

Epigenome-Wide Study Identifies Epigenetic Outliers in Normal Mucosa of Patients with Colorectal Cancer



Jayashri Ghosh¹, Bryant M. Schultz¹, Joe Chan¹, Claudia Wultsch^{2,3}, Rajveer Singh², Imad Shureiqi⁴, Stephanie Chow⁵, Ahmet Doymaz⁶, Sophia Varriano⁷, Melissa Driscoll⁸, Jennifer Muse⁹, Frida E. Kleiman⁶, Konstantinos Krampis^{2,10,11}, Jean-Pierre J. Issa¹², and Carmen Sapienza¹

ABSTRACT

Nongenetic predisposition to colorectal cancer continues to be difficult to measure precisely, hampering efforts in targeted prevention and screening. Epigenetic changes in the normal mucosa of patients with colorectal cancer can serve as a tool in predicting colorectal cancer outcomes. We identified epigenetic changes affecting the normal mucosa of patients with colorectal cancer. DNA methylation profiling on normal colon mucosa from 77 patients with colorectal cancer and 68 controls identified a distinct subgroup of normally-appearing mucosa with markedly disrupted DNA methylation at a large number of CpGs, termed as “Outlier Methylation Phenotype” (OMP) and are present in 15 of 77 patients with cancer versus 0 of 68 controls ($P < 0.001$). Similar findings were also seen in publicly available datasets. Comparison of normal colon mucosa transcription profiles of patients with OMP cancer with those of patients with non-OMP cancer indicates genes whose promoters are hypermethylated in the OMP patients are also transcriptionally downregulated, and that many of the genes most affected are involved in interactions between epithelial cells, the mucus

layer, and the microbiome. Analysis of 16S rRNA profiles suggests that normal colon mucosa of OMPs are enriched in bacterial genera associated with colorectal cancer risk, advanced tumor stage, chronic intestinal inflammation, malignant transformation, nosocomial infections, and KRAS mutations. In conclusion, our study identifies an epigenetically distinct OMP group in the normal mucosa of patients with colorectal cancer that is characterized by a disrupted methylome, altered gene expression, and microbial dysbiosis. Prospective studies are needed to determine whether OMP could serve as a biomarker for an elevated epigenetic risk for colorectal cancer development.

Prevention Relevance: Our study identifies an epigenetically distinct OMP group in the normal mucosa of patients with colorectal cancer that is characterized by a disrupted methylome, altered gene expression, and microbial dysbiosis. Identification of OMPs in healthy controls and patients with colorectal cancer will lead to prevention and better prognosis, respectively.

Introduction

Despite the availability of an effective screening test, colorectal cancer remains the third-leading cause of cancer-related deaths in men and women in the United States (1).

Over the years, scientists have discovered various molecular markers like gene mutations (*KRAS*, *BRAF*, and *APC* genes); CpG island methylator phenotype (CIMP), microsatellite instability, and so on to better understand the heterogeneous

outcomes of colorectal cancer (2). However, it is noteworthy that all of these molecular subtypes are based on investigating the tumor tissues. We, on the other hand, study the normal tissues of patients with colorectal cancer, which could harbor biomarkers to better understand colorectal cancer outcomes.

We have identified site-specific DNA methylation differences in normal colon mucosa that distinguish patients with cancer from patients without cancer with high sensitivity and

¹Fels Cancer Institute for Personalized Medicine, Lewis Katz School of Medicine, Temple University, Philadelphia, Pennsylvania. ²Bioinformatics and Computational Genomics Laboratory, Hunter College, City University of New York, New York, New York. ³Sackler Institute for Comparative Genomics, American Museum of Natural History, New York, New York. ⁴Department of Internal Medicine, University of Michigan, Ann Arbor, Michigan. ⁵Nutrition Department, School of Urban Public Health at Hunter College, New York, New York. ⁶Department of Chemistry, Hunter College, City University of New York, New York, New York. ⁷The Graduate Center, City University of New York, New York, New York. ⁸Northwell Health Imbert Cancer Center, Bayshore, New York. ⁹The Charles Bronfman Institute for Personalized Medicine, Icahn School of Medicine at Mount Sinai, New York, New York. ¹⁰Department of Biological Sciences, Hunter College, City University of New York, New York, New York. ¹¹Institute of Computational Biomedicine, Weill

Cornell Medical College, New York, New York. ¹²Coriell Institute for Medical Research, Camden, New Jersey.

J.-P.J. Issa and C. Sapienza are the co-senior authors of this article.

Corresponding Author: Carmen Sapienza, Fels Cancer Institute for Personalized Medicine, Lewis Katz School of Medicine, Temple University, 3307 N. Broad Street, Room 300, Philadelphia, PA 19140. Phone: 215-707-7373; E-mail: sapienza@temple.edu

Cancer Prev Res 2022;15:755–66

doi: 10.1158/1940-6207.CAPR-22-0258

This open access article is distributed under the Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International (CC BY-NC-ND 4.0) license.

©2022 The Authors; Published by the American Association for Cancer Research

specificity (3). This observation, validated in both an independent population (4) and an animal model (5), suggests that these cancer patient “signature” methylation differences in normal tissues accumulate over time as a result of aging, environmental exposures and, perhaps, genetic influences. Our earlier observation (3) that the largest category of genes affected by differential methylation was those involved in carbohydrate and lipid metabolism is consistent with long-standing epidemiological evidence (6) that dietary factors affect colorectal cancer risk.

Colorectal cancer risk is not distributed uniformly across the population but is higher in patients of African descent than Caucasians, Hispanics or Asians (7). African American (AA) patients with colorectal cancer also appear less likely to develop microsatellite-unstable cancers (a form of colorectal cancer with improved outcome) than their Caucasian counterparts (7). In addition, AA patients who are asymptomatic are more likely to have proximal, large, precancerous adenomatous polyps present on colonoscopy screening (8). While there are likely to be socioeconomic factors involved in disparities in cancer incidence and outcomes, it is also possible that race-associated differences in biology contribute (9).

Our current study was designed to investigate differences in the normal colon epigenome of patients with colorectal cancer by performing genome-wide DNA methylation profiling on 77 patients with colorectal cancer (42 AA and 35 Caucasians) and age-, sex- and race-matched controls (34 AA and 34 Caucasian). We also performed normal colon transcription profiling on selected patients with cancer, as well as microbiome analysis via 16S rRNA sequencing. Our hypothesis is that environmental factors interact with the normal colon epigenome to engender epigenetic changes that predispose to cancer, and that these changes are greater and/or more frequent in AA than in Caucasians. We hypothesize, further, that environmental factors, principally diet, exert much of their effect on the normal colon epigenome through interactions with the microbiome.

Materials and Methods

Samples

Normal colon tissues (fresh frozen) of 77 patients with colorectal cancer were purchased from Fox Chase Cancer Centre biobank. These normal colon tissues adjacent (~10 cm away) to tumors were collected from colorectal cancer patients as described previously (3, 4). Similarly, normal colon tissues (fresh frozen) from age, sex, location, and race matched healthy controls ($n = 70$) were collected during routine screening colonoscopies after informed consent. Controls with previous colonoscopic finding of polyps were excluded.

Written informed consent from the patients was obtained and the study was conducted in accordance with Declaration of Helsinki ethical guidelines and the study was approved by Temple University’s institutional review board.

Sample processing

DNA extraction

Genomic DNA was extracted from colon tissue samples using Invitrogen PureLink genomic DNA kit as per the manufacturer’s protocol.

RNA extraction

RNA was extracted using Qiagen’s RNeasy Plus mini kit. Briefly, nearly 30 mg of colon tissue was homogenized followed by isolation and purification using standard manufacturer’s protocol.

Quantification and quality check

Extracted DNAs and RNAs were quantified using Thermo Fisher’s NanoDrop. RNA integrity was checked on Agilent 2100 Bioanalyzer.

Statistical analyses

All the statistical analyses were done using different packages in R. Plots were made using both R and GraphPad Prism (version 8).

DNA methylation

Illumina EPIC array

Extracted DNA was sent to external Genomic Facility at Penn State University to be run on Illumina’s EPIC array. Prior to array run, extracted DNA was treated with bisulfite using Zymo EZ DNA methylation kit. Bisulfite-treated DNA is processed further to run Illumina’s EPIC array as described previously (10). The output data are generated in the .idat files. Two healthy samples failed to hybridize during the initial array processing.

Data processing

Raw data files from 77 patients with colorectal cancer and 68 healthy controls were preprocessed using minfi’s *preprocessIllumina* function to mimic Genome Studio’s background correction and normalization steps in the R environment. Probe normalization was also done via the *preprocessIllumina* function that equally recreates Genome Studio’s method of normalizing variability in red/green signal using paired red/green control probes in a reference sample. Beta values obtained after these preprocessing steps were used for all the subsequent analyses.

Quality control

The quality of the samples was checked using the minfi getQC test.

Batch effect

As the samples were run in batches, batch effect was checked using correlation and Bland Altman analyses for the replicate samples [both intraplate or interwell replicates (same samples in different wells of the same plate) as well as interplate replicates (same samples in different plates)].

Cell composition/purity

Epithelial cell purity between the tissues from healthy controls and patients with colorectal cancer was estimated by leukocyte unmethylation for purity (LUMP) as described previously (11).

CpG selection for methylation analyses

SNP associated and cross-reactive CpGs (12, 13) and 59 SNP CpGs were excluded from analysis. Poor performing probes (missing values in $\geq 20\%$ of the samples) were also excluded. This resulted in 819,239 CpGs that were included for the analyses.

Cluster analysis

Unsupervised clustering using bootstrap method was performed using the *pyclust* package in R.

Principal component analysis

Principal component analysis was done by using *prcomp* function in R.

Outlier analysis

Outliers or individuals with Outlier Methylation Phenotype (OMP) were identified by following a two-step procedure (14, 15). In the first step, each of the 819,239 CpG sites was analyzed for the presence of outliers [methylation levels beyond 1.5 times the interquartile range below the first quartile (“hypomethylated outliers”) or above the third quartile (“hypermethylated outliers”) of the distribution]. In the second step, the distribution of outlier CpGs was plotted for each sample and similar outlier calculations as in Step 1 were done, to identify individuals with extremely large number of outlier CpGs compared with rest of the population. Outliers of Step 2 were considered as the individuals with OMP.

Differential methylation analysis

Between group comparisons were done using two-sided *t* test for methylation values. Bonferroni correction was used to correct for multiple testing. We checked the location/feature of each of the 819239 CpGs and corrected for 108,498 features (because methylation levels are highly correlated at CpGs within the same feature, and are, thus, not independent) resulting in *P* values less than $4.6E-07$ as the cut off for significance. A cut off (0.05) for magnitude of difference in beta values was also introduced. Hence, CpGs with *P* value less than $4.6E-07$ and magnitude of difference >0.05 were considered to be significant. Differential methylation analyses were done using two-sided *t* test in R. Differential methylated regions (DMR) were identified using “DMRcate” package in R.

Gene expression

RNA sequencing

RNA-sequencing libraries were prepared using Illumina’s TruSeq stranded mRNA kit by following the standard manufacturer’s protocol. Libraries were sequenced in Illumina HiSeq 4000 at GENEWIZ.

Data processing

Sequencing data quality was assessed by FastQC. Sequencing reads were trimmed using Trim Galore and aligned using mapping software STAR (16). Transcripts were counted using HTSeq.

Differential gene expression analysis

We used R package DESeq2, version 3.13 (17) for differential gene expression analyses.

Gene Ontology analysis

We used the R package ReportingTools (<https://bioconductor.org/packages/release/bioc/html/ReportingTools.html>) to generate Gene Ontology (GO) pathways (18).

Microbiome

16S rRNA sequencing

16S rRNA libraries were generated using a modified Illumina 16S protocol that increases input DNA to 62.5 ng. Barcoded libraries were generated with Nextera XT adapters per Illumina’s 16S protocol. Purified libraries were quantified via Qubit and analyzed on the Agilent DNA Bioanalyzer in order to generate 10 mmol/L pooled libraries to be sequenced on the MiSeq platform.

Amplicon sequence variants

We pre-processed raw 16S rRNA sequences generated for 70 colon tissue samples collected from AA patients using QIIME2, version 2019.1 (19). We obtained a total of 3,845,964 quality-screened DNA sequences, with an average count of 54,942 sequence reads per sample. We applied the DADA2 algorithm (20) via the *q2-dada2* plugin to denoise the sequence data and generate unique amplicon sequence variants (ASV). Taxonomic classification of representative ASVs was conducted using the classify-sklearn naïve Bayes classifier (21) against the Greengenes, version 13_8 99% reference database (22).

Taxonomic composition and differential abundance

We used R package *phyloseq*, version 1.24.2 (23) to describe the taxonomic composition of each cohort at the phylum and genus level. In addition, differential abundance analysis using R package DESeq2, version 3.13 (17) was applied to identify bacterial taxa that were significantly different between the cohorts studied. Differential abundances in bacterial species were assessed using a log2foldchange value, and cohort comparisons were conducted applying the Wald test with the Benjamini-Hochberg correction.

Microbiome diversity

A rarefied sampling depth of 14,214 DNA reads per sample and R package *phyloseq*, version 1.24.2 (23) were further used to assess microbiome diversity across sampling cohorts. Diversity within samples (alpha diversity) was estimated as observed number of ASVs and Shannon diversity index and significance of differences was tested using nonparametric Wilcoxon rank sum tests. Rarefied samples were also used to calculate

Bray–Curtis beta diversity (dissimilarity between samples), and nonmetric multidimensional scaling (NMDS) was performed. Significance of differences in beta diversity between cohorts was assessed by permutational analysis of variance (PERMANOVA) and permutation tests for homogeneity in multivariate dispersion (PERMDISP) in R package *vegan*, version 2.5–6 (24) with 999 permutations.

Data availability statement

The datasets generated during this study are available in the GEO repository (GSE 199057).

Results

Quality control of methylation dataset

All the samples passed the quality control (QC) test on minfi (Supplementary Fig. S1A). No batch effects were observed for the processed methylation data. All the replicates (irrespective of intraplate or interplate) were strongly correlated ($R^2 = 0.99$). Similarly, all the replicates (Supplementary Fig. S1B–S1E) showed similar results on Bland Altman analyses wherein nearly 45K CpGs (5%) were outside agreement boundary irrespective of whether those were intraplate (Supplementary Fig. S1B and S1C) or interplate (Supplementary Fig. S1D and S1E) replicates. Furthermore, there was no difference ($P = 0.4626$) in the cell purity of normal tissues from healthy controls and patients with colorectal cancer on LUMP analysis (Supplementary Fig. S1F).

Identification of an outlier methylation group in normal tissues of cancer patients

We performed unsupervised hierarchical cluster analysis, using methylation data from 819,239 CpGs to determine whether our study population could be subdivided on the basis of the normal colon epigenome. Interestingly, we observed a group of 14 colorectal cancer individuals (11 AA and 3 Caucasians) and a Caucasian colorectal cancer patient clustering separately (highlighted in yellow in Fig. 1A) from rest of the dataset. We also performed principal component analysis to determine whether quantitative variation at multiple sites might distinguish the study groups (Fig. 1B). Patients without cancer were less variable compared to the colon cancer groups of both races. The Caucasian healthy (CH) group had the least variability followed by the AA healthy group (AH) group. The AA cancer (AC) group had the highest variability followed by the Caucasian cancer group (CC). Very high variability in the cancer groups was exacerbated by the samples (11AC, 4CC) at the right side of the PCA plot (values >590 in PC1, samples within the black ellipse). It is noteworthy, that these 15 individuals are the ones that cluster separately in Fig. 1A.

Definition of an OMP group

Because both PCA and cluster analysis suggested the existence of a group with dramatically disrupted normal

tissue methylomes, we applied the same metric we have used previously (14, 15) to identify individuals with “Outlier Methylation Phenotype” (OMP) (Fig. 2). Although this method (see Materials and Methods) transforms a fundamentally quantitative trait (methylation values) into a discrete classifier (OMP status), it simplifies further analysis of factors that may contribute to this phenotype. In other words, converting a quantitative variable to a categorical variable simplifies the downstream analysis for better characterization of this group (OMP). We plotted the number of CpGs in which an individual was hyper- (Fig. 2A) or hypo-methylated (Fig. 2B) at greater than 1.5-times the interquartile range to identify those individuals who were OMPs.

None of the CH individuals were hyper- or hypo-methylated outliers (Fig. 2A and B), whereas two of the AH individuals were hypo-methylated outliers. Fifteen AC patients were hyper- or hypo-methylated outliers and 11 were bidirectional (both hyper- and hypo-methylated) outliers. Among CC patients, seven samples were hyper-methylated and five samples were hypo-methylated outliers, of which only four patients were bidirectional outliers. Individuals who were outliers in both hyper- and hypo-methylated plots were classified as OMP (14, 15). Further justification for classifying only bidirectional outliers as OMPs (11 AC and 4 CC) is that these individuals are the same patients who form separate groups in the cluster analyses (Fig. 1A) and are furthest from the other patients with colorectal cancer in the PCA analysis (Fig. 1B).

We also analyzed whether the OMP (red box) and non-OMP (blue box) clusters (Fig. 1A) were based on any particular feature (like age, sex). As shown in Supplementary Table S1, these two clusters showed significant differences in cancer status. All other variables (age, sex, location) were not significantly different. Furthermore, we did see a borderline association ($P = 0.05$) for race but the significance was lost after correcting for multiple (four) tests.

Validation of OMP group in publicly available colorectal cancer datasets

We selected three colorectal cancer datasets from Gene Expression Omnibus (GEO) which had 450K methylation array data for both healthy controls and normal tissues from patients with colorectal cancer. We performed outlier analysis in each of the datasets and identified individuals with OMPs (or bidirectional outliers) as described above. As shown in Supplementary Table S2A, all of the datasets show higher frequency of OMPs in the colorectal cancer group compared to healthy controls. Additionally, the largest dataset (GSE132804) had significantly higher frequency of OMPs in patients with colorectal cancer compared to controls. Furthermore, we also analyzed if any confounding variables in dataset GSE132804 influenced the OMP output. Supplementary Table S2B clearly indicates that the two groups (cancer and controls) were matched for age, sex, and location, justifying that OMP is not an outcome of unbalanced covariables. This validates our

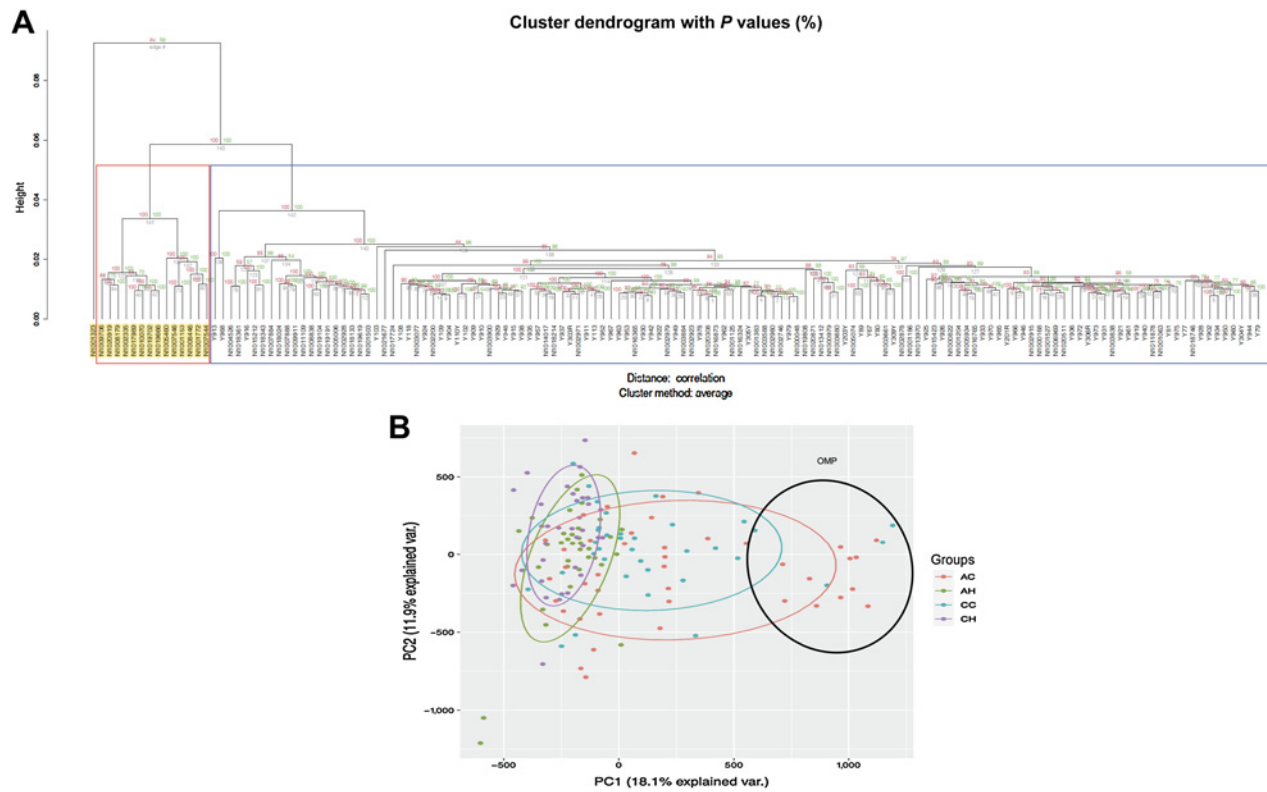


Figure 1. Analysis of methylation data. **A**, Unsupervised cluster analysis of study samples. Hierarchical cluster plots using unsupervised cluster analysis showing separate cluster for the OMPs. **B**, Principal component analyses. Principal component analyses of study groups using 819,239 CpGs.

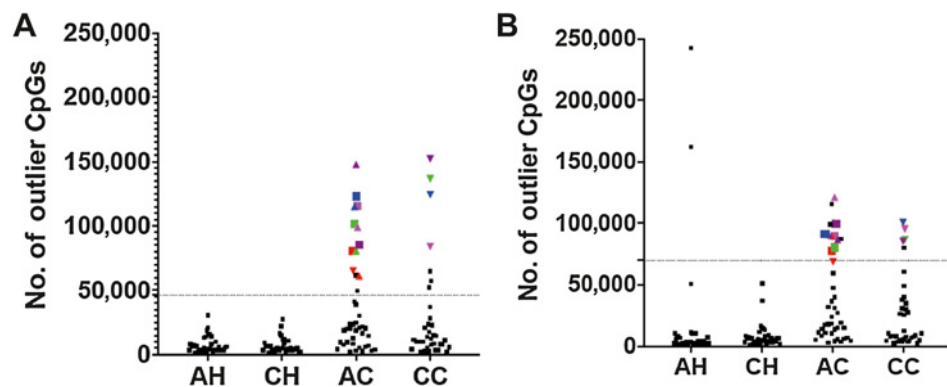
finding that normal tissues of patients with colorectal cancer are more prone to have disrupted epigenome or OMP characteristic compared with healthy controls.

Effect of OMPs on differential methylation in normal colon mucosa of colorectal cancer patients

AA and Caucasian patients with colorectal cancer, combined, showed significantly different methylation at 85,178 CpGs (10.40%) compared with healthy controls (Fig. 3A). On race-stratified subgroup analysis, the AC patients had

26,803 differentially methylated CpGs compared with the AH controls (Fig. 3B), whereas the CC patients had 12,016 (Fig. 3C) differentially methylated CpGs compared with the (CH) controls. More than 60% of the differentially methylated CpGs (7,341 CpGs) in the Caucasian patients with colorectal cancer were also differentially methylated in AA patients with colorectal cancer (Fig. 3B and C), suggesting that many of the cancer-associated methylation alterations were common to both AC and CC patients. However, AA patients with colorectal cancer had a much larger number of abnormally

Figure 2. Identification of samples with OMP. Number of CpGs in which a sample is hypermethylated outlier (A) or Hypomethylated outlier (B). Dotted line indicates outlier boundary. Each symbol is a sample. Symbols above the dotted lines are outliers in respective plots. Colored symbols indicate samples that are outliers in both the plots and are termed as “OMPs”. Samples represented by colored symbols are OMPs. Same color and shape show the same individuals in both the plots.



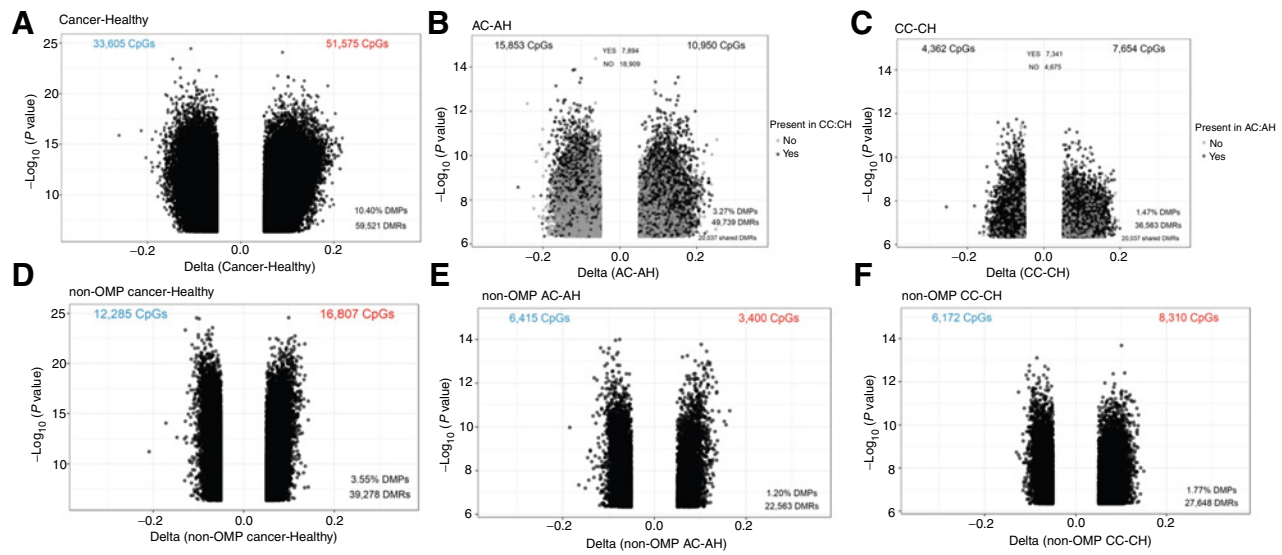


Figure 3. Volcano plots showing differential methylation analyses. Colon cancer versus healthy in all samples (A), AAs (B), and Caucasians (C). Non-outlier colon cancer versus healthy in all samples (D), AAs (E), and Caucasians (F). DMPs, differentially methylated positions; DMRs, differentially methylated regions.

methylated CpG sites compared with their healthy controls (an additional >14,000 CpG sites), than did their Caucasian counterparts.

We also analyzed whether any confounding effects between the control and colorectal cancer groups account for these differences. **Table 1** shows the demographic profile of the analyzed samples. None of the variables were significantly different between cancer and control groups. Hence, all the differentially methylated probes (overall or race-specific) are associated with colorectal cancer.

Because we had identified a group of patients with dramatically disrupted normal colon methylation profiles, and the groups was composed of largely AA patients, we asked whether the increased number of differences between AC and AH groups compared with the CC and CH groups were driven by the OMPs by excluding them from the analysis. When this was done, the number of differentially methylated CpGs was reduced by more than 50% in overall cancer versus healthy comparison (**Fig. 3D**). A similar trend was observed in AC versus AH (**Fig. 3E**). However, we did not observe a reduction in abnormally methylated CpGs between the CC versus CH groups (**Fig. 3F**), suggesting that OMPs in the AC group contributed much more variance than in the CC group.

It is noteworthy that DNA methylation profiles of normal colon mucosa between the controls and colorectal cancer patients of AA and Caucasian races are mostly similar. We observed a very small fraction of race-associated differences in site-specific CpG methylation between either healthy controls (0.10%, or 794 sites) or between cancer patients (0.02%, or 193 sites) (Supplementary Fig. S2). These observations suggest that racial disparities in colon cancer incidence and outcome are not

Table 1. Demographic profile of analyzed samples.

| | All samples | | |
|------------------|-------------------|------------------|----------------|
| | Cancer (n = 77) | Control (n = 68) | P ^a |
| Age (mean± SD) | 57.67 ± 9.68 | 56.81 ± 8.81 | 0.5757 |
| Sex | | | 1.0000 |
| Males | 38 | 33 | |
| Females | 39 | 35 | |
| Race | | | 0.6198 |
| Caucasian | 35 | 34 | |
| African American | 42 | 34 | |
| Location | | | 1.0000 |
| Distal | 42 | 37 | |
| Proximal | 35 | 31 | |
| | Caucasians | | |
| | Cancer (n = 35) | Control (n = 34) | P ^a |
| Age (mean± SD) | 56.80 ± 9.33 | 56.06 ± 9.61 | 0.6905 |
| Sex | | | 1.0000 |
| Males | 15 | 15 | |
| Females | 20 | 19 | |
| Location | | | 1.0000 |
| Distal | 19 | 19 | |
| Proximal | 16 | 15 | |
| | African Americans | | |
| | Cancer (n = 42) | Control (n = 34) | P ^a |
| Age (mean± SD) | 58.40 ± 10.02 | 57.56 ± 8.00 | 0.7461 |
| Sex | | | 1.0000 |
| Males | 23 | 18 | |
| Females | 19 | 16 | |
| Location | | | 1.0000 |
| Distal | 23 | 18 | |
| Proximal | 19 | 16 | |

^at test for age and Fisher exact test for other variables.

a result of large numbers of methylation differences at different CpG sites, with the caveat that not all CpG sites are interrogated by the Illumina platform used.

Each of the above analyses (cluster, PCA, outlier, differential methylation) indicates the presence of a highly epigenetically disrupted group of patients with colorectal cancer, of which the majority are AA. We examined this OMP group of patients, further, to determine what factors might influence this phenotype, and whether it might contribute to observed racial disparity in colorectal cancer incidence and outcome.

Differential expression of genes with differentially methylated promoter CpGs in AA OMPs

A working hypothesis on racial disparities in colon cancer developed from our analysis of normal tissue DNA methylation is that OMPs, although not unique to AAs, are more prevalent among AAs and OMPs may be at higher risk of cancer. It is noteworthy that of the 178,469 CpGs that were differentially methylated between OMP cancer patients and non-OMP cancer patients (Supplementary Fig. S3), 40,961 CpGs were present in the promoter regions of 11,357 genes.

Again, because the majority (~75%) of OMPs were African American and we wished to characterize this group further, we compared gene expression levels between OMPs (AO) and non-OMPs (AC) among AA patients with colorectal cancer for whom we were able to obtain normal colon RNA samples by bulk RNA sequencing (3 OMPs vs. 5 non-OMPs). More than 17% (1,964) of the promoter differentially methylated genes also exhibited differential expression levels (Supplementary Fig. S4). The majority (1,151 genes) of the differentially expressed genes were hypermethylated in the promoters of OMPs. As expected, most of these hypermethylated genes (1,021 or 88.7%) were downregulated in the OMPs compared with non-OMPs (Supplementary Fig. S4).

Supplementary Table S3 lists the differentially expressed genes. Multiple genes linked to mucins (*MUC17*, *MUC3A*, *MUC12*, *MUC4*, *MUC5B*, *MUC20*, *MUC2*, *MUC13*, *MUC1*); claudins (*CLDN8*, *CLDN3*, *CLDN4*, *CLDN7*, *CLDN12*, *CLDN9*), cadherins (*CDHR2*, *CDHR5*, *CDH1*, *CDH17*, *CDHR1*) and other transmembrane junction proteins (*DSC2*, *CGN*, *CAPN13*, *CDHR2*, *TMPRSS2*, *AMN*) were differentially expressed (down regulated) in OMPs. In addition, among those genes that were hypo-methylated (Supplementary Table S3), the proinflammatory cytokine genes *IL6* and *IL11* were both upregulated. The top significant biological processes (Supplementary Table S4) associated with the differentially expressed genes included xenobiotic processes (response to xenobiotic stimulus, xenobiotic metabolic process), leading us to perform an analysis of gut microbiome components.

Differential microbiome in OMPs (AO) compared to non-OMPs (AC) in AA patients with colorectal cancer

Similar to expression analysis, additional microbiome analysis was restricted to AA patients and included 35 AH, 25 AC and 10 AO patients. In total, we identified 18,522 ASVs

across all samples analyzed. At the phylum level and across all cohorts, the microbiota was dominated by ASVs assigned to the phyla Firmicutes, Bacteroidetes, Proteobacteria, Actinobacteria, Fusobacteria, and Verrucomicrobia (Fig. 4A). At the genus level, ASVs were assigned primarily to the genera *Bacteroides*, *Oscillospira*, *Clostridium*, *Coproccoccus*, *Prevotella*, and *Ruminococcus* (Fig. 4B).

Although neither alpha nor beta diversity estimates were significantly different between AH, AC, and AO cohorts (Wilcoxon rank sum tests $P > 0.05$, Supplementary Fig. S5A; PERMANOVA, $P = 0.084$, Supplementary Fig. S5B), differential abundance analysis (Supplementary Fig. S5C) revealed that significant differences among cohorts were detected in the phyla Actinobacteria, Bacteroidetes, Firmicutes, and Proteobacteria. More specifically, we detected an increased abundance of the *Eubacterium* genus in AC tissues when directly compared with taxonomic profiles of the AO cohort, whereas the genera *Fusobacterium*, *Phascolarctobacterium*, *Bacteroides*, *Roseburia*, *Dialister*, *Stenotrophomonas*, and *Ruminococcus* were more prevalent in the AO cohort (Fig. 4C).

Discussion

We performed genome-wide DNA methylation profiling on normal colon mucosa from African American and Caucasian colorectal cancer patients and age-, sex-, and race-matched controls. Our hypothesis was that colorectal cancer incidence and outcome were associated with underlying differences in the normal tissue epigenome.

Unsupervised hierarchical cluster analysis (Fig. 1A) and principal component analysis (Fig. 1B) both suggested the existence of a separate group of colorectal cancer patients with dramatically disrupted normal tissue methylomes. Interestingly, we were also able to identify this same group of epigenetically disrupted individuals by using a simple metric of outlier determination as used previously by our group (14, 15). We have termed this group of patients with colorectal cancer as OMP (see also Fig. 2). We also identified OMPs in publicly available colorectal cancer datasets. OMP frequencies varied from 1% to 2% in controls and 8% to 30% in patients with colorectal cancer. This clearly suggests that patients with colorectal cancer are more prone to develop OMPs compared with controls. Furthermore, the varying percentage of OMPs among patients with colorectal cancer (<10% in GSE48684 and GSE131013; and ~30% in GSE 132804) in these datasets could be explained by smaller sample size (24 patients with colorectal cancer in GSE48684) or difference in ethnicity (Spanish population in GSE131013). Unfortunately, these datasets do not have any AA samples, so we could not perform race-specific analyses.

While we identified many differences in average site-specific methylation between patients with colorectal cancer and controls, confirming and extending our previous studies (3, 4), the major difference we identified between AA and Caucasian

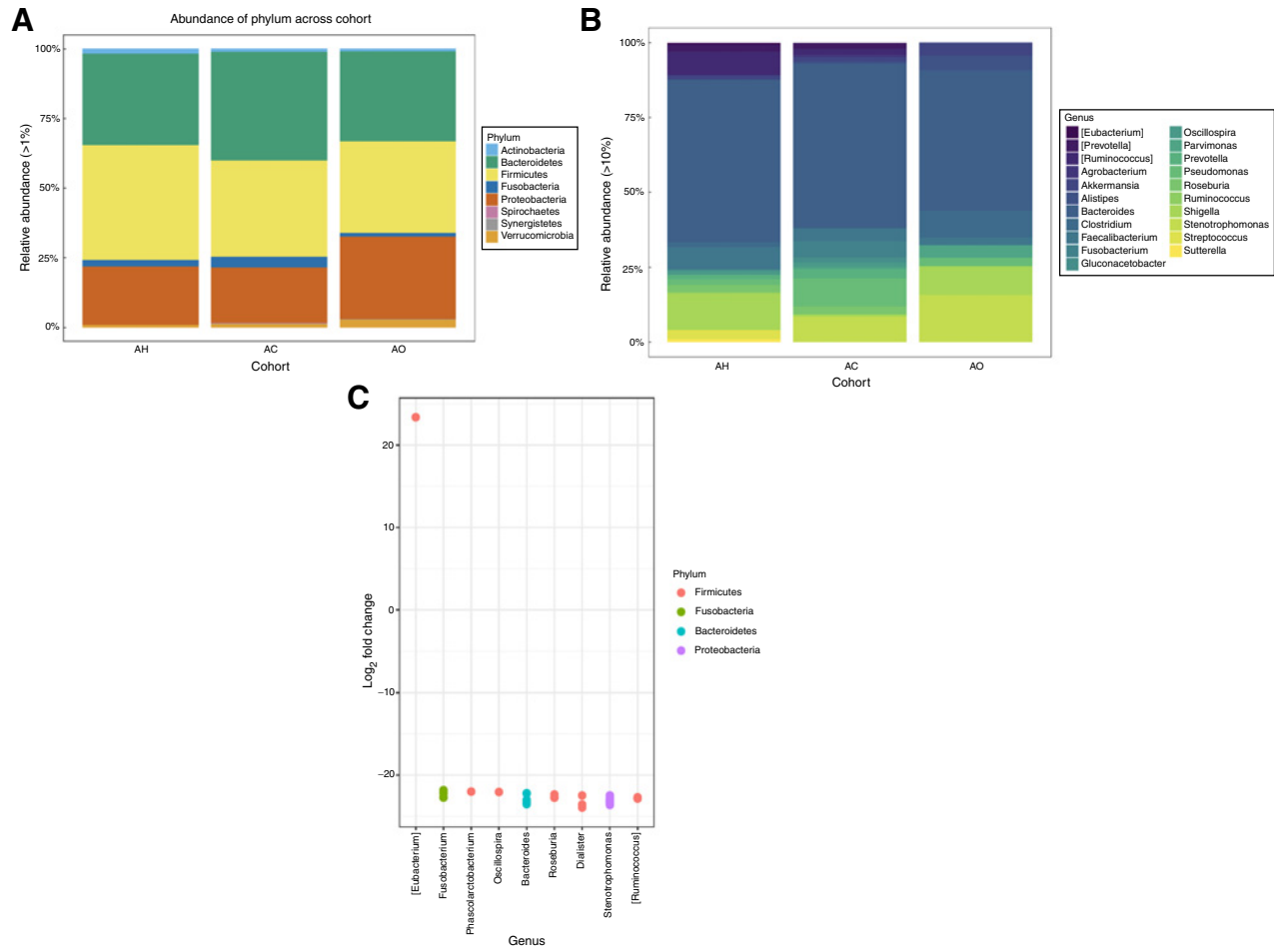


Figure 4. Analysis of the microbiome from AA samples. **A**, Taxonomic composition of colon tissue microbiomes at the phylum level. **B**, Taxonomic composition of colon tissue microbiomes at the genus level. **C**, Differential abundance analysis between microbiome samples of AC and AO cohorts. Each point represents ASV belonging to respective bacteria species. ASVs were considered significant if their false discovery rate-corrected *P* value was < 0.05. Multiple points visualized under the same genus represent ASVs that are classified within the same genus but differ by one or more nucleotides. Taxa in square brackets are annotations for proposed taxonomy supplied by the Greengenes database.

patients with colorectal cancer was in the number of patients with OMP (15). AAs (26%) were more than twice as likely to be OMPs compared with Caucasian (11%) patients with colorectal cancer. Furthermore, AA patients with colorectal cancer displayed higher abnormality in methylation profiles (AC vs. AH) than their Caucasian counterparts (CC vs. CH). However, methylation differences between the AC and AH groups were greatly reduced on excluding the OMPs, suggesting a substantial role for OMPs in causing epigenetic disbalances in AA patients with colorectal cancer.

Because the frequency of OMPs appears higher among African Americans (Fig. 2), and OMPs have sometimes been associated with undesirable outcomes in other diseases (14), as well as cancer (25, 26), a greater frequency of OMPs among AA patients with colorectal cancer could be associated with racial disparities in colorectal cancer incidence and outcome. However, too few OMPs have been identified to determine whether

this unusual molecular phenotype is associated with any clinical outcome or any established molecular subtype in patients with colorectal cancer. However, it is noteworthy, that our previous study on OMP in TCGA data showed that OMP is independent of CIMP (15). Another important aspect of cancer including colorectal cancer is the significance of epigenetic aging in tumorigenesis, and its potential use for cancer risk prediction (27). It would be interesting to further evaluate if OMPs have epigenetic age drift in normal tissues, which could be used as a predictive and prognostic tool. Nevertheless, determining the cause of OMP in normal tissues is of interest because of its potential to affect gene expression in normal colon mucosa, as well as the potential for environmental factors to influence this phenotype.

Our analysis of gene expression, comparing normal colon mucosa of OMP cancer patients with non-OMP cancer patients, indicated that the major pathways differentially

affected in OMP patients were involved in repression of genes mediating the interaction between the intestinal epithelium/mucus barrier and the microbiome. For instance, a number of genes from the cadherin superfamily, claudins and other transmembrane junction proteins were downregulated in the OMP group. Cadherins and claudins are integral parts of adherens and tight junctions, respectively. Cadherins are important cell adhesion molecules and loss of cell adhesion, specifically by downregulation of E-cadherin (*CDH1*) has been associated with malignant characteristics including tumor progression, loss of differentiation, invasion and metastasis (28). On the other hand, claudins are transmembrane proteins that maintain the barrier functioning of tight junctions (29). Clearly, loss of expression of these and other transmembrane junction proteins leads to deregulation of normal tissue function and development of epithelium-related diseases, including cancer (30). Furthermore, genes belonging to the mucin family were downregulated in OMP cancer patients. Aberrant mucin expression is linked to chronic inflammation and colorectal cancer, as mucus functions as a physical barrier and influences microbial composition by providing nutrients and attachment sites for the microbial community (31).

Analysis of the microbiome further showed differential abundance of several genera between OMPs versus non-OMP colorectal cancer patients. The genus *Eubacterium* was found to be in lower abundance in OMPs in our study. Interestingly, the abundance of *Eubacterium hallii*, and *Eubacterium ventriosum* were found to be significantly higher in healthy samples than in colorectal cancer samples (32). *E. hallii* utilizes glucose and the fermentation intermediates acetate and lactate to form butyrate and hydrogen, which are important in maintaining intestinal metabolic balance (33).

Fusobacterium and *Bacteroides*, which are among the most prominent colorectal cancer-associated bacteria, were highly abundant in OMPs compared with non-OMPs (34). *Fusobacterium* is also known to be associated with microsatellite instability (MSI), hypermethylation and malignant transformation of epithelial cells (35). On the other hand, *Bacteroides fragilis* cause a series of inflammatory reactions due to *B. fragilis* toxin (BFT), which leads to chronic intestinal inflammation and tissue injury and plays a crucial role leading to colorectal cancer (36).

Other genera found to be in higher abundance in OMPs, such as *Phascolarctobacterium*, *Roseburia*, *Ruminococcus*, *Dialister* and *Stenotrophomonas* have also been reported to be in higher abundance in patients with colorectal cancer in other studies (37–40). Furthermore, *Ruminococcus gnavus* has been positively associated with KRAS mutations (a known colorectal cancer mutation) (41). Recent studies have also highlighted the role of *Dialister pneumosintes* in advanced colorectal cancer patients (42). *Stenotrophomonas maltophilia* is a nosocomial pathogen which is found in higher abundance in colorectal cancer patients after radio or chemotherapy (43).

A recent study (44) showed that the overall microbial composition in normal adjacent tissues is relatively similar to their tumor tissues, with the exceptions of some bacteria which show different prevalence between these two tissue types. This suggests that some of the microbiome changes that we observe may be affected by the presence of an adjacent neoplasm.

AA race is widely understudied and underrepresented in both publicly available datasets (like TCGA) and tissue biobanks. We were limited by the number of African American biospecimens available in the biobank. It is to be noted that some of the largest colorectal cancer biobanks and Consortia have negligible representation of African Americans.

Although our sample size was insufficient to clinically characterize (like tumor grade, side of tumor, age, sex) the OMP group, analysis of the microbiome clearly reflected that normal colon mucosa of OMPs are enriched in bacterial genera associated with colorectal cancer risk, advanced tumor stage, chronic intestinal inflammation, malignant transformation, nosocomial infections, and KRAS mutations. These observations suggest that OMP patients may have microbial dysbiosis that is distinct from that of non-OMP patients.

In conclusion, we identified a distinct group of highly abnormally methylated colorectal cancer patients, termed “OMPs”, and validated their existence using multiple statistical approaches and in multiple datasets. This epigenetically disrupted OMP group was more prevalent among AA patients with colorectal cancer than Caucasian patients with colorectal cancer. Furthermore, we showed that the vast majority of methylation differences between AA patients with colorectal cancer and healthy controls are driven by this OMP group. We were also able to demonstrate downregulation of crucial genes in the OMP group, especially mucins and transmembrane junction genes. Finally, microbiome analysis showed higher abundance of microbial genera that are associated with colorectal cancer risk, malignancy and advanced tumor stage in OMP cancer patients compared with non-OMP cancer patients.

Whether these differences might be a cause or effect of normal colon OMP is unclear. Such questions are only likely to be answered by examination of a much larger number of OMP patients. In this regard, a major consideration for future studies is the relative rarity of OMP individuals, and a major weakness of the present study is the small number of OMP individuals examined. If OMPs are, in fact, more prevalent among patients of African ancestry, examination of a much larger number of such patients might shed additional light on the significance of this phenotype, as well as whether it might be associated with observed racial disparities in colon cancer incidence and outcome.

Authors' Disclosures

B.M. Schultz reports grants from NIH during the conduct of the study. J.-P.J. Issa reports grants from NIH during the conduct of the study; personal fees from Daiichi outside the submitted work. C. Sapienza reports

grants from U54 from NIH, grants from PA Cure, and grants from R21 from NIH during the conduct of the study. No disclosures were reported by the other authors.

Disclaimer

The content is solely the responsibility of the authors and does not necessarily represent the official views of the NCI or the NIH.

Authors' Contributions

J. Ghosh: Data curation, formal analysis, validation, investigation, visualization, methodology, writing—original draft, writing—review and editing. **B.M. Schultz:** Formal analysis, validation, investigation, methodology, writing—review and editing. **J. Chan:** Formal analysis, investigation, methodology, writing—review and editing. **C. Wultsch:** Formal analysis, investigation, methodology, writing—original draft, writing—review and editing. **R. Singh:** Formal analysis, validation, writing—review and editing. **I. Shureiqi:** Resources, writing—review and editing. **S. Chow:** Formal analysis, investigation, methodology, writing—review and editing. **A. Doymaz:** Investigation, methodology, writing—review and editing. **S. Varriano:** Investigation, methodology, writing—review and editing. **M. Driscoll:** Formal analysis, investigation, methodology, writing—review and editing. **J. Muse:** Formal analysis, investigation, methodology, writing—review and editing. **F.E. Kleiman:** Conceptualization, supervision, funding acquisition, writing—review and editing. **K. Krampis:** Formal analysis, supervision, funding acquisition, methodology, writing—review and editing. **J.-P.J. Issa:** Conceptualization, resources, supervision, funding

acquisition, project administration, writing—review and editing. **C. Sapienza:** Conceptualization, resources, supervision, funding acquisition, visualization, writing—original draft, project administration, writing—review and editing.

Acknowledgments

This work was supported by TUFCCC/HC Regional Comprehensive Cancer Health Disparity Partnership, Award Number U54 CA221704 (5) from the NCI (to C. Sapienza, J.-P.J. Issa, F.E. Kleiman, K. Krampis). Work in the Issa laboratory is supported by NIH grants CA214005 (to J.-P.J. Issa). Work in the Sapienza laboratory is also supported by two PA Cure grant 914103047 (to C. Sapienza) from Pennsylvania Department of Health and R21 CA264213 (to C. Sapienza and J. Ghosh) from the NIH.

The publication costs of this article were defrayed in part by the payment of publication fees. Therefore, and solely to indicate this fact, this article is hereby marked “advertisement” in accordance with 18 USC section 1734.

Note

Supplementary data for this article are available at Cancer Prevention Research Online (<http://cancerprevres.aacrjournals.org/>).

Received May 26, 2022; revised July 13, 2022; accepted August 23, 2022; published first August 25, 2022.

References

- Siegel RL, Miller KD, Goding Sauer A, Fedewa SA, Butterly LF, Anderson JC, et al. Colorectal cancer statistics, 2020. *CA Cancer J Clin* 2020;70:145–64.
- Guinney J, Dienstmann R, Wang X, de Reyniès A, Schlicker A, Soneson C, et al. The consensus molecular subtypes of colorectal cancer. *Nat Med* 2015;21:1350–6.
- Silviera ML, Smith BP, Powell J, Sapienza C. Epigenetic differences in normal colon mucosa of cancer patients suggest altered dietary metabolic pathways. *Cancer Prev Res* 2012;5:374–84.
- Cesaroni M, Powell J, Sapienza C. Validation of methylation biomarkers that distinguish normal colon mucosa of cancer patients from normal colon mucosa of patients without cancer. *Cancer Prev Res* 2014;7:717–26.
- Leclerc D, Pham DN, Lévesque N, Truongcao M, Foulkes WD, Sapienza C, et al. Oncogenic role of PDK4 in human colon cancer cells. *Br J Cancer* 2017;116:930–6.
- Giovannucci E, Willett WC. Dietary factors and risk of colon cancer. *Ann Med* 1994;26:443–52.
- Carethers JM. Clinical and genetic factors to inform reducing colorectal cancer disparities in African Americans. *Front Oncol* 2018;8:531.
- Lieberman DA, Williams JL, Holub JL, Morris CD, Logan JR, Eisen GM, et al. Race, ethnicity, and sex affect risk for polyps >9 mm in average-risk individuals. *Gastroenterology* 2014;147:351–8.
- Zavala VA, Bracci PM, Carethers JM, Carvajal-Carmona L, Coggins NB, Cruz-Correa MR, et al. Cancer health disparities in racial/ethnic minorities in the United States. *Br J Cancer* 2021;124:315–32.
- Mani S, Ghosh J, Lan Y, Senapati S, Ord T, Sapienza C, et al. Epigenetic changes in preterm birth placenta suggest a role for ADAMTS genes in spontaneous preterm birth. *Hum Mol Genet* 2019;28:84–95.
- Aran D, Sirota M, Butte AJ. Systematic pan-cancer analysis of tumour purity. *Nat Commun* 2015;6:8971.
- Pidsley R, Zotenko E, Peters TJ, Lawrence MG, Risbridger GP, Molloy P, et al. Critical evaluation of the Illumina MethylationEPIC BeadChip microarray for whole-genome DNA methylation profiling. *Genome Biol* 2016;17:208.
- Zhou W, Laird PW, Shen H. Comprehensive characterization, annotation and innovative use of Infinium DNA methylation BeadChip probes. *Nucleic Acids Res* 2017;45:e22.
- Ghosh J, Mainigi M, Coutifaris C, Sapienza C. Outlier DNA methylation levels as an indicator of environmental exposure and risk of undesirable birth outcome. *Hum Mol Genet* 2016;25:123–9.
- Ghosh J, Schultz B, Coutifaris C, Sapienza C. Highly variant DNA methylation in normal tissues identifies a distinct subclass of cancer patients. *Adv Cancer Res* 2019;142:1–22.
- Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* 2013;29:15–21.
- Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol* 2014;15:550.
- Huntley MA, Larson JL, Chaivorapol C, Becker G, Lawrence M, Hackney JA, et al. ReportingTools: an automated result processing and presentation toolkit for high-throughput genomic analyses. *Bioinformatics* 2013;29:3220–1.
- Bolyen E, Rideout JR, Dillon MR, Bokulich NA, Abnet CC, Al-Ghalith GA, et al. Reproducible, interactive, scalable and extensible microbiome data science using QIIME 2. *Nat Biotechnol* 2019;37:852–7.
- Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJ, Holmes SP. DADA2: High-resolution sample inference from Illumina amplicon data. *Nat Methods* 2016;13:581–3.

21. Bokulich NA, Kaehler BD, Rideout JR, Dillon M, Bolyen E, Knight R, et al. Optimizing taxonomic classification of marker-gene amplicon sequences with QIIME 2's q2-feature-classifier plugin. *Microbiome* 2018;6:90.
22. McDonald D, Price MN, Goodrich J, Nawrocki EP, DeSantis TZ, Probst A, et al. An improved Greengenes taxonomy with explicit ranks for ecological and evolutionary analyses of bacteria and archaea. *ISME J* 2012;6:610–8.
23. McMurdie PJ, Holmes S. phyloseq: an R package for reproducible interactive analysis and graphics of microbiome census data. *PLoS One* 2013;8:e61217.
24. Oksanen J, Blanchet FG, Kindt R, Legendre P, Minchin PR, O'hara R, et al. Package 'vegan'. Volume 2(9): Community ecology package, version; 2013. p. 1–295.
25. Teschendorff AE, Gao Y, Jones A, Ruebner M, Beckmann MW, Wachter DL, et al. DNA methylation outliers in normal breast tissue identify field defects that are enriched in cancer. *Nat Commun* 2016;7:10478.
26. Panjarian S, Madzo J, Keith K, Slater CM, Sapienza C, Jelinek J, et al. Accelerated aging in normal breast tissue of women with breast cancer. *Breast Cancer Res* 2021;23:58.
27. Yu M, Hazelton WD, Luebeck GE, Grady WM. Epigenetic aging: more than just a clock when it comes to cancer. *Cancer Res* 2020;80:367–74.
28. Christou N, Perraud A, Blondy S, Jauberteau MO, Battu S, Mathonnet M. E-cadherin: a potential biomarker of colorectal cancer prognosis. *Oncol Lett* 2017;13:4571–6.
29. Chiba H, Osanai M, Murata M, Kojima T, Sawada N. Transmembrane proteins of tight junctions. *Biochim Biophys Acta* 2008;1778:588–600.
30. Bujko M, Kober P, Mikula M, Ligaj M, Ostrowski J, Siedlecki JA. Expression changes of cell-cell adhesion-related genes in colorectal tumors. *Oncol Lett* 2015;9:2463–70.
31. Coleman OI, Haller D. Microbe-mucus interface in the pathogenesis of colorectal cancer. *Cancers* 2021;13:616.
32. Ai D, Pan H, Li X, Gao Y, Liu G, Xia LC. Identifying gut microbiota associated with colorectal cancer using a zero-inflated lognormal model. *Front Microbiol* 2019;10:826.
33. Engels C, Ruscheweyh HJ, Beerenwinkel N, Lacroix C, Schwab C. The common gut microbe eubacterium hallii also contributes to intestinal propionate formation. *Front Microbiol* 2016;7:713.
34. Ternes D, Karta J, Tsenkova M, Wilmes P, Haan S, Letellier E. Microbiome in colorectal cancer: how to get from meta-omics to mechanism? *Trends Microbiol* 2020;28:401–23.
35. Zhou Z, Chen J, Yao H, Hu H. *Fusobacterium* and colorectal cancer. *Front Oncol* 2018;8:371.
36. Cheng WT, Kantilal HK, Davamani F. The mechanism of bacteroides fragilis toxin contributes to colon cancer formation. *Malays J Med Sci* 2020;27:9–21.
37. Flemer B, Lynch DB, Brown JM, Jeffery IB, Ryan FJ, Claesson MJ, et al. Tumour-associated and non-tumour-associated microbiota in colorectal cancer. *Gut* 2017;66:633–43.
38. Loftus M, Hassouneh SA, Yooseph S. Bacterial community structure alterations within the colorectal cancer gut microbiome. *BMC Microbiol* 2021;21:98.
39. Weir TL, Manter DK, Sheflin AM, Barnett BA, Heuberger AL, Ryan EP. Stool microbiome and metabolome differences between colorectal cancer patients and healthy adults. *PLoS One* 2013;8:e70803.
40. Peters BA, Dominianni C, Shapiro JA, Church TR, Wu J, Miller G, et al. The gut microbiota in conventional and serrated precursors of colorectal cancer. *Microbiome* 2016;4:69.
41. Hong BY, Ideta T, Lemos BS, Igarashi Y, Tan Y, DiSiena M, et al. Characterization of mucosal dysbiosis of early colonic neoplasia. *NPJ Precis Oncol* 2019;3:29.
42. Osman MA, Neoh HM, Ab Mutalib NS, Chin SF, Mazlan L, Raja Ali RA, et al. *Parvimonas micra*, *Peptostreptococcus stomatis*, *Fusobacterium nucleatum* and *Akkermansia muciniphila* as a four-bacteria biomarker panel of colorectal cancer. *Sci Rep* 2021;11:2925.
43. Mori G, Rampelli S, Orena BS, Rengucci C, De Maio G, Barbieri G, et al. Shifts of faecal microbiota during sporadic colorectal carcinogenesis. *Sci Rep* 2018;8:10329.
44. Nejman D, Livyatan I, Fuks G, Gavert N, Zwang Y, Geller LT, et al. The human tumor microbiome is composed of tumor type-specific intracellular bacteria. *Science* 2020;368:973–80.

