

1 **Smoking is Associated with Hypermethylation of the APC 1A Promoter in**
2 **Colorectal Cancer: the ColoCare Study**

3
4 Timothy M. Barrow^{1,2,3*}, Hagen Klett^{1,3,4*}, Reka Toth^{1,5}, Jürgen Böhm^{1,5}, Biljana
5 Gigic^{1,5,6}, Nina Habermann^{1,5}, Dominique Scherer^{1,5,7}, Petra Schrotz-King^{1,5}, Stephanie
6 Skender^{1,5}, Clare Abbenhardt-Martin^{1,5}, Lin Zielske^{1,5}, Martin Schneider⁶, Alexis Ulrich⁶,
7 Peter Schirmacher^{1,8}, Esther Herpel^{1,8}, Hermann Brenner^{1,5,9}, Hauke Busch^{1,3,4†},
8 Melanie Boerries^{1,3,4†}, Cornelia M. Ulrich^{1,5,10†}, Karin B. Michels^{1,2,3,11†}

9
10 ¹German Cancer Consortium (DKTK), Heidelberg, Germany

11 ²Institute for Prevention and Cancer Epidemiology, Faculty of Medicine and Medical
12 Center, University of Freiburg, Freiburg, Germany

13 ³German Cancer Research Center (DKFZ), Heidelberg, Germany

14 ⁴Institute of Molecular Medicine and Cell Research, Faculty of Medicine and Medical
15 Center, University of Freiburg, Freiburg, Germany

16 ⁵Division of Preventive Oncology, National Center for Tumor Diseases (NCT) and
17 German Cancer Research Center (DKFZ), Heidelberg, Germany

18 ⁶Department of Surgical Oncology, University Clinic Heidelberg, Heidelberg, Germany

19 ⁷Institute of Medical Biometry and Informatics, University of Heidelberg, Heidelberg,
20 Germany

21 ⁸Department of General Pathology, University Clinic Heidelberg, Heidelberg, Germany

22 ⁹Division of Clinical Epidemiology and Aging Research, German Cancer Research
23 Center (DKFZ), Heidelberg, Germany

24 ¹⁰Population Sciences, Huntsman Cancer Institute, Salt Lake City, Utah, USA.

1 ¹¹Department of Epidemiology, UCLA Fielding School of Public Health, Los Angeles,
2 CA, USA.

3 *These authors contributed equally to the work (First authors)

4 †These authors contributed equally to the work (Last authors)

5

6 *Correspondence to: Prof. Karin Michels, Institute for Prevention and Cancer
7 Epidemiology, Faculty of Medicine and Medical Center, University of Freiburg, Freiburg,
8 Germany. Tel: +49 270 77360, Fax: +49 270 77340, Email:
9 *tumorepidemiologie@uniklinik-freiburg.de.*

10

11 The authors report no conflicts of interest.

12

13 Running title: Smoking and APC promoter hypermethylation in CRC

14

15 Word count: 3,999

16

17 Raw data will be supplied for publication as Supplementary Materials.

1 **Abstract**

2 Smoking tobacco is a known risk factor for the development of colorectal cancer, and for
3 mortality associated with the disease. While smoking has been reported to be associated
4 with changes in DNA methylation in blood and in lung tumour tissues, there has been
5 scant investigation of how epigenetic factors may be implicated in the increased risk of
6 developing colorectal cancer. To identify epigenetic changes associated with smoking
7 behaviours, we performed epigenome-wide analysis of DNA methylation in colorectal
8 tumours from 36 never smokers, 47 former smokers and 13 active smokers, and adjacent
9 mucosa from 49 never smokers, 64 former smokers and 18 active smokers. Our analyses
10 identified 15 CpG sites within the *APC* 1A promoter that were significantly
11 hypermethylated and 14 CpG loci within the *NFATC1* gene body that were significantly
12 hypomethylated ($p_{LIS} < 1 \times 10^{-5}$) in tumours of active smokers. The *APC* 1A promoter was
13 hypermethylated in 7 of 36 tumours from never smokers (19%), 12 of 47 tumours from
14 former smokers (26%), and 8 of 13 tumours from active smokers (62%). Promoter
15 hypermethylation was positively associated with duration of smoking (Spearman rank
16 correlation, $\rho=0.26$, $p=0.03$) and was confined to tumours, with hypermethylation never
17 observed in adjacent mucosa. Further analysis of adjacent mucosa revealed significant
18 hypomethylation of four loci associated with the *TNXB* gene in tissue from active smokers.
19 Our findings provide exploratory evidence for hypermethylation of the key tumour
20 suppressor gene *APC* being implicated in smoking-associated colorectal carcinogenesis.
21 Further work is required to establish the validity of our observations in independent
22 cohorts.

23

24 Keywords: Smoking; Tobacco; Colorectal cancer; Epigenetics; DNA methylation; APC.

25

1 **Introduction**

2 Smoking tobacco is a risk factor for many forms of cancer, including colorectal cancer
3 (CRC). Ever-smokers, which includes both current and former smokers, have an 18%
4 increase in risk of developing the disease relative to individuals who have never smoked
5 [1], and the risk is greatest for the development of tumours in the rectum. In addition to
6 increased incidence, active smokers have a 23% greater risk of CRC-related mortality [2]
7 and patients who are former smokers still display increased risk of all-cause mortality [3].
8 The duration and intensity of smoking are known to modify risk, with individuals who have
9 smoked for ≥ 30 years and those with ≥ 20 pack-years of smoking each displaying a 40%
10 increase in risk of CRC-related mortality [3]. However, the mechanisms by which smoking
11 tobacco increases CRC risk have not been elucidated. It has been hypothesised that the
12 carcinogenic products of cigarette smoke may reach the colorectum through the blood
13 and be implicated in the early initiation of cancer, as opposed to furthering the
14 development of existing adenomas [4].

15 Smoking is associated with alterations in DNA methylation, an epigenetic modifier
16 of gene expression, in healthy individuals. Such epigenetic events display tissue-
17 specificity [5] and differ by ethnicity [6,7], and can serve as markers of long-term exposure
18 to tobacco smoke [8]. Several studies examining the blood of smokers have reported
19 differential methylation of loci within the aryl hydrocarbon receptor repressor (*AHRR*)
20 gene [6,9], a putative tumour-suppressor which mediates the detoxification of products in
21 cigarette smoke, and the coagulation factor II (thrombin) receptor-like 3 (*F2RL3*) gene
22 [6,8-10], implicated in blood clotting. Associations have been identified between smoking-
23 related changes in DNA methylation of *AHRR*, *F2RL3* and LINE1 elements measured in
24 blood and the risk of cancer [11] and mortality from the disease [12].

1 Further to these observations in healthy individuals, there is evidence that smoking
2 is associated with epigenetic changes in tumour tissue. Epigenome-wide association
3 studies have identified distinct methylation profiles in lung tumours from smokers and
4 non-smokers [13], while candidate-gene approaches have identified smoking-related
5 changes in the methylation of cyclin-dependent kinase inhibitor 2A (*CDKN2A/p16*) and
6 runt-related transcription factor 3 (*RUNX3*) in bladder tumours [14,15] and *CDKN2A/p16*
7 and O-6-methylguanine-DNA methyltransferase (*MGMT*) in lung tumours [16]. Smoking-
8 related epigenetic events may occur early in carcinogenesis, as demonstrated by their
9 observation in stage I non-small cell lung cancers [17]. However, the evidence for
10 smoking-associated epigenetic dysregulation in CRC is currently limited. Smoking has
11 been reported as associated with microsatellite instability and positive CpG island
12 methylator phenotype (CIMP) status [18], but there has otherwise been scant research of
13 DNA methylation in colorectal tumours by smoking status.

14 In this study, we investigated whether epigenetic factors may be implicated in the
15 increased risk of CRC among tobacco smokers by analysing epigenetic patterns in
16 colorectal tumours and neighbouring mucosa in relation to smoking behaviours. We
17 utilised the Illumina HumanMethylation450 microarray platform to analyse DNA
18 methylation in samples taken from a total of 137 colorectal cancer patients, 51 of whom
19 had never smoked ('never smokers'), 68 who had been smokers but had ceased at least
20 two years prior to cancer diagnosis ('former smokers'), and 18 who smoked at the point
21 of diagnosis ('active smokers'). We report that promoter 1A of the *APC* gene, commonly
22 inactivated in CRC, is hypermethylated in the tumours of active smokers. Methylation of
23 this region is associated with duration of smoking, and hypermethylation ($\beta > 0.2$) was
24 never observed in adjacent mucosa. Our results suggest that the increased risk of CRC

1 development among smokers may progress through epigenetic inactivation of the key
2 tumour suppressor gene *APC*.

3

4

1 **Material and Methods**

2 **The ColoCare Study**

3 The ColoCare consortium is a multicentre initiative of interdisciplinary research on
4 outcomes associated with colorectal cancer, with sites at the Fred Hutchison Cancer
5 Research Center (Seattle, USA), Moffit Cancer Center (Tampa, USA), and from 2010 at
6 the German Cancer Research Center (Heidelberg, Germany). This study exclusively
7 focussed upon patients recruited in Heidelberg. ColoCare has been approved by the
8 ethics committee of the University of Heidelberg medical faculty. Patients were enrolled
9 to this prospective cohort at the point of diagnosis, having given informed consent, with
10 biospecimens and data collected at regularly scheduled intervals of 3, 6, 12, 24 and 36
11 months post-surgery. Medical factors were abstracted from patients' charts and records
12 from the University Hospital of Heidelberg. Data on dietary habits, exercise and physical
13 activity, smoking habits, medication, socio-demographic information, and quality of life
14 were collected via questionnaires. To date, 500 patients have been recruited at the
15 Heidelberg site.

16

17 **Tissue samples**

18 Tissue samples were collected from patients undergoing surgery at the University
19 Hospital of Heidelberg, and were reviewed by pathologists to ensure their quality and
20 origin. Tumour samples were collected from 36 patients who had never smoked, 47 who
21 were former smokers, and 13 who were active smokers at the point of diagnosis. Mucosa
22 was taken from adjacent to tumours from 49 never smokers, 64 former smokers and 18
23 active smokers. A summary of patient characteristics is provided in Table 1.

24

25 **DNA isolation**

1 DNA was extracted from fresh-frozen tissue using the QIAamp AllPrep DNA/RNA mini kit
2 (Qiagen) according to the manufacturer's instructions.

3
4 **Illumina Infinium HumanMethylation450 BeadChip microarrays**
5 DNA microarrays were performed at the Genomics and Proteomics Core Facility at the
6 German Cancer Research Center (Heidelberg, Germany). 1.0µg of Genomic DNA was
7 bisulfite-converted using the EZ DNA Methylation kit (Zymo Research) according to the
8 manufacturer's instructions. The microarrays were then performed according to the
9 Illumina Infinium HD Methylation protocol.

10
11 **Microarray data analysis**

12 Microarray data was pre-processed using the Illumina Genome Studio software program
13 before analysis using the R minfi package. Background correction and dye-bias
14 normalisation were performed using noob [19], and functional normalisation was
15 performed to remove batch effects and inner technical variability and adjust for Type I/II
16 probe fluorescence effect, as described elsewhere [20]. Prior to background correction
17 and normalisation, probes with detection p values >0.01 in 10% of samples (n=662) or
18 bead counts less than three in 10% of samples (n=162) were removed. Probes with SNPs
19 within 10bp of the target CpG with minor allele frequencies of >0.01 (n=19,099) and
20 mapping to the X and Y chromosomes (n=11,150) were removed. Subsequently, a total
21 of 456,144 probes were taken forward for analysis.

22 Loci that are differentially methylated by smoking behaviours were identified by
23 fitting a linear least-squares regression model across the conditions followed by
24 computing moderated t-statistics for every CpG site, as described in the limma pipeline
25 [21]. Due to the non-independent structure of the univariate t-statistics, we used a non-

1 homogenous hidden Markov model (NHMM) to incorporate the dependence coming from
2 the chromosomal positions of CpGs in the test statistics, as proposed and described
3 elsewhere [22]. Briefly, t-statistics are z-score transformed and distances (base pairs)
4 between CpGs were calculated and used as dependence structure in the NHMM. The
5 NHMM parameters were estimated by expectation maximisation with randomised initial
6 values. To avoid local maxima in the maximisation algorithm we used 30 initialisations
7 and chose the initialisation with the smallest Bayesian information criteria (BIC). This
8 provides a reproducible local index of significance (LIS), as previously defined [23], and
9 can be interpreted as dependence corrected p value (pLIS). For computational efficiency
10 we performed the analysis by chromosome and pooled the results afterwards, with
11 significance defined as $pLIS < 1 \times 10^{-5}$. The pLIS scores were computed using the R
12 package NHMMfdr. Comparisons were made between never smokers and active
13 smokers and between never smokers and former smokers in tumour and adjacent
14 mucosa tissues. To identify loci that are differentially methylated between tumours and
15 adjacent mucosa in a smoking-specific manner, we compared the differences in active
16 smokers of tumour and mucosa with the differences among never smokers of tumour and
17 mucosa. All analyses were adjusted for age and sex in the linear regression model.

18 The methylation microarray dataset is available from the NCBI Gene Expression
19 Omnibus repository (accession number: GSE101764).

20

21 **Identification of probe-associated SNPs**

22 To account for false positives stemming from genetic variation, we used the UCSC
23 Genome Browser and NCBI dbSNP databases [24,25] to identify single nucleotide
24 polymorphisms (SNPs) within the 50-mer probes of the microarray for sites identified as
25 significantly differentially methylated by smoking behaviours. The unconverted DNA

1 sequences ('SourceSeq') for each significantly-different probe in tumour tissue and
2 adjacent mucosa were extracted from the GenomeStudio output file and were used to
3 perform a BLAT search using the UCSC Genome Browser [24]. The minor allele
4 frequencies for all SNPs located within the probe sequences were identified using the
5 UCSC Genome Browser and the NCBI dbSNP database [24,25]. Data from across all
6 ethnicities or, where available, European populations was recorded, using estimates
7 from studies with the largest sample sizes.

8

9 **Statistical analyses**

10 Associations between DNA methylation and smoking habits were calculated using data
11 on pack-years and duration of smoking for each patient, and time since cessation among
12 former smokers. Detailed data on smoking habits was available for 87 patients from whom
13 tumour tissue was taken and 115 patients providing adjacent mucosa. Associations
14 between DNA methylation (beta values) and intensity (pack-years) and duration (years)
15 of smoking were identified using Spearman's rank correlation coefficient, as were
16 associations with time since cessation of smoking (years). Associations between tumour
17 location and *APC* promoter 1A hypermethylation were calculated using Fisher's exact
18 test. Statistical significance defined as $p < 0.05$.

19

20

1 **Results**

2 **Characteristics of the patients**

3 Details of the CRC patients from whom samples of colorectal tumours and adjacent
4 mucosa were obtained are provided in Table 1. Tumour tissue was obtained from 36
5 never smokers, 47 former smokers and 13 active smokers, while adjacent mucosa was
6 taken from 49 patients who were never smokers, 64 who were former smokers and 18
7 active smokers. Matched pairs of tumour and adjacent mucosa tissue were available for
8 89 of the patients (33 never smokers, 43 former smokers and 13 active smokers). The
9 mean level of smoking was 18.7 pack-years among active smokers and 12.7 pack-years
10 among former smokers. The mean duration of smoking was 37.6 years among active
11 smokers and 19.6 years among former smokers.

12

13 **The *APC* promoter 1A is hypermethylated in the tumours of active smokers**

14 Epigenome-wide analysis of DNA methylation in 96 colorectal tumours and 131 samples
15 of adjacent mucosa was performed using the Illumina Infinium HumanMethylation450
16 BeadChip microarray platform at the German Cancer Research Center Genomics and
17 Proteomics Core Facility (Heidelberg, Germany). An overview of performed analyses with
18 the different comparisons is shown in Figure 1.

19 We identified 21 CpG sites where methylation was significantly different between
20 tumours from patients who had never smoked and those who were active smokers at the
21 point of diagnosis (*Figure 2a*). These mapped to 14 loci within the *NFATC1* gene, 6 within
22 the *APC* gene, and 1 within *LAMB1* (*Table 2*). The 14 loci that mapped to the *NFATC1*
23 gene were distributed throughout the gene body and predominantly located in CpG
24 islands. In contrast, each of the six loci associated with *APC* corresponded to the 1A
25 promoter region and were within a span of 83 bp. Median beta values at each of the six

1 CpG sites were 0.41–0.53 higher in active smokers in comparison to never smokers. No
2 CpG sites were differentially methylated between tumours from former smokers and
3 never smokers.

4 5 **Smoking-specific differential methylation between tumours and adjacent mucosa**

6 We performed further analysis to identify genes that may be implicated in smoking-
7 associated carcinogenesis by identifying loci that are differentially methylated between
8 tumours and adjacent mucosa among active smokers but not never smokers. We
9 identified 148 loci that were significantly differentially methylated between these
10 conditions (*Figure 2b, Supplementary Table 1*). This included all six of the loci
11 previously identified within the *APC* 1A promoter and 9 of the 14 sites previously
12 identified within the *NFATC1* gene body. The nine sites with greatest statistical
13 significance all mapped to the *APC* 1A promoter, and a further six significantly
14 differentially methylated sites were also identified within this region. The average beta
15 values in tumours and adjacent mucosa from active smokers differed by >0.24 at each
16 of the 15 sites of the *APC* 1A promoter, while differing by <0.10 in the same tissues
17 from never smokers. Other genes prominently identified by this analysis included
18 receptor-type tyrosine-protein phosphatase N2 (*PTPRN2*) and sidekick cell adhesion
19 molecule 1 (*SDK1*).

20

21 ***APC* promoter 1A methylation and tumour pathology**

22 Our epigenome-wide analysis identified the *APC* promoter 1A as the leading target for
23 smoking-associated methylation changes. This was confirmed by cross-validation
24 analysis, which identified this region as the most predictive to distinguish between

1 tumours from never and active smokers (*Supplementary Figure 1*). We sought to further
2 characterise methylation of this region by tumour pathology and smoking behaviours.
3 Expanded analysis across the 15 significantly differentially methylated loci mapping to
4 the *APC* 1A promoter revealed distinct hypermethylation in some patients (*Figure 3a*).
5 Defining hypermethylation as mean beta values of >0.2 , in accordance with our
6 observed values across all tumours, the *APC* 1A promoter was hypermethylated in 7 of
7 36 tumours from never smokers (19%), 12 of 47 tumours from former smokers (26%),
8 and 8 of 13 tumours from active smokers (62%). Across all smoking behaviours,
9 hypermethylation was observed at all AJCC stages, including 4 of 8 stage I tumours
10 (*Figure 3B*), and was more common in tumours located in the rectum (14 of 38 tumours,
11 37%) and distal colon (8 of 25, 32%) than in the proximal colon (2 of 14, 14%), but not
12 significantly so (Fisher's exact test, $p=0.18$ and $p=0.28$ respectively). We identified no
13 associations between methylation at the six differentially methylated loci within the *APC*
14 1A promoter and alcohol consumption (grams/day) or BMI (both $p > 0.05$).
15 Hypermethylation of the 1A promoter was significantly more frequent among women
16 (Fisher's exact test, $p=0.02$) and was associated with younger age (Spearman rank
17 correlation, $\rho=-0.28$, $p=0.01$).

18

19 **Methylation of the *APC* 1A promoter is associated with duration of smoking**

20 To explore the relation between the intensity and duration of smoking with methylation
21 of the *APC* 1A promoter, we utilised data for the 72 former and active smokers in this
22 study regarding intensity (pack-years) and smoking duration (length of time for which
23 the patient smoked). Additionally, for the 47 former smokers, the relation with the length
24 of time between cessation of smoking and cancer diagnosis was also assessed.

1 Greater duration of smoking was significantly and positively associated with increased
2 methylation at cg14479889 ($\rho=0.27$, $p=0.03$) and trended towards significance at each
3 of the other five differentially methylated loci ($\rho>0.19$, $p<0.09$) (*Table 3*). Most notably,
4 the average methylation (beta values) across the 15 differentially methylated loci
5 mapped to this promoter region was significantly and positively associated with duration
6 of smoking ($\rho=0.26$, $p=0.03$). No significant associations were observed with pack-years
7 of smoking ($p>0.29$) or time between cessation of smoking and cancer diagnosis among
8 former smokers ($p>0.16$).

9

10 **The *APC* promoter 1A is not hypermethylated in the mucosa adjacent to tumours**

11 We examined *APC* promoter 1A methylation in mucosa adjacent to tumours, to
12 determine whether hypermethylation of this region exists as a field defect. Matched
13 tumour and adjacent mucosal tissue were available for 24 of the 27 patients with
14 tumoural hypermethylation of the 1A promoter (irrespective of smoking status). No
15 promoter hypermethylation was observed in the adjacent mucosa from any of the 24
16 patients (*Figure 3c*, average beta < 0.11), or individually at any of the six differentially
17 methylated loci ($\beta<0.13$) (*Supplementary Figure 2*).

18

19 ***TNXB* is differentially methylated in the adjacent mucosa of smokers**

20 To gain insight into how smoking may act upon the colon, such as through carcinogenic
21 compounds from cigarette smoke carried in the blood or chronic inflammation, we
22 performed epigenome-wide analyses of DNA methylation in adjacent mucosa by
23 smoking behaviours. We identified four sites within a 500 bp region that map to the

1 tenascin XB (*TNXB*) gene body that were significantly hypomethylated in mucosa from
2 active smokers (*Supplementary Table 2*). No differentially methylated loci were
3 observed between former and never smokers.

4

1 **Discussion**

2 In this study, we investigated how epigenetic factors may be implicated in conferring the
3 increased risk of colorectal cancer among smokers by performing epigenome-wide
4 analysis of DNA methylation in samples of tumours and adjacent mucosa by smoking
5 behaviours. We report that smoking at the time of diagnosis is significantly associated
6 with hypermethylation of the 1A promoter of *APC*, a key tumour suppressor gene that
7 has been extensively studied with regard to colorectal cancer. Hypermethylation was
8 unique to tumour tissue and was associated with the duration for which the patient has
9 smoked. We observed that hypermethylation of this promoter was more common in the
10 rectum and distal colon, in concordance with evidence that the association between
11 smoking and CRC risk is greatest for developing tumours in the rectum [1,26]. Our
12 findings may implicate the epigenetic silencing of *APC* in smoking-associated colorectal
13 carcinogenesis. However, due to the relatively small number of patients who were
14 active smokers at diagnosis, our results should be considered exploratory at this stage.
15 We have been unable to validate our observations in an independent cohort due to the
16 absence of publicly-available datasets incorporating smoking history, and insufficient
17 numbers of active smokers at diagnosis within other studies. Further work in external
18 cohorts is required to examine the validity of our observations.

19 *APC* is a tumour suppressor gene and regulator of the Wnt signalling pathway,
20 which acts via regulation of β -catenin degradation and localisation. Loss of *APC*
21 function has been proposed as a key early event in the development of sporadic
22 colorectal cancer [27], with inactivation frequently occurring through mutations,
23 especially in the mutation cluster region [28], and promoter methylation [29]. Expression
24 of the 1A mRNA isoform of *APC* is regulated in part through methylation of promoter 1A
25 (chr5:112,072,710-112,073,585) [30], and this region is aberrantly methylated in

1 colorectal, breast and lung tumours, resulting in transcriptional silencing and increased
2 activation of the Wnt signalling pathway [29,31]. We observed significantly greater
3 methylation of this region in tumours from patients who were active smokers at the point
4 of diagnosis, thereby linking smoking behaviours to silencing of this key tumour
5 suppressor gene. It has been reported elsewhere that smoking is associated with
6 mutations in *TP53* and *BRAF* but not *APC* [32], which together with our study may
7 suggest that inactivation of this gene more commonly occurs through epigenetic
8 dysregulation in smoking-associated CRC than through genetic changes. Median
9 promoter methylation levels (beta values) were approximately 0.5 higher in active
10 smokers (*Figure 3*), consistent with monoallelic methylation of the promoter. Although
11 evidence from the mouse model suggests that inactivation of both alleles is required for
12 tumourigenesis [33], monoallelic methylation of the *APC* promoter 1A is a frequent
13 event in human colorectal tumours [31,34] and cancer cell lines [35], and has been
14 reported in gastric tumours [36].

15 Interestingly, hypermethylation of the *APC* promoter 1A was never present in
16 mucosa adjacent to the tumours. Methylation at each interrogated CpG site within
17 promoter 1A was very highly conserved in adjacent mucosa, while in direct contrast
18 there was substantial variation in promoter methylation between tumour samples
19 (*Figure 3C, Supplementary Figure 2*). Cancer is associated with significantly greater
20 variability in DNA methylation than is found in healthy tissue, and this loss of stability
21 and increased stochastic variation may facilitate malignant cells to adapt to changes in
22 their microenvironments [37,38]. Genetic and epigenetic alterations implicated in
23 carcinogenesis are sometimes present in the surrounding tissue as field defects [39,40],
24 and increased variation in DNA methylation has been observed in cytologically-normal
25 cells from individuals later diagnosed with cervical cancer [41]. However, we observed

1 that methylation of the *APC* promoter 1A was still highly conserved in adjacent mucosa,
2 in line with studies reporting an absence of *APC* hypermethylation in colonic mucosa
3 [31]. Mutations in *APC* are sufficient to induce polyp formation in mice [42,43] and
4 humans [44], and we therefore speculate that this absence of *APC* hypermethylation in
5 adjacent mucosa may be due to the key role for loss of *APC* in driving carcinogenesis.
6 Indeed, we observed hypermethylation of the 1A promoter in half of stage I tumours
7 (*Figure 3B*). This hypothesis is further supported by evidence of *APC* promoter
8 methylation being an early event in colorectal carcinogenesis that is detectable in small
9 (<15 mm) adenomas [31].

10 We observed a significant association between promoter 1A hypermethylation
11 and duration of smoking (*Table 3*), but further work is required to expand upon the
12 relation between the intensity and duration of exposure and epigenetic events in CRC.
13 Indeed, promoter hypermethylation was not observed in tumours from any of the 8
14 former smokers who had smoked for >35 years, and our epigenome-wide analysis did
15 not identify any differentially methylated sites between never smokers and former
16 smokers in tumours or adjacent mucosa. The cessation of smoking is known to reduce
17 the risk of CIMP-high colorectal cancer and patients who quit >10 years prior to
18 diagnosis display similar risk of CIMP-high tumours to never smokers [18]. Furthermore,
19 it is known that methylation of the *AHRR* and *F2RL3* genes returns to normal levels with
20 increasing time since cessation [9]. Therefore, as only 9 of the 40 former smokers in this
21 study for whom there is relevant data ceased smoking <10 years prior to diagnosis, we
22 speculate that the time since cessation may also be a significant factor in the risk of
23 hypermethylation of *APC* promoter 1A.

24 Our epigenome-wide analysis also identified the *NFATC1* gene body as being
25 hypomethylated in tumours from smokers (*Figure 2*). This gene encodes a transcription

1 factor implicated in T cell activation. Epigenetic dysregulation of this gene has been
2 observed in hepatocellular carcinoma [45] and lymphomas [46] while hypomethylation
3 has been reported in healthy individuals with lower socioeconomic status [47]. Our
4 study is the first to report hypomethylation of *NFATC1* in colorectal tumours.
5 Overexpression of the gene is associated with worse prognosis in stage II and III
6 colorectal cancer patients, which may occur through the promotion of cell migration and
7 metastasis [48,49]. As the ColoCare Study began to recruit patients in October 2010,
8 we are currently unable to determine whether *NFATC1* methylation is associated with
9 patient prognosis in this cohort. We will be able to address this question in time as
10 further data regarding patient outcomes is collected.

11 Our data suggests that smoking is not associated with the accumulation of
12 widespread epigenetic defects in the adjacent mucosa. Methylation of the *APC*
13 promoter 1A occurs independently of other epigenetic events in CRC [31], and we
14 identified only one gene, *TNXB*, as differentially methylated in the adjacent mucosa of
15 active smokers (*Table 3*). This may be considered to be in contrast to the findings of
16 Paun *et al* [50], who reported disruption of normal gene methylation profiles in the
17 normal rectal mucosa of smokers. We speculate that this may be the product of our
18 analyses identifying genes implicated in malignant transformation due to our
19 comparison of tumours and adjacent mucosa, while Paun *et al* examined rectal mucosa
20 prior to the advent of tumour formation. To our knowledge, ours is the first study to
21 observe differential methylation of *TNXB* by smoking behaviours. Further work is
22 required to investigate how this extracellular matrix glycoprotein could be implicated in
23 smoking-associated carcinogenesis.

24 Further to the inability to confirm our findings in an independent cohort, the
25 comparatively low number of patients who actively smoked at the point of diagnosis is a

1 limitation of this study, and one which could inhibit the identification of associations
2 between smoking and methylation. We therefore incorporated the chromosomal position
3 into test statistics by means of a NHMM, which also served to reduce the probability of
4 secluded differentially methylated CpGs and hence most likely false positives. A particular
5 strength of this study is the analysis of both tumour tissue and adjacent mucosa, which
6 has enabled us to gain greater insight by identifying epigenetic events associated with
7 smoking that are uniquely found in tumour tissue (hypermethylation of the *APC* promoter
8 1A) and to establish an absence of field defects associated with smoking in the
9 neighbouring mucosa.

10 In conclusion, we report exploratory evidence for hypermethylation of the *APC*
11 promoter 1A being implicated in the development of colorectal tumours among smokers.
12 Methylation of this region was significantly associated with smoking at the point of
13 diagnosis and with the duration of time for which the patient smoked, and
14 hypermethylation was confined to tumours. Further work is required to validate our
15 observations in independent cohorts, and to identify implications for patient prognosis.

16

1 **Acknowledgements**

2 The authors would like to thank all ColoCare study participants and the entire ColoCare
3 study team in Heidelberg, especially Dr Werner Diehl for data acquisition and
4 documentation, and Judith Kammer, Susanne Jakob and Torsten Koelsch for patient
5 recruitment and tissue collection. We are grateful to Dr. Melanie Bewerunge-Hudler and
6 the Genomics and Proteomics Core Facility at the German Cancer Research Center
7 (Heidelberg, Germany) for the performance of the Illumina Infinium
8 HumanMethylation450 BeadChip microarrays. The ColoCare Study and Consortium
9 has been designed and first implemented at the Fred Hutchinson Cancer Research
10 Center, Seattle, USA (PIs: Ulrich/Grady) and protocols have been used with permission
11 in Heidelberg, Germany (PI: Ulrich). The ColoCare Study site in Heidelberg has been
12 funded by the Matthias Lackas Foundation, the German Consortium for Translational
13 Cancer Research (DKTK) and the Division of Preventive Oncology at the German
14 Cancer Research Center. Hagen Klett and Melanie Boerries were additionally funded by
15 the German Ministry of Education and Research (BMBF) within the e:Med consortium
16 “DeCaRe-Delineating Cardiac Regeneration”.

17

18 **Statement of author contributions**

19 CU conceived the cohort study. TB, RT, NH, CU and KM conceived the investigation
20 into smoking. JB, LZ, MS, AU, PS and EH organised and performed the sample
21 collection. RT, BG, DS, SS, CAM, PSK and HBrenner were involved in data collection
22 and organisation. TB, HK, RT, HBusch and MB analysed the DNA methylation data. TB,
23 RT, CU and KM performed data interpretation. TB wrote the manuscript, with figures
24 generated by TB and HK. All authors were involved in writing and had final approval of
25 the submitted manuscript.

1

2 **References**

- 3 1. Botteri E, Iodice S, Bagnardi V, *et al.* 2008. Smoking and colorectal cancer: A meta-
4 analysis. *JAMA* **300**(23): 2765-2778.
- 5 2. Parajuli R, Bjerkaas E, Tverdal A, *et al.* 2014. Cigarette smoking and colorectal
6 cancer mortality among 602,242 Norwegian males and females. *Clin Epidemiol* **6**: 137-
7 145.
- 8 3. Walter V, Jansen L, Hoffmeister M, *et al.* 2014. Smoking and survival of colorectal
9 cancer patients: Systematic review and meta-analysis. *Ann Oncol* **25**(8): 1517-1525.
- 10 4. Giovannucci E & Martínez ME 1996. Tobacco, colorectal cancer, and adenomas: A
11 review of the evidence. *J Natl Cancer Inst* **88**(23): 1717-1730.
- 12 5. Novakovic B, Ryan J, Pereira N, *et al.* 2014. Postnatal stability, tissue, and time
13 specific effects of AHRR methylation change in response to maternal smoking in
14 pregnancy. *Epigenetics* **9**(3): 377-386.
- 15 6. Elliott HR, Tillin T, McArdle WL, *et al.* 2014. Differences in smoking associated DNA
16 methylation patterns in South Asians and Europeans. *Clin Epigenetics* **6**(1): 4.
- 17 7. Leng S, Liu Y, Thomas CL, *et al.* 2013. Native American ancestry affects the risk for
18 gene methylation in the lungs of Hispanic smokers from New Mexico. *Am J Respir Crit*
19 *Care Med* **188**(9): 1110-1116.
- 20 8. Zhang Y, Yang R, Burwinkel B, *et al.* 2014. F2RL3 methylation as a biomarker of
21 current and lifetime smoking exposures. *Environ Health Perspect* **122**(2): 131-137.
- 22 9. Shenker NS, Polidoro S, van Veldhoven K, *et al.* 2013. Epigenome-wide association
23 study in the European Prospective Investigation into Cancer and Nutrition (EPIC-Turin)
24 identifies novel genetic loci associated with smoking. *Hum Mol Genet* **22**(5): 843-851.

- 1 10. Breitling LP, Yang R, Korn B, *et al.* 2011. Tobacco-smoking-related differential DNA
2 methylation: 27K discovery and replication. *Am J Hum Genet* **88**(4): 450-457.
- 3 11. Andreotti G, Karami S, Pfeiffer RM, *et al.* 2014. LINE1 methylation levels associated
4 with increased bladder cancer risk in pre-diagnostic blood DNA among US (PLCO) and
5 European (ATBC) cohort study participants. *Epigenetics* **9**(3): 404-415.
- 6 12. Zhang Y, Yang R, Burwinkel B, *et al.* 2014. F2RL3 methylation in blood DNA is a
7 strong predictor of mortality. *Int J Epidemiol* **43**(4): 1215-1225.
- 8 13. Tan Q, Wang G, Huang J, *et al.* 2013. Epigenomic analysis of lung adenocarcinoma
9 reveals novel DNA methylation patterns associated with smoking. *Onco Targets Ther* **6**:
10 1471-1479.
- 11 14. Marsit CJ, Karagas MR, Schned A, *et al.* 2006. Carcinogen exposure and epigenetic
12 silencing in bladder cancer. *Ann N Y Acad Sci* **1076**: 810-821.
- 13 15. Wolff EM, Liang G, Cortez CC, *et al.* 2008. RUNX3 methylation reveals that bladder
14 tumors are older in patients with a history of smoking. *Cancer Res* **68**(15): 6208-6214.
- 15 16. Liu Y, Lan Q, Siegfried JM, *et al.* 2006. Aberrant promoter methylation of p16 and
16 MGMT genes in lung tumors from smoking and never-smoking lung cancer patients.
17 *Neoplasia* **8**(1): 46-51.
- 18 17. Lokk K, Vooder T, Kolde R, *et al.* 2012. Methylation markers of early-stage non-
19 small cell lung cancer. *PLoS One* **7**(6): e39813.
- 20 18. Nishihara R, Morikawa T, Kuchiba A, *et al.* 2013. A prospective study of duration of
21 smoking cessation and colorectal cancer risk by epigenetics-related tumor classification.
22 *Am J Epidemiol* **178**(1): 84-100.
- 23 19. Triche TJ, Weisenberger DJ, Van Den Berg D, *et al.* 2013. Low-level processing of
24 illumina infinium DNA methylation beadarrays. *Nucleic Acids Res* **41**(7): e90.

- 1 20. Fortin JP, Labbe A, Lemire M, *et al.* 2014. Functional normalization of 450k
2 methylation array data improves replication in large cancer studies. *Genome Biol*
3 **15**(12):503.
- 4 21. Smyth GK 2004. Linear models and empirical bayes methods for assessing
5 differential expression in microarray experiments. *Stat Appl Genet Mol Biol* **3**: Article3.
- 6 22. Kuan PF & Chiang DY 2012. Integrating prior knowledge in multiple testing under
7 dependence with applications to detecting differential DNA methylation. *Biometrics*
8 **68**(3): 774-783.
- 9 23. Sun W & Tony Cai T 2009. Large-scale multiple testing under dependence. *Journal*
10 *of the Royal Statistical Society: Series B (Statistical Methodology)* **71**(2): 393-424.
- 11 24. Kent WJ, Sugnet CW, Furey TS, *et al.* 2002. The human genome browser at UCSC.
12 *Genome Res* **12**(6): 996-1006.
- 13 25. Sherry ST, Ward M & Sirotkin K 1999. DbSNP-database for single nucleotide
14 polymorphisms and other classes of minor genetic variation. *Genome Res* **9**(8): 677-
15 679.
- 16 26. Cheng J, Chen Y, Wang X, *et al.* 2015. Meta-analysis of prospective cohort studies
17 of cigarette smoking and the incidence of colon and rectal cancers. *Eur J Cancer Prev*
18 **24**(1): 6-15.
- 19 27. Fearon ER & Vogelstein B 1990. A genetic model for colorectal tumorigenesis. *Cell*
20 **61**(5): 759-767.
- 21 28. Miyoshi Y, Nagase H, Ando H, *et al.* 1992. Somatic mutations of the APC gene in
22 colorectal tumors: Mutation cluster region in the APC gene. *Hum Mol Genet* **1**(4): 229-
23 233.

- 1 29. Hiltunen MO, Alhonen L, Koistinaho J, *et al.* 1997. Hypermethylation of the APC
2 (adenomatous polyposis coli) gene promoter region in human colorectal carcinoma. *Int*
3 *J Cancer* **70**(6): 644-648.
- 4 30. Segditsas S, Sieber OM, Rowan A, *et al.* 2008. Promoter hypermethylation leads to
5 decreased APC mRNA expression in familial polyposis and sporadic colorectal tumours,
6 but does not substitute for truncating mutations. *Exp Mol Pathol* **85**(3): 201-206.
- 7 31. Esteller M, Sparks A, Toyota M, *et al.* 2000. Analysis of adenomatous polyposis coli
8 promoter hypermethylation in human cancer. *Cancer Res* **60**(16): 4366-4371.
- 9 32. Chen K, Xia G, Zhang C, *et al.* 2015. Correlation between smoking history and
10 molecular pathways in sporadic colorectal cancer: A meta-analysis. *Int J Clin Exp Med*
11 **8**(3): 3241-3257.
- 12 33. Samuel MS, Suzuki H, Buchert M, *et al.* 2009. Elevated Dnmt3a activity promotes
13 polyposis in *Apc*(min) mice by relaxing extracellular restraints on Wnt signaling.
14 *Gastroenterology* **137**(3): 902-13, 913.e1-11.
- 15 34. Arnold CN, Goel A, Niedzwiecki D, *et al.* 2004. APC promoter hypermethylation
16 contributes to the loss of APC expression in colorectal cancers with allelic loss on 5q.
17 *Cancer Biol Ther* **3**(10): 960-964.
- 18 35. Lind GE, Thorstensen L, Løvig T, *et al.* 2004. A CpG island hypermethylation profile
19 of primary colorectal carcinomas and colon cancer cell lines. *Mol Cancer* **3**: 28.
- 20 36. Clément G, Bosman FT, Fontollet C, *et al.* 2004. Monoallelic methylation of the
21 APC promoter is altered in normal gastric mucosa associated with neoplastic lesions.
22 *Cancer Res* **64**(19): 6867-6873.
- 23 37. Barrow TM, Barault L, Ellsworth RE, *et al.* 2015. Aberrant methylation of imprinted
24 genes is associated with negative hormone receptor status in invasive breast cancer. *Int*
25 *J Cancer* **137**(3): 537-547.

- 1 38. Hansen KD, Timp W, Bravo HC, *et al.* 2011. Increased methylation variation in
2 epigenetic domains across cancer types. *Nat Genet* **43**(8): 768-775.
- 3 39. Ellsworth DL, Ellsworth RE, Love B, *et al.* 2004. Genomic patterns of allelic
4 imbalance in disease free tissue adjacent to primary breast carcinomas. *Breast Cancer*
5 *Res Treat* **88**(2): 131-139.
- 6 40. Shen L, Kondo Y, Rosner GL, *et al.* 2005. MGMT promoter methylation and field
7 defect in sporadic colorectal cancer. *J Natl Cancer Inst* **97**(18): 1330-1338.
- 8 41. Teschendorff AE, Jones A, Fiegl H, *et al.* 2012. Epigenetic variability in cells of
9 normal cytology is associated with the risk of future morphological transformation.
10 *Genome Med* **4**(3): 24.
- 11 42. Oshima M, Oshima H, Kitagawa K, *et al.* 1995. Loss of Apc heterozygosity and
12 abnormal tissue building in nascent intestinal polyps in mice carrying a truncated Apc
13 gene. *Proc Natl Acad Sci U S A* **92**(10): 4482-4486.
- 14 43. Shibata H, Toyama K, Shioya H, *et al.* 1997. Rapid colorectal adenoma formation
15 initiated by conditional targeting of the Apc gene. *Science* **278**(5335): 120-123.
- 16 44. Lamlum H, Papadopoulou A, Ilyas M, *et al.* 2000. APC mutations are sufficient for
17 the growth of early colorectal adenomas. *Proc Natl Acad Sci U S A* **97**(5): 2225-2228.
- 18 45. Song MA, Tiirikainen M, Kwee S, *et al.* 2013. Elucidating the landscape of aberrant
19 DNA methylation in hepatocellular carcinoma. *PLoS One* **8**(2): e55761.
- 20 46. Akimzhanov A, Krenacs L, Schlegel T, *et al.* 2008. Epigenetic changes and
21 suppression of the nuclear factor of activated T cell 1 (NFATC1) promoter in human
22 lymphomas with defects in immunoreceptor signaling. *Am J Pathol* **172**(1): 215-224.
- 23 47. Stringhini S, Polidoro S, Sacerdote C, *et al.* 2015. Life-course socioeconomic status
24 and DNA methylation of genes regulating inflammation. *Int J Epidemiol* **44**(4): 1320-
25 1330.

- 1 48. Jauliac S, López-Rodríguez C, Shaw LM, *et al.* 2002. The role of NFAT transcription
2 factors in integrin-mediated carcinoma invasion. *Nat Cell Biol* **4**(7): 540-544.
- 3 49. Tripathi MK, Deane NG, Zhu J, *et al.* 2014. Nuclear factor of activated t-cell activity
4 is associated with metastatic capacity in colon cancer. *Cancer Res* **74**(23): 6947-6957.
- 5 50. Paun BC, Kukuruga D, Jin Z, *et al.* 2010. Relation between normal rectal
6 methylation, smoking status, and the presence or absence of colorectal adenomas.
7 *Cancer* **116**(19): 4495-4501.

Table 1: Clinical and demographic characteristics of the patients

		Adjacent mucosa			Tumour		
		Never	Former	Active	Never	Former	Active
Patients	<i>n</i>	49	64	18	36	47	13
Age	Mean	63.7	65.5	56.6	63.9	65.1	59.6
	SD	11.8	10.9	12.5	11.5	10.3	9.8
	Range	34 - 82	38 - 89	22 - 79	41 - 82	38 - 89	35 - 79
Gender	Male	24	48	10	21	37	6
	Female	25	16	8	15	10	7
Stage	I	8	8	4	3	5	1
	II	15	24	7	12	17	7
	III	13	18	5	12	14	4
	IV	11	14	2	9	11	1
Pack-years	Mean (years)	-	11.5	18.7	-	12.6	16.4
	0 – 9 (n)	-	31	5	-	22	4
	10 – 19 (n)	-	15	4	-	11	3
	≥20 (n)	-	11	6	-	11	4
Duration	Mean (years)	-	18.8	31.8	-	19.6	37.6
	0 – 9 (n)	-	18	2	-	11	0
	10 – 19 (n)	-	17	3	-	11	2
	20 – 29 (n)	-	13	2	-	8	0
	≥30 (n)	-	15	17	-	11	9

Table 2: CpG sites with differential methylation by smoking status in tumours

Probe ID	Chromosomal location	Gene	Gene region	Island status	Mean β -value		LIS p value
					Never	Active	
cg08571859	chr5:112073350	APC	TSS1500	Open sea	0.11	0.36	7.4×10^{-6}
cg14511739	chr5:112073373	APC	TSS200	Open sea	0.11	0.39	1.2×10^{-6}
cg22035501	chr5:112073426	APC	TSS200	Open sea	0.12	0.42	4.4×10^{-7}
cg11613015	chr5:112073433	APC	TSS200	Open sea	0.10	0.34	9.0×10^{-7}
cg14479889	chr5:112073426	APC	TSS200	Open sea	0.12	0.38	1.6×10^{-6}
cg16970232	chr5:112073433	APC	TSS200	Open sea	0.13	0.40	2.4×10^{-6}
cg04744624	chr7: 107641770	LAMB1	Body	N_Shore	0.23	0.41	8.8×10^{-6}
cg15138382	chr18: 77186504	NFATC1	Body / 5' UTR	Island	0.89	0.75	1.8×10^{-6}
cg05302701	chr18: 77196320	NFATC1	Body / 5' UTR	Island	0.81	0.68	4.4×10^{-6}
cg18092363	chr18: 77202678	NFATC1	Body / 5' UTR	Island	0.94	0.85	7.9×10^{-6}
cg26550337	chr18: 77203542	NFATC1	Body / 5' UTR	Island	0.81	0.70	1.5×10^{-6}
cg26100137	chr18: 77203667	NFATC1	Body / 5' UTR	Island	0.97	0.90	2.3×10^{-6}
cg22279865	chr18: 77204561	NFATC1	Body / 5' UTR	S_Shore	0.93	0.87	4.5×10^{-6}
cg00445548	chr18: 77207209	NFATC1	Body / 5' UTR	Island	0.93	0.82	7.1×10^{-6}
cg02675550	chr18: 77208807	NFATC1	Body / 5' UTR	Island	0.86	0.75	5.6×10^{-6}
cg21242663	chr18: 77208881	NFATC1	Body	Island	0.90	0.81	8.1×10^{-6}
cg25595641	chr18: 77208991	NFATC1	Body	Island	0.94	0.86	1.8×10^{-6}
cg21806238	chr18: 77210990	NFATC1	Body	Island	0.92	0.84	5.1×10^{-6}
cg16253249	chr18: 77211212	NFATC1	Body	Island	0.81	0.74	1.8×10^{-6}
cg03239925	chr18: 77230795	NFATC1	Body	Island	0.74	0.63	7.0×10^{-6}
cg22324981	chr18: 77283493	NFATC1	Body	N_Shore	0.80	0.58	9.2×10^{-7}

Table 3: Associations between DNA methylation and smoking intensity and duration in tumours

Probe ID	Chromosomal location	Gene	Pack-years		Duration		Time since cessation	
			ρ	p	ρ	p	ρ	p
cg08571859	chr5:112073350	APC	-0.07	0.31	0.22	0.06	-0.05	0.41
cg14511739	chr5:112073373	APC	-0.01	0.47	0.19	0.09	-0.05	0.41
cg22035501	chr5:112073426	APC	-0.08	0.29	0.19	0.09	-0.11	0.31
cg11613015	chr5:112073433	APC	-0.02	0.44	0.20	0.08	-0.14	0.27
cg14479889	chr5:112073426	APC	-0.05	0.37	0.27	0.03	-0.21	0.17
cg16970232	chr5:112073433	APC	-0.01	0.48	0.21	0.07	-0.10	0.32
Promoter 1A	chr5:112,072,710 - 112,073,585	APC	0.02	0.44	0.26	0.03	-0.20	0.19

Figure legends

Table 1: Clinical and demographic characteristics of the patients. Data are provided regarding the age (mean, standard deviation, and range), gender, tumour stage and pack-years of smoking for the patients according to smoking status at the point of cancer diagnosis.

Table 2: CpG sites with differential methylation by smoking status in tumours. Loci with significantly different methylation between tumours from never smokers and active smokers are listed, including Illumina annotation data. Median beta values are provided, along with pLIS values.

Table 3: Associations between DNA methylation and smoking intensity and duration in tumours. Spearman's rank correlation coefficients were calculated for each of the significantly different loci in tumour tissue, using data from former (n=47) and active (n=13) smokers. Correlations were calculated between methylation (beta values) and the pack-years of smoking or duration (years) of smoking. Additionally, for former smokers, correlations between methylation and time since cessation were calculated. ρ and p values are provided, with significant values highlighted in bold.

Figure 1: Overview of analyses by smoking behaviours in tumours and adjacent mucosa. Differentially methylated sites between smokers and never smokers were identified in tumour tissue and in adjacent mucosa. Further analyses were performed to identify sites displaying smoking-specific differential methylation between tumours and adjacent mucosa.

Figure 2: Manhattan plots showing differentially methylated sites between never and active smokers. Results of the analyses between tumours from never and active smokers (A) and differential methylation between tumours and adjacent mucosa unique to active smokers (B). Genesymbols of the genes associated with the most significantly different sites are provided. The threshold (line) represents statistical significance ($p_{LIS} < 1 \times 10^{-5}$)

Figure 3: Methylation of the *APC* promoter 1A in tumours and matched adjacent mucosa. Mean methylation levels (beta values) for each patient were calculated across the 15 CpG sites mapping to the 1A promoter that were identified as differentially methylated by smoking status (*Figure 2*). A: promoter methylation in tumours by patient smoking status. Mean values by smoking status are indicated by horizontal lines. B: promoter methylation in tumours by AJCC stage in all patients. Mean values by stage are indicated by horizontal lines. C: promoter methylation in matched samples of tumours and adjacent mucosa from 89 patients (33 never smokers, 43 former smokers, and 13 active smokers). Lines indicate matched samples from the same patient.

Supplementary Table 1: CpG sites with smoking-specific differential methylation between tumours and adjacent mucosa. Loci that are differentially methylated between tumour and adjacent mucosa tissue among active smokers, but not never smokers, are listed. Illumina annotation data, median beta values in adjacent mucosa

and tumour tissue for never and active smokers, and pLIS values are provided.

Statistical significance was defined as $pLIS < 1 \times 10^{-5}$.

Supplementary Table 2: CpG sites with differential methylation by smoking status in adjacent mucosa. Loci with significantly different methylation between adjacent mucosa from never smokers and active smokers are listed, including Illumina annotation data. Median beta values are provided, along with pLIS values.

Supplementary Figure 1: Prediction performances of the top-six methylation features. Cross-validation analysis with support vector machine was used to identify regions predictive of smoking behaviour in tumours. The area under the curve of the receiver operating characteristic are presented for the top-six predictive features.

Supplementary Figure 2: Methylation of differentiated methylated sites within the *APC* promoter 1A in matched tumours and adjacent mucosa. Methylation (beta values) at each of the 15 significantly differentially methylated loci identified in tumours (*Figure 2*) in match tumour and adjacent mucosa samples from 89 patients. Lines indicate matched samples from the same patient.

Supplementary Materials and Methods: Cross-validation analysis

Figure 1: Overview of analyses by smoking behaviours in tumours and adjacent mucosa

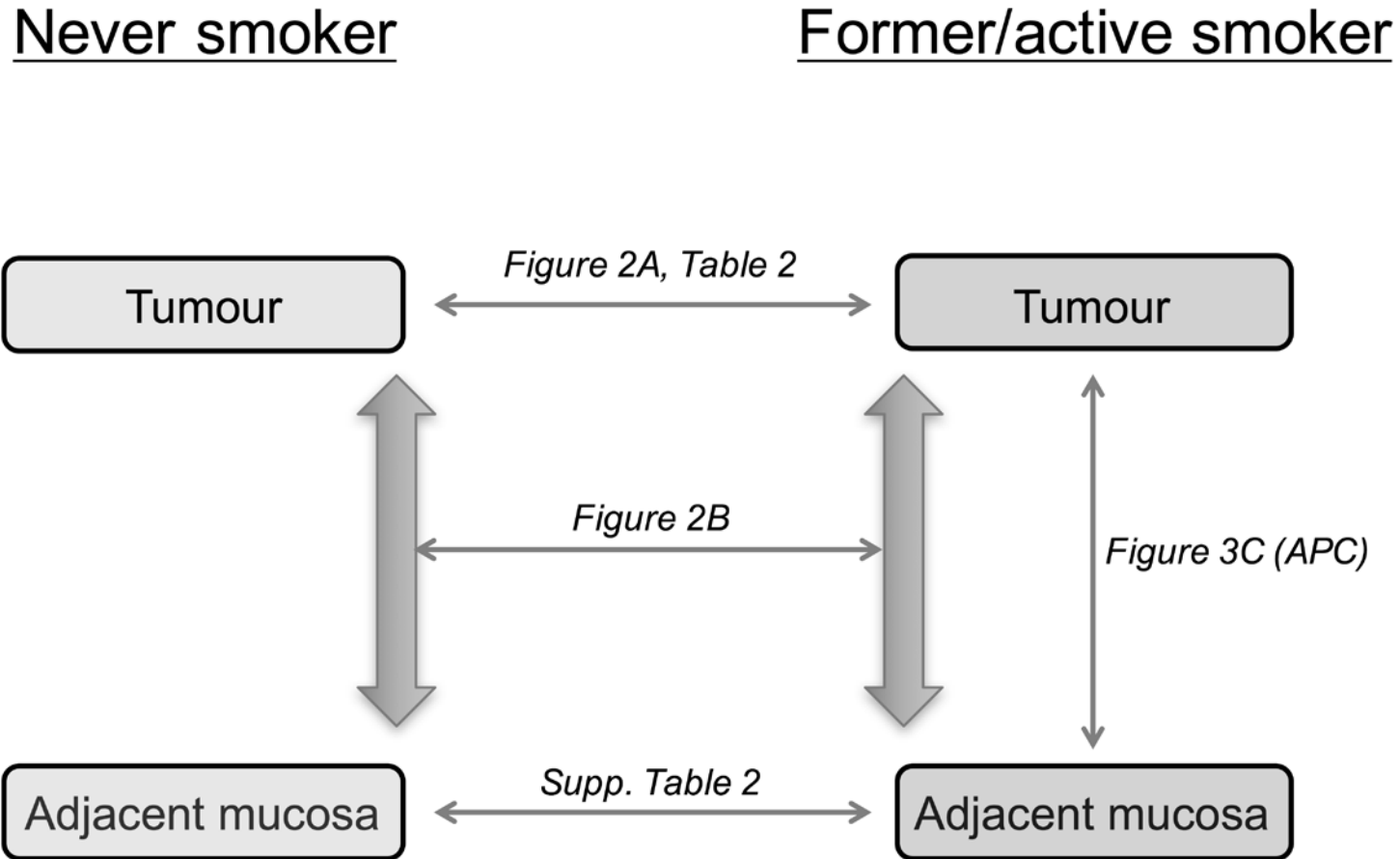


Figure 2: Manhattan plots showing differentially methylated sites between never and active smokers

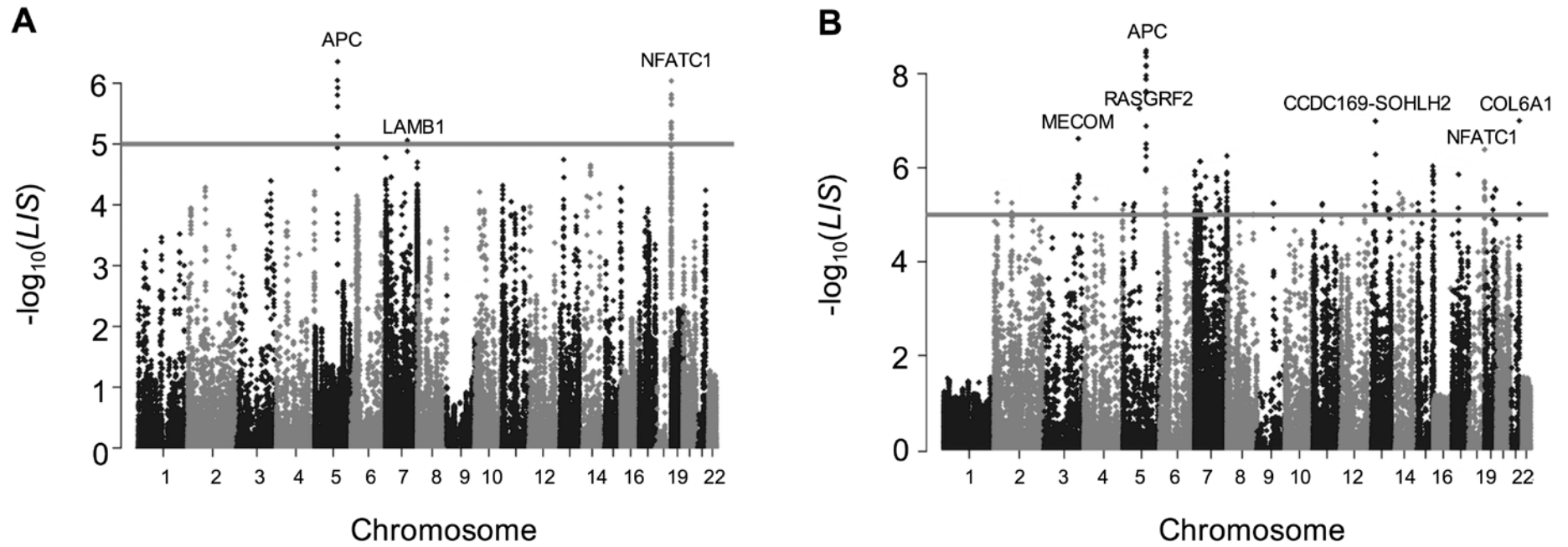
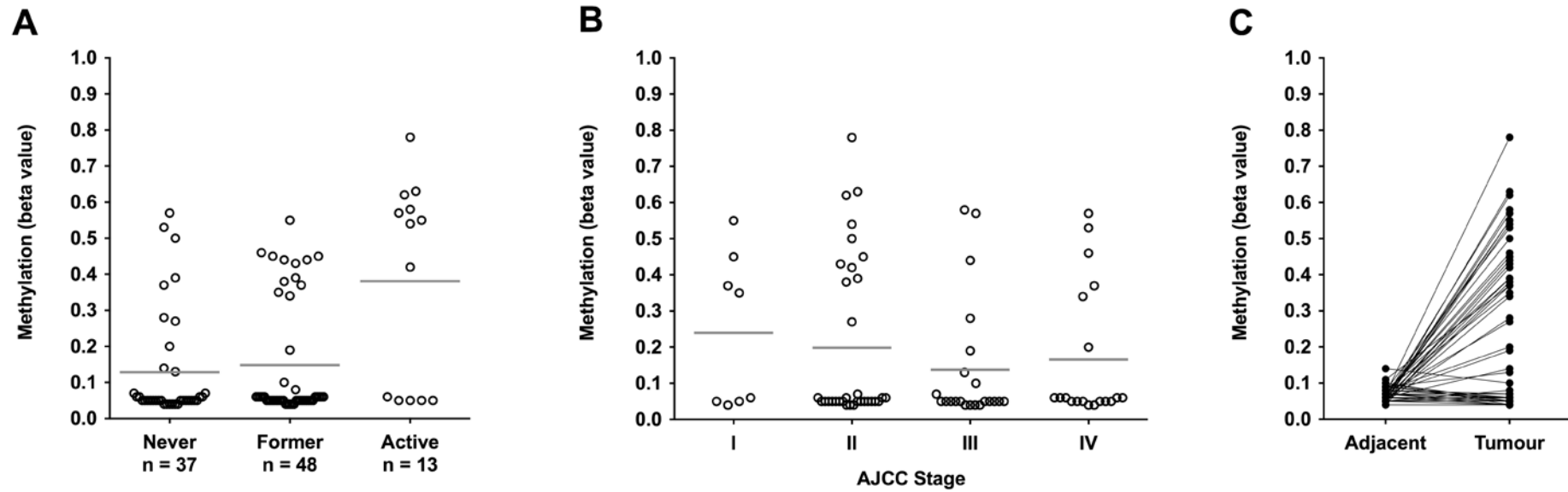


Figure 3: Methylation of the *APC* promoter 1A in tumours and matched adjacent mucosa



Supplementary Materials and Methods

Cross-validation analysis

To confirm robustness of our findings from the NHMM, we used cross-validation analysis with support vector machine to find smoking associated methylation regions according to their predictive power. Average beta values were calculated at promoter (TSS1500 – 1stExon) and gene body for different island status (Island, Shore, Shelf, OpenSea) for each gene according to the Illumina 450k annotation. This combined to 78,405 features of which we considered the top 2% with the highest standard deviation across tumor samples for further analysis (n=1,569). These features were individually taken to predict smoking status (never / active) with machine learning. Prediction performances, as the area under the curve (AUC) of the receiver operating characteristic (ROC), were calculated by three times repeated 10-fold cross-validation using a linear support vector machine kernel (cost=1). To account for imbalanced sample groups (never smokers = 36; active smokers = 13), we averaged results of 24 down-samplings of the never-smoking group, i.e., 13 of 36 randomly selected never-smoking samples were considered for cross-validation analysis. All analysis was performed with R caret package.

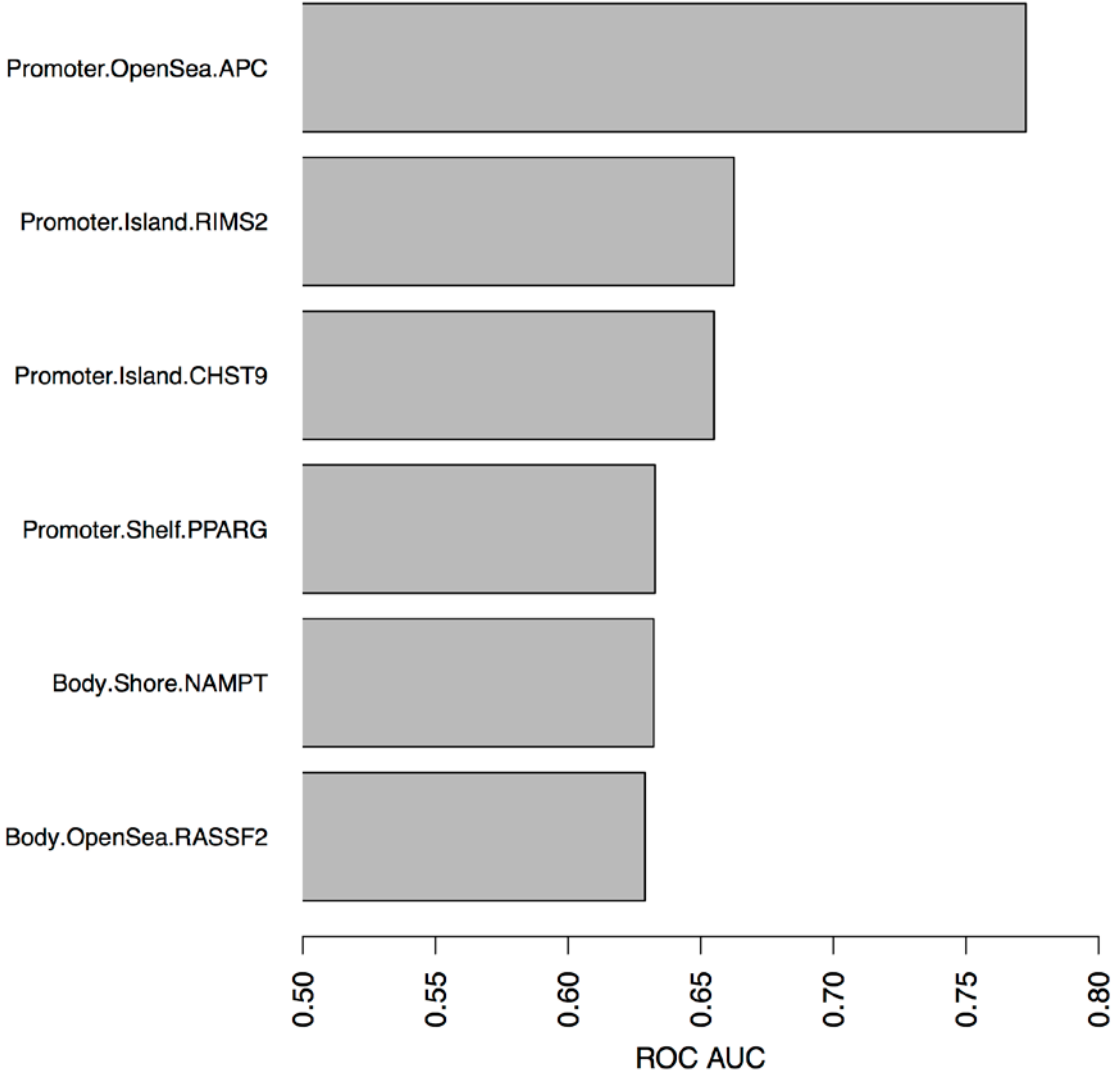
We identified the *APC* promoter OpenSea (30 CpGs) to be the most predictive region to differentiate active from never smokers in tumor tissues (AUC = 0.77). This supported the results from the NHMM and confirms the robustness in our data set. Other features performed considerably worse (AUC < 0.67). The *NFATC1* body island consists of 102 CpGs, of which only 12 were found differentially methylated in the NHMM analysis. Hence, the smoking associated differences averaged out for location and island status in *NFATC1*. This explains why *NFATC1* body island was not included in the most variable features and not tested in cross-validation analysis.

Supplementary Table 1: CpG sites with smoking-specific differential methylation between tumours and adjacent mucosa.

Supplementary Table 2: CpG sites with differential methylation by smoking status in adjacent mucosa

Probe ID	Chromosomal location	Gene	Gene region	Island status	Mean β -value		LIS p value
					Never	Active	
cg17662683	chr6: 32064146	TNXB	Body	Island	0.44	0.38	4.6×10^{-6}
cg20414186	chr6: 32064491	TNXB	Body	Island	0.30	0.24	5.6×10^{-6}
cg24882324	chr6: 32064508	TNXB	Body	Island	0.43	0.37	7.2×10^{-6}
cg12694372	chr6: 32064582	TNXB	Body	Island	0.50	0.42	8.4×10^{-6}

Supplementary Figure 1: Prediction performances of the top-six methylation features.



Supplementary Figure 2: Methylation of differentiated methylated sites within the *APC* promoter 1A in matched tumours and adjacent mucosa

