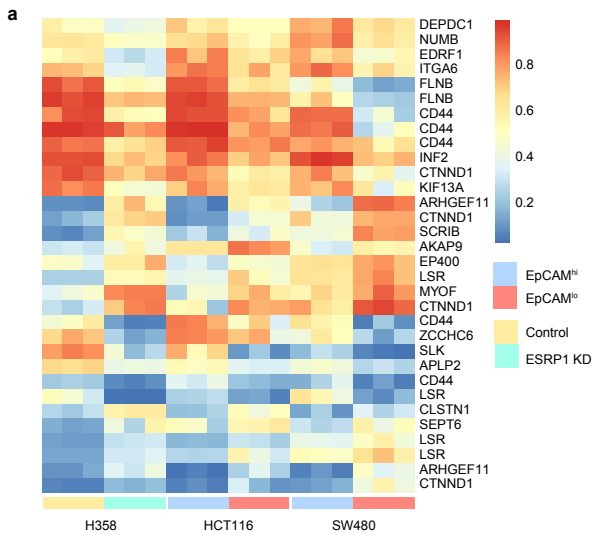log2 Fold Change

EpCAM^lo
EpCAM^hi

ZEB1
ESRP1

Gene Rank

**b**

ESRP1

Fold change

EpCAM^hi
EpCAM^lo
Bulk

**
**
**
**

HCT116
SW480

**c**

HCT116
SW480

MW    EpCAM^hi  EpCAM^lo  Bulk    EpCAM^hi  EpCAM^lo  Bulk

72                                                          ESRP1

42
KDa                                                         β-actin

**d**

ZEB1

Fold change

ESRP1

HCT116 ZEB1 OE
SW480 ZEB1 OE

HCT116 ZEB1 KD
SW480 ZEB1 KD

**e**

ESRP1

Fold change

ZEB1

HCT116 ESRP1 OE
SW480 ESRP1 OE
HCT116 ESRP1 KD
SW480 ESRP1 KD

**f**

HCT116          HCT116 ESRP1 OE-1     HCT116 ESRP1 OE-2

CD44

8.3%            0.4%                  0.2%

84.0%           97.4%                 98%

EPCAM

SW480           SW480 ESRP1 OE-1      SW480 ESRP1 OE-2

CD44

36.7%  45.1%    1.4%    84%          0.4%

95.6%

EPCAM

Fig 2



**a** Heatmap with gene labels: DEPDC1, NUMB, EDRF1, ITGA6, FLNB, FLNB, CD44, CD44, CD44, INF2, CTNND1, KIF13A, ARHGEF11, CTNND1, SCRIB, AKAP9, EP400, LSR, MYOF, CTNND1, CD44, ZCCHC6, SLK, APLP2, CD44, LSR, CLSTN1, SEPT6, LSR, LSR, ARHGEF11, CTNND1

Color scale: 0.8, 0.6, 0.4, 0.2

Legend: EpCAM^hi, EpCAM^lo, Control, ESRP1 KD

Sample labels: H358, HCT116, SW480

**b** CD44 / NUMB depth plots for ESRP1 KD, SW480, HCT116

Legend: Control, ESRP1 KD, EpCAM^hi, EpCAM^lo

Regions: CD44, CD44v, CD44s, NUMB, PRR region

**c** Bar charts: CD44s, CD44v6, NUMB1/3, NUMB2/4 (Fold change, HCT116, SW480)

Legend: EpCAM^hi, EpCAM^lo, Bulk

**d** Western blots: HCT116, SW480 (Bulk, EpCAM^hi, EpCAM^lo)

MW markers: 95, 150, 72, 65, 42 KDa

Protein labels: CD44v, CD44s, CD44v6, NUMB 1/3, NUMB 2/4, β-actin

Fig 3

# Fig 4

**a**



HCT116 parental — CD44 / EPCAM — 14.0% / 75.3%
CD44s OE — 56.4% / 30.8%
CD44v6 OE — 9.7% / 85.5%

SW480 parental — 35.8% / 50.6%
CD44s OE — 94.8% / 1.6%
CD44v6 OE — 15.8% / 70.0%

Legend: □ Intermediate ▨ EpCAM$^{lo}$ ■ EpCAM$^{hi}$

**b**



HCT116 parental — 12.4% / 78.4%
NUMB1 OE — 1.5% / 94.2%
NUMB2 OE — 89.0% / 5.3%
NUMB3 OE — 4.6% / 87.8%
NUMB4 OE — 49.7% / 45.0%
NUMB1/3 OE — 1.5% / 93.8%
NUMB2/4 OE — 51.2% / 41.3%

**c**



SW480 parental — 32.8% / 55.6%
NUMB1 OE — 23.7% / 69.2%
NUMB2 OE — 58.0% / 38.0%
NUMB3 OE — 13.5% / 75.6%
NUMB4 OE — 85.5% / 8.0%
NUMB1/3 OE — 12.6% / 75.0%
NUMB2/4 OE — 95.9% / 1.9%

**d**



Bulk    CD44v6    CD44s

EpCAM$^{lo}$    NUMB1/3    NUMB2/4

**e**



Number of liver metastases

HCT116, EpCAM$^{lo}$, CD44v6 OE, CD44s OE, NUMB1/3 OE, NUMB2/4 OE

Fig 5

Fig 6

**a**

Gene correlation analysis in primary CRC tumors

Gene Group ● CD44s OE ● CD44v6 OE ● Other

**b**

Figure with panels CD44 and NUMB

isoform
● CD44s
● CD44v6
● NUMB1
● NUMB2
● NUMB3
● NUMB4

Pearson Correlation in TCGA COAD

Fig. S1

**a**



**b**

| Gene Name | log2FC_HCT116 | log2FC_SW480 |
|-----------|---------------|--------------|
| ESRP1 | -2.22645 | -6.89236 |
| ESRP2 | -0.896173 | -1.89007 |
| NOVA2 | 0.757111 | 2.2421 |
| RBM47 | -1.36442 | -0.900345 |
| MBNL2 | 0.784795 | 0.355627 |
| MBNL3 | -0.637449 | -2.0392 |
| RBM24 | 1.52564 | NS |
| RBM19 | -0.496814 | NS |
| HNRNPAB | -0.452521 | NS |
| HNRNPF | -0.352169 | NS |
| RBM43 | 0.490142 | NS |
| U2AF2 | -0.385724 | NS |
| RBM14 | NS | -0.458869 |
| QKI | NS | 0.643307 |
| SRSF5 | NS | 0.412868 |
| HNRNPH1 | NS | 0.244149 |

**c**

Fig. S2

Fig. S3

Fig. S4

Fig. S5

Fig. S6