

SUPPLEMENTARY METHODS

Study Population / Retrieval of Tumor Specimens

Briefly, beginning in 1997 in the HPFS and 2001 in the NHS, we began retrieving, from treating hospital pathology departments, representative pathological specimens from the primary tumor for participants whom we confirmed to have developed colorectal cancer. We successfully obtained specimens for 76% of cases over 16 years of follow-up in HPFS and 58% of cases over 22 years of follow-up in NHS. Tissue sections from all cases were reviewed by a pathologist (S.O.) unaware of other data (11).

Immunohistochemical Assessment and Molecular Assays

Appropriate positive and negative controls were included in every run of immunohistochemistry. For CTSB staining, antigen retrieval was performed, and deparaffinized tissue sections in Antigen Retrieval Citra Solution (Biogenex Laboratories, San Ramon, CA) were treated with microwave in a pressure cooker (5 min). Tissue sections were incubated with 5% normal rabbit serum (Vector Laboratories, Burlingame, CA) in phosphate-buffered saline (30 min). Primary antibody against CTSB [Goat polyclonal anti-CTSB (S-12), 1:1250 dilution; Santa Cruz Biotechnology, Inc., Santa Cruz, CA] was applied, and the slides were maintained at room temperature for 15min, followed by goat secondary antibody (Vector Laboratories) (30 min), an avidin–biotin complex conjugate (Vector Laboratories) (30 min), diaminobenzidine (5 min) and methyl-green counterstain. A pathologist (Y.B.), blinded to any other participant data, recorded cytoplasmic CTSB expression as absent, weak, moderate or strong expression with the

percentage of positive tumor cells. In our previous data from murine models, we demonstrated the ability of NIRF cathepsin-specific molecular agents to identify tumors with weak to strong levels of CTSB expression using immunohistochemistry. Thus, for further analysis in this study, we defined tumors with weak to strong cytoplasmic expression of CTSB as CTSB-positive and tumors with absent cytoplasmic expression of CTSB as CTSB-negative. Methods of immunohistochemistry for p53 and cyclooxygenase-2 (COX-2) were previously described (11, 12).

For methylation analyses, bisulfite DNA treatment and real-time PCR (MethyLight) were validated and performed (12, 13, 17-20). We quantified DNA methylation in 8 CpG island methylator phenotype (CIMP)-specific promoters [*CACNA1G*, *CDKN2A* (p16), *CRABP1*, *IGF2*, *MLH1*, *NEUROG1*, *RUNX3* and *SOCS1*]. CIMP-high was defined as the presence of $\geq 6/8$ methylated promoters, CIMP-low as 1/8-5/8 methylated promoters, and CIMP-0 as the absence (0/8) of methylated promoters, according to the previously established criteria. In order to accurately quantify relatively high methylation levels in long interspersed nuclear elements-1 (LINE-1) repetitive elements, we utilized Pyrosequencing as previously described (14, 21).

We extracted genomic DNA from tumors and performed PCR and Pyrosequencing targeted for *KRAS* (codons 12 and 13), *BRAF* (codon 600) and *PIK3CA* (exons 9 and 20) (13, 15, 16, 20). MSI status was determined using D2S123, D5S346, D17S250, BAT25, BAT26, BAT40, D18S55, D18S56, D18S67 and D18S487 (12, 13). MSI-high was defined as the presence of instability in $\geq 30\%$ of the markers, MSI-low as instability in 1-29% of the markers, and microsatellite stable (MSS) as the absence of unstable markers.

Statistical Analysis

Our multivariate models initially included sex, age at diagnosis (continuous), body mass index (BMI, <30 vs. ≥ 30 kg/m²), family history of colorectal cancer in any first degree relative (present vs. absent), tumor location (proximal vs. distal), tumor stage (I-II vs. III-IV), tumor grade (low vs. high), mucinous component (0% vs. $>0\%$), signet ring cell component (0% vs. $>0\%$), CIMP (high vs. low/0), MSI (high vs. low/MSS), LINE-1 methylation (continuous), *BRAF*, *KRAS*, *PIK3CA*, p53, and COX-2, using patients (N=409) with available molecular data. Backward stepwise elimination with a threshold of $p=0.20$ was used to select variables in the final model. For cases with missing information in any of molecular markers [tumor grade (0.3%), CIMP (1.4%), MSI (0.7%), *BRAF* (1.2%), *KRAS* (0.7%), p53 (0.7%), and COX-2 (0.5%)], we included those cases in a majority category of that missing variable in the initial model. After the selection was done, we assigned separate missing indicator variables to those cases with missing information in any of the categorical covariates in the final model. We confirmed that excluding cases with missing information in any of the molecular markers did not substantially alter results (data not shown). For our Cox proportional hazards model, proportionality of hazards assumption was satisfied by evaluating time-dependent variables, which were the cross product of the CTSB variable and survival time ($p=0.62$ for colon cancer-specific mortality; $p=0.67$ for overall mortality). Statistical interaction was assessed by using the Wald test to assess significance of the cross product of CTSB variable and another variable of interest.