

## Identification of Genetic Susceptibility Loci for Colorectal Tumors in a Genome-Wide Meta-analysis

ULRIKE PETERS,<sup>1,2,\*</sup> SHUO JIAO,<sup>1,\*</sup> FREDRICK R. SCHUMACHER,<sup>3,\*</sup> CAROLYN M. HUTTER,<sup>1,2,\*</sup> AARON K. ARAGAKI,<sup>1</sup> JOHN A. BARON,<sup>4</sup> SONJA I. BERNDT,<sup>5</sup> STÉPHANE BÉZIEAU,<sup>6</sup> HERMANN BRENNER,<sup>7</sup> KATJA BUTTERBACH,<sup>7</sup> BETTE J. CAAN,<sup>8</sup> PETER T. CAMPBELL,<sup>9</sup> CHRISTOPHER S. CARLSON,<sup>1,2</sup> GRAHAM CASEY,<sup>3</sup> ANDREW T. CHAN,<sup>10,11</sup> JENNY CHANG-CLAUDE,<sup>12</sup> STEPHEN J. CHANOCK,<sup>5</sup> LIN S. CHEN,<sup>13</sup> GERHARD A. COETZEE,<sup>3</sup> SIMON G. COETZEE,<sup>3</sup> DAVID V. CONTI,<sup>3</sup> KEITH R. CURTIS,<sup>1</sup> DAVID DUGGAN,<sup>14</sup> TODD EDWARDS,<sup>15</sup> CHARLES S. FUCHS,<sup>11,16</sup> STEVEN GALLINGER,<sup>17</sup> EDWARD L. GIOVANNUCCI,<sup>11,18</sup> STEPHANIE M. GOGARTEN,<sup>19</sup> STEPHEN B. GRUBER,<sup>3</sup> ROBERT W. HAILE,<sup>3</sup> TABITHA A. HARRISON,<sup>1</sup> RICHARD B. HAYES,<sup>20</sup> BRIAN E. HENDERSON,<sup>3</sup> MICHAEL HOFFMEISTER,<sup>7</sup> JOHN L. HOPPER,<sup>21</sup> THOMAS J. HUDSON,<sup>22,23</sup> DAVID J. HUNTER,<sup>18</sup> REBECCA D. JACKSON,<sup>24</sup> SUN HA JEE,<sup>25</sup> MARK A. JENKINS,<sup>21</sup> WEI-HUA JIA,<sup>26</sup> LAURENCE N. KOLONEL,<sup>27</sup> CHARLES KOOPERBERG,<sup>1</sup> SÉBASTIEN KÜRY,<sup>6</sup> ANDREA Z. LACROIX,<sup>1</sup> CATHY C. LAURIE,<sup>19</sup> CECELIA A. LAURIE,<sup>19</sup> LOIC LE MARCHAND,<sup>27</sup> MATHIEU LEMIRE,<sup>22</sup> DAVID LEVINE,<sup>19</sup> NORALANE M. LINDOR,<sup>28</sup> YAN LIU,<sup>29</sup> JING MA,<sup>11</sup> KAREN W. MAKAR,<sup>1</sup> KEITARO MATSUO,<sup>30</sup> POLLY A. NEWCOMB,<sup>1,2</sup> JOHN D. POTTER,<sup>1,31</sup> ROSS L. PRENTICE,<sup>1</sup> CONGHUI QU,<sup>1</sup> THOMAS ROHAN,<sup>32</sup> STEPHANIE A. ROSSE,<sup>1,2</sup> ROBERT E. SCHOEN,<sup>33</sup> DANIELA SEMINARA,<sup>34</sup> MARTHA SHRUBSOLE,<sup>15</sup> XIAO-OU SHU,<sup>15</sup> MARTHA L. SLATTERY,<sup>35</sup> DARIN TAVERNA,<sup>14</sup> STEPHEN N. THIBODEAU,<sup>36</sup> CORNELIA M. ULRICH,<sup>1,2,37</sup> EMILY WHITE,<sup>1,2</sup> YONGBING XIANG,<sup>38</sup> BRENT W. ZANKE,<sup>39</sup> YI-XIN ZENG,<sup>26</sup> BEN ZHANG,<sup>15</sup> WEI ZHENG,<sup>40</sup> and LI HSU,<sup>1</sup> on behalf of the Colon Cancer Family Registry and the Genetics and Epidemiology of Colorectal Cancer Consortium

<sup>1</sup>Public Health Sciences Division, Fred Hutchinson Cancer Research Center, Seattle, Washington; <sup>2</sup>Department of Biostatistics, <sup>3</sup>School of Public Health, University of Washington, Seattle, Washington; <sup>4</sup>Keck School of Medicine, University of Southern California, Los Angeles, California; <sup>5</sup>Department of Medicine, School of Medicine, University of North Carolina, Chapel Hill, North Carolina; <sup>6</sup>Division of Cancer Epidemiology and Genetics, <sup>34</sup>Division of Cancer Control and Population Sciences, National Cancer Institute, Bethesda, Maryland; <sup>7</sup>Service de Génétique Médicale, CHU Nantes, Nantes, France; <sup>8</sup>Division of Clinical Epidemiology and Aging Research, <sup>12</sup>Division of Cancer Epidemiology, German Cancer Research Center, Heidelberg, Germany; <sup>9</sup>Division of Research, Kaiser Permanente Medical Care Program, Oakland, California; <sup>10</sup>American Cancer Society, Atlanta, Georgia; <sup>11</sup>Division of Gastroenterology, Massachusetts General Hospital and Harvard Medical School, Boston, Massachusetts; <sup>13</sup>Channing Division of Network Medicine, Brigham and Women's Hospital and Harvard Medical School, Boston, Massachusetts; <sup>14</sup>Department of Health Studies, University of Chicago, Chicago, Illinois; <sup>15</sup>Translational Genomics Research Institute, Phoenix, Arizona; <sup>16</sup>Department of Medicine, Vanderbilt University School of Medicine, Nashville, Tennessee; <sup>17</sup>Department of Medical Oncology, Dana Farber Cancer Institute, Boston, Massachusetts; <sup>18</sup>Department of Surgery, Toronto General Hospital, Toronto, Ontario, Canada; <sup>19</sup>School of Public Health, Harvard University, Boston, Massachusetts; <sup>20</sup>Division of Epidemiology, New York University School of Medicine, New York, New York; <sup>21</sup>Melbourne School of Population Health, University of Melbourne, Melbourne, Victoria, Australia; <sup>22</sup>Ontario Institute for Cancer Research, Toronto, Ontario, Canada; <sup>23</sup>Departments of Medical Biophysics and Molecular Genetics, University of Toronto, Toronto, Ontario, Canada; <sup>24</sup>Division of Endocrinology, Diabetes, and Metabolism, Ohio State University, Columbus, Ohio; <sup>25</sup>Institute for Health Promotion, Yonsei University, Seoul, Korea; <sup>26</sup>Cancer Center, Sun Yat-sen University, Guangzhou, China; <sup>27</sup>Epidemiology Program, University of Hawaii Cancer Center, Honolulu, Hawaii; <sup>28</sup>Department of Health Sciences Research, Mayo Clinic, Scottsdale, Arizona; <sup>29</sup>Stephens and Associates, Carrollton, Texas; <sup>30</sup>Division of Epidemiology and Prevention, Aichi Cancer Center Research Institute, Nagoya, Japan; <sup>31</sup>Centre for Public Health Research, Massey University, Wellington, New Zealand; <sup>32</sup>Department of Epidemiology and Population Health, Albert Einstein College of Medicine, Bronx, New York; <sup>33</sup>Department of Epidemiology, University of Pittsburgh Medical Center, Pittsburgh, Pennsylvania; <sup>35</sup>Department of Internal Medicine, University of Utah Health Sciences Center, Salt Lake City, Utah; <sup>36</sup>Departments of Laboratory Medicine and Pathology and Laboratory Genetics, Mayo Clinic, Rochester, Minnesota; <sup>37</sup>Division of Preventive Oncology, National Center for Tumor Diseases and German Cancer Research Center, Heidelberg, Germany; <sup>38</sup>Department of Epidemiology, Shanghai Cancer Institute, Shanghai, China; <sup>39</sup>Division of Hematology, Faculty of Medicine, The University of Ottawa, Ottawa, Ontario, Canada; and <sup>40</sup>Department of Medicine, Vanderbilt Epidemiology Center, Vanderbilt University School of Medicine, Nashville, Tennessee

See Covering the Cover synopsis on page 665.

**BACKGROUND & AIMS:** Heritable factors contribute to the development of colorectal cancer. Identifying the genetic loci associated with colorectal tumor formation could elucidate the mechanisms of pathogenesis. **METHODS:** We conducted a genome-wide association study that included 14 studies, 12,696 cases of colorectal tumors (11,870 cancer, 826 adenoma), and 15,113 controls of European descent. The 10 most statistically significant, previously unreported findings were followed up in 6 studies; these included 3056 colorectal tumor cases (2098 cancer, 958 adenoma) and 6658 controls of European and Asian descent. **RESULTS:** Based on the combined analysis, we identified a locus that reached the conventional genome-wide significance level at less than  $5.0 \times 10^{-8}$ : an intergenic region on chromosome 2q32.3, close to *nucleic*

*acid binding protein 1* (most significant single nucleotide polymorphism: rs11903757; odds ratio [OR], 1.15 per risk allele;  $P = 3.7 \times 10^{-8}$ ). We also found evidence for 3 additional loci with  $P$  values less than  $5.0 \times 10^{-7}$ : a locus within the *laminin gamma 1* gene on chromosome 1q25.3 (rs10911251; OR, 1.10 per risk allele;  $P = 9.5 \times 10^{-8}$ ), a

\*Authors share co-first authorship.

**Abbreviations used in this paper:** CCFR, Colon Cancer Family Registry; CCND, cyclin D2; CEU, Utah residents with Northern and Western European ancestry from the CEPH collection; CEPH, Centre d'étude du polymorphisme humain; GECCO, Genetics and Epidemiology of Colorectal Cancer Consortium; GWAS, genome-wide association study; HWE, Hardy Weinberg Equilibrium; LAMC1, laminin gamma 1; MAF, minor allele frequency; NABP1, nucleic acid binding protein 1; ORs, odds ratios; QC, quality control; SNP, single nucleotide polymorphism; TBX3, T-box 3.

© 2013 by the AGA Institute

0016-5085/\$36.00

<http://dx.doi.org/10.1053/j.gastro.2012.12.020>

locus within the *cyclin D2* gene on chromosome 12p13.32 (rs3217810 per risk allele; OR, 0.84;  $P = 5.9 \times 10^{-8}$ ), and a locus in the *T-box 3* gene on chromosome 12q24.21 (rs59336; OR, 0.91 per risk allele;  $P = 3.7 \times 10^{-7}$ ). **CONCLUSIONS:** In a large genome-wide association study, we associated polymorphisms close to *nucleic acid binding protein 1* (which encodes a DNA-binding protein involved in DNA repair) with colorectal tumor risk. We also provided evidence for an association between colorectal tumor risk and polymorphisms in *laminin gamma 1* (this is the second gene in the laminin family to be associated with colorectal cancers), *cyclin D2* (which encodes for cyclin D2), and *T-box 3* (which encodes a T-box transcription factor and is a target of Wnt signaling to  $\beta$ -catenin). The roles of these genes and their products in cancer pathogenesis warrant further investigation.

**Keywords:** Colon Cancer; Genetics; Risk Factors; SNP.

Colorectal cancer has a sizable heritable component; a large twin study estimated that 35% of colorectal cancer risk may be explained by heritable factors.<sup>1</sup> Over the past several years, genome-wide association studies (GWAS), which focus on common single-nucleotide polymorphisms (SNPs), successfully have discovered low-penetrance loci for colorectal cancer.<sup>2-12</sup> These analyses have highlighted genes within the known *transforming growth factor- $\beta$*  and *Wnt* signaling pathways (eg, *bone morphogenetic protein 2 & 4*, *SMAD7*), as well as regions and genes not previously strongly implicated in colorectal cancer (eg, *zinc finger protein 90*, *laminin alpha 5*, *disco-interacting protein 2*), thereby highlighting pathways previously not understood to be involved in colorectal carcinogenesis.<sup>2-12</sup>

To identify additional common genetic risk factors for colorectal tumors, we conducted a genome-wide scan across 14 independent studies including nearly 28,000 subjects and follow-up evaluation of nearly 10,000 independent subjects. We included both colorectal cancer cases and colorectal adenoma cases. Colorectal adenoma is a well-defined colorectal cancer precursor<sup>13</sup> and the majority of colorectal cancers develop through the adenoma-cancer sequence.<sup>14</sup> It has been estimated that the 10-year cumulative rate for advanced adenoma to transition to colorectal cancer is between 10% and 45%, depending on age and sex.<sup>13,15,16</sup> Accordingly, the 2 phenotypes have overlapping etiology.<sup>17</sup> Inclusion of adenoma cases can increase sample size, and hence statistical power, to identify genetic risk factors related to early events in the adenoma-carcinoma process, during which risk factor intervention strategies may offer the greatest potential benefit for cancer prevention.

## Materials and Methods

### Study Participants

Each study is described in detail in the Supplementary Materials and Methods section and the number of cases and

controls as well as age and sex distributions are listed in Supplementary Table 1. In brief, colorectal cancer cases were defined as colorectal adenocarcinoma and confirmed by medical records, pathologic reports, or death certificate. Colorectal adenoma cases were confirmed by medical records, histopathology, or pathologic reports. Controls for adenoma cases had a negative colonoscopy (except for the Nurses' Health Study and the Health Professionals Follow-up Study controls matched to cases with distal adenoma, which either had a negative sigmoidoscopy or colonoscopy examination). All participants provided written informed consent and studies were approved by their respective institution's Institutional Review Boards.

### Genotyping

**GWAS in the Genetics and Epidemiology of Colorectal Cancer Consortium and the Colon Cancer Family Registry.** We conducted a meta-analysis of GWAS from 13 studies within the Genetics and Epidemiology of Colorectal Cancer Consortium (GECCO) (10,729 cases and 13,328 controls) and additional GWAS within the Colon Cancer Family Registry (CCFR) (1967 cases and 1785 controls). Details on genotyping, quality assurance/quality control, and imputation can be found in the Supplementary Materials and Methods section. Average sample and SNP call rates and concordance rates for blinded duplicates are listed in Supplementary Table 2. In brief, all analyses were restricted to European ancestry. Genotyped SNPs were excluded based on call rate ( $<98\%$ ), lack of Hardy-Weinberg equilibrium in controls (HWE,  $P < 1 \times 10^{-4}$ ), and low minor allele frequency (MAF). Because imputation of genotypes is established as standard practice in the analysis of genotype array data, we imputed the autosomal SNPs of all studies to the Utah residents with Northern and Western European ancestry from the Centre d'étude du polymorphisme humain (CEPH) collection (CEU) population in HapMap II (available at: <http://hapmap.ncbi.nlm.nih.gov/>). Imputed SNPs were restricted based on MAF ( $\geq 1\%$ ) and imputation accuracy ( $R^2 > 0.3$ ). After imputation and quality control (QC), a total of 2,708,280 SNPs were used in the meta-analysis of GECCO studies and CCFR. In our detailed result table (Supplementary Table 3), we list for each SNP the number of studies with directly genotyped or imputed data and the mean imputation  $R^2$ . These data show, as expected, that imputed SNPs tend to show very similar results as SNPs that were directly genotyped if the correlation is high between SNPs.

**Follow-up studies.** We selected the 10 most statistically significant regions (excluding known GWAS loci) based on the  $P$  value from the GECCO and CCFR meta-analysis for further follow-up evaluation in colorectal cancer studies in the Asian colorectal cancer consortium and a US-based colorectal adenoma study. Details on genotyping, quality assurance/quality control, and imputation can be found in the Supplementary Materials and Methods section. After quality control exclusions, 2098 colorectal cancer cases and 5749 controls, and 958 colorectal adenoma cases and 909 controls remained in the analysis.

### Statistical Analysis

**GWAS in GECCO and CCFR.** Statistical analyses of the GECCO and CCFR samples were conducted centrally at the coordinating center on individual-level data to ensure a consistent analytic approach. For each study, we estimated the association between SNPs and risk for colorectal cancer by calculating  $\beta$  values, odds ratios (ORs), standard errors, 95% confidence intervals, and  $P$  values using logistic regression models with

log-additive genetic effects. Each directly genotyped SNP was coded as 0, 1, or 2 copies of the risk allele. For imputed SNPs, we used the expected number of copies of the risk allele (the dosage), which has been shown to provide unbiased estimates in the association test for imputed SNPs.<sup>18</sup> We adjusted for age, sex (when appropriate), center (when appropriate), smoking status (Physicians' Health Study only), batch effects (The french Association STudy Evaluating RISK for sporadic colorectal cancer), and the first 3 principal components from EIGENSTRAT (available at: <http://genepath.med.harvard.edu/~reich/EIGENSTRAT.htm>) to account for population substructure. Because CCFR set 2 is a family-based study, we used a conditional logistic regression stratified by family identification while adjusting for age and sex. When analyzing genotyped SNPs on the X chromosome we need to account for different genotype variances between males and females. Therefore, we used the 1 degree of freedom modified Cochran-Armitage test<sup>19</sup> to test for associations. This method has been shown to have robust and powerful performance across a wide range of scenarios.<sup>20</sup> We used logistic regression to model SNP  $\times$  SNP interaction effects for a log-additive model, in which the interaction term is the product of the 2 SNPs.

Quantile-quantile plots were assessed to determine whether the distribution of the *P* values in each study was consistent with the null distribution (except for the extreme tail). We also calculated the genomic inflation factor ( $\lambda$ ) to measure the overdispersion of the test statistics from the association tests by dividing the median of the squared *Z* statistics by 0.455, the median of a chi-squared distribution with 1 degree of freedom. The inflation factor  $\lambda$  was between 0.999 and 1.044 for individual studies based on all SNPs including both directly genotyped and imputed, indicating there is little evidence of residual population substructure, cryptic relatedness, or differential genotyping between cases and controls. This result was consistent with the visual inspection of the study-specific quantile-quantile plots.

We conducted inverse-variance weighted, fixed-effects meta-analysis to combine  $\beta$  estimates and standard errors across individual studies. In this approach, we weighed the  $\beta$  estimate of each study by its inverse variance and calculated a combined estimate by summing the weighted  $\beta$  estimates and dividing by the summed weights. For imputed SNPs, it has been shown that the inverse variance is approximately proportional to the imputation quality.<sup>18</sup> Thus, the inverse variance weighting scheme automatically incorporates imputation quality in the meta-analysis for imputed SNPs. We calculated the heterogeneity *P* values based on Cochran's *Q* statistic<sup>21</sup> and investigated sources for heterogeneity if the *P* value was less than .05 for the 10 most significant SNPs. For the most significant SNPs highlighted in this article, we also examined recessive and unrestricted genetic models and compared models by calculating the Akaike information criterion. We used PLINK (available at: <http://pngu.mgh.harvard.edu/~purcell/plink/>)<sup>22</sup> and R (available at: <http://www.r-project.org/>)<sup>23</sup> to conduct the statistical analysis and summarized results graphically using LocusZOOM (available at: <http://csg.sph.umich.edu/locuszoom/>).<sup>24</sup>

**Follow-up studies.** The 10 most significant SNPs from the GWAS meta-analysis described earlier were analyzed in the follow-up studies (*P* values from GWAS meta-analysis  $2.5 \times 10^{-7}$  to  $6.5 \times 10^{-6}$ ). For the Asian colorectal cancer follow-up study, genotyped SNPs and dosage data of imputed SNPs were analyzed using the program mach2dat (available at: <http://www.sph.umich.edu/csg/abecasis/MACH/download/>).<sup>25</sup> The association between SNP and colorectal cancer risk was assessed

using logistic regression with log-additive genetic effects after adjusting for age and sex. Meta-analyses were performed using the inverse-variance method based on a fixed-effects model, and calculations were implemented in the METAL package (available at: <http://www.sph.umich.edu/csg/abecasis/metal/>).<sup>26</sup> Because the MAF in the Asian follow-up population was very low for the locus on chromosome 14q23.1 (MAF, 0–0.006 in Han Chinese individuals from Beijing, China), we excluded this SNP from the follow-up evaluation in the Asian studies. Given potential differences in the linkage disequilibrium structure between European and Asian descent subjects, we also included all SNPs correlated with these 10 selected SNPs ( $r^2 > 0.5$  in CEU).

For the adenoma follow-up study (all European descent), the association between each genetic marker and risk for colorectal adenoma was estimated by calculating ORs and 95% confidence intervals, using a log-additive genetic model. SNPTTESTv2.2.0 (available at: [https://mathgen.stats.ox.ac.uk/genetics\\_software/snptest/snptest.html](https://mathgen.stats.ox.ac.uk/genetics_software/snptest/snptest.html)) with the “-method score” option<sup>27</sup> was used for logistic regression with the frequentist test, and the model was adjusted for age and sex.

For a combined analysis of GWAS and follow-up results, we conducted inverse-variance weighted fixed-effects meta-analysis to combine ORs from log-additive models across individual studies and measured heterogeneity using Cochran's *Q* statistic, as discussed earlier.

**Criterion for genome-wide significance.** Based on an increasing number of articles<sup>28–33</sup> providing a detailed discussion on the appropriate genome-wide significance threshold, which all arrive at similar values in the range of  $5 \times 10^{-7}$  to  $5 \times 10^{-8}$  for European populations, we decided to use a *P* value of  $5 \times 10^{-8}$  as the genome-wide significance threshold. In addition, we reported on SNPs with *P* values between less than  $5 \times 10^{-7}$  and greater than  $5 \times 10^{-8}$  as a potentially novel SNP that merited additional follow-up evaluation.

**Heritability estimates.** We estimated the additive heritability of colorectal cancer explained by all genotyped SNPs using the method by Yang et al<sup>34</sup> and implemented in the Genome-wide Complex Trait Analysis tool.<sup>35</sup> We set the prevalence of colorectal cancer to 0.004, based on Surveillance, Epidemiology and End Results incidence and National Center for Health Statistics mortality statistics.<sup>36</sup> We used all genotyped SNPs of Darmkrebs: Chancen der Verhütung durch Screening set II and Diet, Activity, and Lifestyle Study set I given the sizable sample set, different genotyping platforms, and inclusion of both sexes (Supplementary Table 1). We also estimated the heritability of previously and newly identified variants by using the method of So et al.<sup>37</sup> Furthermore, we used the method described by Park et al<sup>38</sup> to estimate the total number of loci expected to be identified for colorectal cancer based on the observed effect sizes and power for identifying the loci known to date (Table 1 and Supplementary Table 4).

**Functional annotation of findings.** We conducted a functional annotation for all tagging SNPs (and correlated SNPs) highlighted in this article. As detailed in the Supplementary Materials and Methods section, we queried multiple bioinformatic databases based on the University of California, Santa Cruz genome browser.

## Results

Summary results of the GWAS meta-analysis of GECCO and CCFR are shown in the Manhattan plot (Supplementary Figure 1). Several of the previously iden-

**Table 1.** Risk Estimates for Newly Identified SNPs Associated With Colorectal Cancer at a  $P$  Value Less Than  $5 \times 10^{-7}$ 

SNP	Chromosome (gene) <sup>a</sup>	Risk allele	Stage <sup>b</sup>	RAF (range) <sup>c</sup>	OR (95% CI)	$P$ value	$P$ heterogeneity
SNP with $P < 5 \times 10^{-8}$							
rs11903757	2q32.3 ( <i>NABP1</i> )	C	GWAS	0.16 (0.11–0.23)	1.15 (1.09–1.22)	1.38E-06	
			Asian	0.05 (0.04–0.08)	1.16 (0.95–1.41)	1.34E-01	
			Adenoma	0.15	1.23 (1.03–1.47)	2.30E-02	
			Overall		1.16 (1.10–1.22)	<b>3.71E-08</b>	.27
SNPs with $P < 5 \times 10^{-7}$ and $P > 5 \times 10^{-8}$							
rs10911251	1q25.3 ( <i>LAMC1</i> )	A	GWAS	0.57 (0.49–0.63)	1.10 (1.06–1.14)	1.34E-06	
			Asian	0.54 (0.50–0.55)	1.09 (1.01–1.17)	3.20E-02	
			Adenoma	0.58	1.06 (0.93–1.21)	3.66E-01	
			Overall		1.09 (1.06–1.13)	<b>9.45E-08</b>	.69
rs3217810	12p13.32 ( <i>CCND2</i> )	T	GWAS	0.16 (0.18–0.10)	1.19 (1.11–1.28)	3.40E-07	
			Asian	NA	NA	NA	
			Adenoma	0.15	1.31 (1.00–1.71)	5.07E-02	
			Overall		1.20 (1.12–1.28)	<b>5.86E-08</b>	.91
rs3217901	12p13.32 ( <i>CCND2</i> )	G	GWAS	0.41 (0.43–0.39)	1.10 (1.06–1.15)	1.71E-06	
			Asian	0.56 (0.54–0.58)	1.08 (0.98–1.18)	1.04E-01	
			Adenoma	0.42	1.08 (0.91–1.27)	3.72E-01	
			Overall		1.10 (1.06–1.14)	<b>3.31E-07</b>	.51
rs59336	12q24.21 ( <i>TBX3</i> )	T	GWAS	0.48 (0.51–0.42)	1.10 (1.06–1.14)	7.64E-07	
			Asian	0.60 (0.56–0.62)	1.05 (0.94–1.18)	3.61E-01	
			Adenoma	0.48	1.13 (0.90–1.40)	2.89E-01	
			Overall		1.09 (1.06–1.13)	<b>3.67E-07</b>	.39

CI, confidence interval; RAF, risk allele frequency.

NOTE. Bolded entries signify combined results.

<sup>a</sup>Chromosome position and build from Genome Browser based on build 37.

<sup>b</sup>GWAS ( $n = 12,696$  cases and  $15,113$  controls); Asian ( $2098$  cases and  $5749$  controls); adenoma ( $958$  cases and  $909$  controls); overall ( $15,752$  cases and  $21,771$  controls except for rs3217810, which has  $13,654$  cases and  $16,022$  controls).

<sup>c</sup>Mean and (range) computed from respective studies.

tified GWAS SNPs were highly significantly associated with colorectal cancer, and overall we found a nominal significant association ( $P < .05$ ) in the same direction for 16 of 18 previously identified GWAS loci (Supplementary Table 4). After excluding previously identified regions, we followed up the 10 most significant regions from the GWAS meta-analysis ( $P = 2.5 \times 10^{-7}$  to  $6.5 \times 10^{-6}$ ; Supplementary Table 3). In 4 regions the follow-up studies showed evidence of replication with the association in the same directions as the GWAS and an overall improved significance level (Table 1). Of these 4 regions, 1 region reached the conventional genome-wide significance level at a  $P$  value less than  $5.0 \times 10^{-8}$  in the combined analysis (GWAS + follow-up evaluation). This region was on chromosome 2q32.3 (rs11903757: OR, 1.16 per risk allele;  $P = 3.7 \times 10^{-8}$ ; Table 1 and Supplementary Figure 2). The SNP showed no evidence for heterogeneity ( $P = .27$ ) across all studies. The SNP was correlated strongly ( $r^2 > 0.9$ ) with several SNPs in the same region, which showed similar results (Supplementary Figure 3 and Supplementary Table 3).

The other 3 regions had  $P$  values less than  $5.0 \times 10^{-7}$  (and  $P > 5.0 \times 10^{-8}$ ) in the combined analysis (GWAS + follow-up evaluation). Reporting by chromosomal position, the first of these 3 regions was on chromosome

1q25.3. In this region, the association with rs10911251 had the lowest  $P$  value (OR, 1.09 per risk allele;  $P = 9.5 \times 10^{-8}$ ; Table 1 and Supplementary Figure 2), showing no evidence of heterogeneity ( $P = .69$ ) across studies. This was correlated strongly with a large number of SNPs in the same region showing similar allele frequencies, risk estimates, and  $P$  values spanning across the entire laminin gamma 1 (*LAMC1*) gene (Supplementary Figure 3 and Supplementary Table 3).

The second region with  $P$  values less than  $5.0 \times 10^{-7}$  and greater than  $5.0 \times 10^{-8}$  was on chromosome 12p13, within the *cyclin D2* (*CCND2*) gene. The most statistically significant SNP was rs3217810 (OR, 1.20 per risk allele;  $P = 5.9 \times 10^{-8}$ ; Table 1 and Supplementary Figure 2). Furthermore, only 17.1 kb apart resides a second SNP, rs3217901, which was not strongly correlated with rs3217810 ( $r^2 = 0.052$ – $0.063$ ) and showed a slightly lower significance level (OR, 1.10 per risk allele;  $P = 4.9 \times 10^{-7}$ ). Although the risk allele frequency of rs3217810 in our European descent studies was on average 0.16, this SNP is very uncommon in Asian populations (0.03 in Japanese in Tokyo, Japan, and 0.01 in Han Chinese individuals from Beijing, China) and, hence, the follow-up evaluation of rs3217810 did not include the Asian cases and controls. Both SNPs were not heterogeneous across studies ( $P$  for

heterogeneity = .51 and .91). When we included both SNPs simultaneously in the logistic regression analysis the significance of both SNPs was reduced (Supplementary Table 5).

The third region with  $P$  values less than  $5.0 \times 10^{-7}$  was in the *T-box 3* (*TBX3*) gene on chromosome 12q24.21. The most statistically significant SNP in this region was rs59336 (OR, 1.09 per risk allele;  $P = 3.7 \times 10^{-7}$ ; Table 1 and Supplementary Figure 2). Again, we observed no evidence for heterogeneity across studies ( $P = .39$ ).

We investigated if the 4 regions listed earlier might be more significant (lower  $P$  value) under a different genetic model than the log-additive model. None of the variants was more significant when we modeled the unrestricted, dominant, or recessive mode of inheritance (Supplementary Table 6).

When we stratified results by colorectal adenoma and cancer we observed stronger associations for adenoma compared with cancer for rs11903757 at 2q32.3, similar associations for rs3217810 and rs3217901 at 12p13/*CCDN2* and for rs59336 at 12q24.21/*TBX3*, and a weaker association for rs10911251 at 1q25.3/*LAMC1* (Supplementary Table 7). For previously identified loci, in particular, associations for rs16892766 at 8q23.3/*EIF3H* and rs4939827 at 18q21/*SMAD7* tended to be stronger for adenoma, whereas associations for other loci tended to be similar or weaker compared with cancer (Supplementary Table 4).

We observed no evidence for interaction between the SNPs in the newly identified regions or with SNPs in previously identified regions. The smallest  $P$  value for interaction was .017 for rs59336/*TBX3* and rs11632715/15q13 and was not significant after accounting for multiple comparisons.

As popularized by Yang et al,<sup>34</sup> we estimated that the additive heritability of colorectal cancer explained by all genotyped SNPs would be 14.2% (standard error, 8.2%). The newly identified loci (Table 1) and previously identified loci (Supplementary Table 4) explained about 11% of the additive heritability and cumulatively these newly and previously identified loci explain 1.6% of the variation of colorectal cancer. Based on the study by Park et al<sup>38</sup> we estimated that the total number of loci expected to be identified for colorectal cancer would be between 239 and 500 if the type I error rate was between  $5 \times 10^{-7}$  and  $5 \times 10^{-8}$ .

## Discussion

In this large genome-wide scan meta-analysis and follow-up evaluation of a total of close to 38,000 subjects, we identified an intergenic region on chromosome 2q32.3 close to *nucleic acid binding protein 1* (*NABP1*) that was associated with colorectal tumor risk with  $P$  values less than  $5.0 \times 10^{-8}$ , the conventional genome-wide significance level. Furthermore, we identified 3 regions with  $P$  values less than  $5.0 \times 10^{-7}$ : one on chromosome 1q31 in *LAMC1*, a second on chromosome 12p13 in *CCND2*, and

a third on chromosome 12q24.21 in *TBX3*. All showed highly significant associations with  $P$  values less than  $5 \times 10^{-7}$ .

Our study provides strong support for an intergenic locus on chromosome 2q32.3. The most significant SNPs in this region are in closest proximity to the *NABP1* gene (44 kb centromeric) and the gene *serum deprivation response* (112 kb telomeric), which encodes for the serum-deprivation response phosphatidylserine-binding protein. The SNPs are downstream of *NABP1*, which also is known as *human single-strand DNA binding protein 2* or *oligonucleotide/oligosaccharide binding fold-containing protein 2A*. This protein binds single-stranded DNA via the oligonucleotide/oligosaccharide binding fold domain.<sup>39</sup> Single-stranded DNA binding proteins are important for diverse DNA processes, such as DNA replication, recombination, transcription, and repair.<sup>40–42</sup> Cells depleted of *NABP1* show hypersensitivity to DNA-damaging reagents; *NABP1* participates in repair of DNA double-strand breaks and *ataxia telangiectasia mutated*-dependent signaling pathways,<sup>43</sup> similar to the role of its homolog, *NABP2* (which is also known as human single-strand DNA binding protein 1).<sup>39</sup> Although our functional annotation did not provide further insights on the function of the SNPs, the biologic data described earlier support the importance of *NABP1* with respect to genomic stability, which could explain a link to the development of cancer.<sup>44</sup>

In addition to the genome-wide significant region we observed 3 regions that were slightly less significant with  $P$  values less than  $5 \times 10^{-7}$  but greater than  $5 \times 10^{-8}$ . As has been shown previously,<sup>45</sup> a large fraction of SNPs with borderline genome-wide-significant associations replicated when results from additional studies were added, suggesting that further follow-up evaluation of these regions is warranted. The first of these 3 regions was on chromosome 1q31 and included correlated SNPs showing associations that spanned across the *LAMC1* gene. Interestingly, previous genome-wide scans of colorectal cancer identified a different laminin gene on chromosome 20q13.33, *laminin alpha 5*, as associated with colorectal cancer,<sup>9,11</sup> supporting the importance of this gene family for the development of colorectal cancer. Laminins are extracellular matrix glycoproteins that constitute a major component of the basement membrane in most tissues<sup>46</sup> and in the colon are part of the intestinal epithelial barrier. Laminins are involved in a wide variety of biological functions, such as regulation of cell adhesion, differentiation, migration, signaling, and metastasis.<sup>47–50</sup> Loss of cell-surface laminin anchoring has been found in many cancer cells, particularly those with aggressive subtypes.<sup>51</sup>

*LAMC1* is a large gene spanning 122 kb and containing 28 coding exons. rs10911251 is correlated strongly ( $r^2 > 0.8$ ) with several other SNPs across the gene (Supplementary Figure 3 and Supplementary Table 3). Upon functional annotation, we identified a potential functional candidate (rs10911205) that is correlated strongly with the most significant tagSNP ( $r^2 = 0.73$ ) and located 72 kb upstream within the first intron of *LAMC1*. As shown in

the University of California, Santa Cruz Genome Browser view (Supplementary Figure 4), rs10911205 is located within a highly evolutionarily conserved region and, given its close proximity to the promoter, it is possible that this region influences gene transcription. In addition, the patterns of histone modifications and DNase signals indicating accessibility for transcription factors suggest that this variant may affect cell-type-specific enhancer activity. In summary, given the statistical evidence, support from functional annotation, and evidence from a previous GWAS that identified another laminin gene to be associated with colorectal cancer, we believe there is strong support for the importance of *LAMC1* in the development of colorectal cancer. It is of note that the biologic role of this gene family has not yet been studied substantially in relation to colorectal cancer, supporting the novelty of this finding.

A second region with  $P$  values less than  $5 \times 10^{-7}$  was on chromosome 12p13.32, with 2 independent SNPs both located in the intron of *CCND2*, which belongs to the highly conserved cyclin family, specifically encoding for the protein cyclin D2. Through regulation of CDK4 and CDK6, cyclin D2 affects the cell-cycle transition of the G1/S phase.<sup>52,53</sup> Furthermore, cyclin D2 interacts with tumor-suppressor protein retinoblastoma. Recent studies have identified *CCND2* as a microRNA target gene in different colorectal cancer cell lines.<sup>54,55</sup> Interestingly, genetic variants in *CCND1* also have been related to colorectal cancer<sup>56,57</sup> and a previous GWAS identified a SNP in *CCND1* to be associated with breast cancer.<sup>58</sup>

The third region with  $P$  values less than  $5 \times 10^{-7}$  was identified within the *TBX3* gene, which encodes the T-box transcription factor. *TBX3* is overexpressed in several cancers, including pancreas, liver, breast cancer, and melanoma,<sup>59</sup> playing multiple roles in normal development and cancer.<sup>60</sup> In liver cancer, *TBX3* was identified as a downstream target of the Wnt/ $\beta$ -catenin pathway, mediating  $\beta$ -catenin activities on cell proliferation and survival.<sup>61</sup> The Wnt/ $\beta$ -catenin pathway plays a key role in colorectal cancer development.<sup>62</sup> *TBX5*, another member of the T-Box gene family, has been suggested as an epigenetically inactivated tumor-suppressor gene in colon cancer<sup>63</sup> and provides an additional mechanism by which this gene family may influence colorectal cancer development.

Our study adds further support for all, except 3, previously identified GWAS loci for colorectal cancer. The 3 SNPs (on chromosomes 1q41, 3q26.2, and 6p21) that did not replicate are among the more recently identified GWAS loci<sup>9,12</sup> and have smaller effect sizes (OR for risk allele,  $\leq 1.1$ ) compared with the earlier GWAS findings. As a result, larger sample sizes may be needed to fully replicate these SNPs. Furthermore, it is possible that the effect varies by environmental exposures, which may differ among the study populations. Overall, effect sizes from our study for previous GWAS loci tend to be weaker than in the initial reports, which may be explained by the fact that previous results were subject to the “winner’s curse.”<sup>64</sup>

The large sample size of our GWAS and follow-up studies and availability of individual-level GWAS data are

important advantages of our study. However, the study also had limitations. To increase the sample size, we included Asian descent subjects, who may have different linkage disequilibrium patterns, and the SNPs analyzed may be tagging different underlying causal variants. To address this potential limitation we included all SNPs correlated with the most significant SNPs, which likely will identify any variant that genuinely is associated with colorectal cancer risk across different ancestral groups, as shown for other GWAS loci.<sup>65–68</sup> Given that genotyping platforms only capture a subset of the genome, we used imputation to HapMap II to obtain a better coverage of the common variation across the genome and to generate a common set of SNPs from the different platforms. Because imputed SNPs tend to result in less significant findings depending on their imputation accuracy,<sup>69</sup> we expect that our results provide relative conservative significance levels.<sup>70</sup> Similar to previous GWAS,<sup>2,4,6–10,12</sup> we included colorectal adenoma as the major precursor of colorectal cancer to improve our statistical power and to identify genetic variants that act early in the adenoma-cancer sequence, where adenomas and cancer have a shared etiology. Although the inclusion of adenoma also may add heterogeneity because adenomas will not show an association for genetic variants that act later in the carcinogenic process (ie, on progression from adenoma to cancer) or for variants that act through adenoma-independent pathways, stratified analysis may provide insights into the mediating roles of genes within the normal to adenoma to cancer pathway. We show that for some of the newly and previously identified loci, associations are stronger for adenomas compared with cancer; however, we observed similar or weaker associations for other loci. These results may suggest that some genes are important in early stages of cancer development while others may be more important for the progression from adenoma to cancer. However, given the relatively small number of adenoma cases (only 6.5% of the GWAS and 31% of the follow-up cases were adenoma cases), it is important that our findings are replicated in studies with larger numbers of adenoma cases.

In summary, in this large study, we identified one novel susceptibility locus associated with the risk of colorectal tumor on chromosome 2q32.3 close to *NABP1*, and 3 potential loci with borderline genome-wide significant results within *LAMC1*, *CCND2*, and *TBX3*. These findings are supported by biologic plausibility, functional annotation, and previous GWAS findings within the same gene family, emphasizing the potential relevance of these genes in the etiology of colorectal cancer.

### Supplementary Material

Note: To access the supplementary material accompanying this article, visit the online version of *Gastroenterology* at [www.gastrojournal.org](http://www.gastrojournal.org), and at <http://dx.doi:10.1053/j.gastro.2012.12.020>.

## References

1. Lichtenstein P, Holm NV, Verkasalo PK, et al. Environmental and heritable factors in the causation of cancer—analyses of cohorts of twins from Sweden, Denmark, and Finland. *N Engl J Med* 2000;343:78–85.
2. Tomlinson I, Webb E, Carvajal-Carmona L, et al. A genome-wide association scan of tag SNPs identifies a susceptibility variant for colorectal cancer at 8q24.21. *Nat Genet* 2007;39:984–988.
3. Zanke BW, Greenwood CM, Rangrej J, et al. Genome-wide association scan identifies a colorectal cancer susceptibility locus on chromosome 8q24. *Nat Genet* 2007;39:989–994.
4. Broderick P, Carvajal-Carmona L, Pittman AM, et al. A genome-wide association study shows that common alleles of SMAD7 influence colorectal cancer risk. *Nat Genet* 2007;39:1315–1317.
5. Tenesa A, Farrington SM, Prendergast JG, et al. Genome-wide association scan identifies a colorectal cancer susceptibility locus on 11q23 and replicates risk loci at 8q24 and 18q21. *Nat Genet* 2008;40:631–637.
6. Jaeger E, Webb E, Howarth K, et al. Common genetic variants at the CRAC1 (HMPS) locus on chromosome 15q13.3 influence colorectal cancer risk. *Nat Genet* 2008;40:26–28.
7. Tomlinson IP, Webb E, Carvajal-Carmona L, et al. A genome-wide association study identifies colorectal cancer susceptibility loci on chromosomes 10p14 and 8q23.3. *Nat Genet* 2008;40:623–630.
8. Houlston RS, Webb E, Broderick P, et al. Meta-analysis of genome-wide association data identifies four new susceptibility loci for colorectal cancer. *Nat Genet* 2008;40:1426–1435.
9. Houlston RS, Cheadle J, Dobbins SE, et al. Meta-analysis of three genome-wide association studies identifies susceptibility loci for colorectal cancer at 1q41, 3q26.2, 12q13.13 and 20q13.33. *Nat Genet* 2010;42:973–977.
10. Tomlinson IP, Carvajal-Carmona LG, Dobbins SE, et al. Multiple common susceptibility variants near BMP pathway loci GREM1, BMP4, and BMP2 explain part of the missing heritability of colorectal cancer. *PLoS Genet* 2011;7:e1002105.
11. **Peters U, Hutter CM**, Hsu L, et al. Meta-analysis of new genome-wide association studies of colorectal cancer risk. *Hum Genet* 2012;131:217–234.
12. Dunlop MG, Dobbins SE, Farrington SM, et al. Common variation near CDKN1A, POLD3 and SHROOM2 influences colorectal cancer risk. *Nat Genet* 2012;44:770–776.
13. Brenner H, Hoffmeister M, Stegmaier C, et al. Risk of progression of advanced adenomas to colorectal cancer by age and sex: estimates based on 840,149 screening colonoscopies. *Gut* 2007;56:1585–1589.
14. Kinzler KW, Vogelstein B. Lessons from hereditary colorectal cancer. *Cell* 1996;87:159–170.
15. Eide TJ. Risk of colorectal cancer in adenoma-bearing individuals within a defined population. *Int J Cancer* 1986;38:173–176.
16. Stryker SJ, Wolff BG, Culp CE, et al. Natural history of untreated colonic polyps. *Gastroenterology* 1987;93:1009–1013.
17. Potter JD. Colorectal cancer: molecules and populations. *J Natl Cancer Inst* 1999;91:916–932.
18. Jiao S, Hsu L, Hutter CM, et al. The use of imputed values in the meta-analysis of genome-wide association studies. *Genet Epidemiol* 2011;35:597–605.
19. Clayton D. Testing for association on the X chromosome. *Biostatistics* 2008;9:593–600.
20. Hickey PF, Bahlo M. X chromosome association testing in genome wide association studies. *Genet Epidemiol* 2011;35:664–670.
21. Ioannidis JP, Patsopoulos NA, Evangelou E. Heterogeneity in meta-analyses of genome-wide association investigations. *PLoS ONE* 2007;2:e841.
22. Purcell S, Neale B, Todd-Brown K, et al. PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am J Hum Genet* 2007;81:559–575.
23. R Development Core Team. R: a language and environment for statistical computing. Vienna, Austria: R Foundation for Statistical Computing, 2011.
24. Pruim RJ, Welch RP, Sanna S, et al. LocusZoom: regional visualization of genome-wide association scan results. *Bioinformatics* 2010;26:2336–2337.
25. Li Y, Willer CJ, Ding J, et al. MaCH: using sequence and genotype data to estimate haplotypes and unobserved genotypes. *Genet Epidemiol* 2010;34:816–834.
26. Willer CJ, Li Y, Abecasis GR. METAL: fast and efficient meta-analysis of genomewide association scans. *Bioinformatics* 2010;26:2190–2191.
27. Marchini J, Howie B, Myers S, et al. A new multipoint method for genome-wide association studies by imputation of genotypes. *Nat Genet* 2007;39:906–913.
28. Risch N, Merikangas K. The future of genetic studies of complex human diseases. *Science* 1996;273:1516–1517.
29. International HapMap Consortium. A haplotype map of the human genome. *Nature* 2005;437:1299–1320.
30. Wellcome Trust Case Control Consortium. Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls. *Nature* 2007;447:661–678.
31. Hoggart CJ, Clark TG, De IM, et al. Genome-wide significance for dense SNP and resequencing data. *Genet Epidemiol* 2008;32:179–185.
32. Pe'er I, Yelensk R, Altshuler D, et al. Estimation of the multiple testing burden for genomewide association studies of nearly all common variants. *Genet Epidemiol* 2008;32:381–385.
33. Dudbridge F, Gusnanto A. Estimation of significance thresholds for genomewide association scans. *Genet Epidemiol* 2008;32:227–234.
34. Yang J, Benyamin B, McEvoy BP, et al. Common SNPs explain a large proportion of the heritability for human height. *Nat Genet* 2010;42:565–569.
35. Yang J, Lee SH, Goddard ME, et al. GCTA: a tool for genome-wide complex trait analysis. *Am J Hum Genet* 2011;88:76–82.
36. Howlander N, Noone AM, Krapcho M, et al. SEER cancer statistics review, 1975-2009. Bethesda, MD: National Cancer Institute, 2012.
37. So HC, Gui AH, Cherny SS, et al. Evaluating the heritability explained by known susceptibility variants: a survey of ten complex diseases. *Genet Epidemiol* 2011;35:310–317.
38. Park JH, Wacholder S, Gail MH, et al. Estimation of effect size distribution from genome-wide association studies and implications for future discoveries. *Nat Genet* 2010;42:570–575.
39. Richard DJ, Bolderson E, Cubeddu L, et al. Single-stranded DNA-binding protein hSSB1 is critical for genomic stability. *Nature* 2008;453:677–681.
40. Bochkarev A, Bochkareva E, Frappier L, et al. The crystal structure of the complex of replication protein A subunits RPA32 and RPA14 reveals a mechanism for single-stranded DNA binding. *EMBO J* 1999;18:4498–4504.
41. Wold MS. Replication protein A: a heterotrimeric, single-stranded DNA-binding protein required for eukaryotic DNA metabolism. *Annu Rev Biochem* 1997;66:61–92.
42. Yang H, Jeffrey PD, Miller J, et al. BRCA2 function in DNA binding and recombination from a BRCA2-DSS1-ssDNA structure. *Science* 2002;297:1837–1848.
43. Li Y, Bolderson E, Kumar R, et al. HSSB1 and hSSB2 form similar multiprotein complexes that participate in DNA damage response. *J Biol Chem* 2009;284:23525–23531.
44. Broderick S, Rehmet K, Concannon C, et al. Eukaryotic single-stranded DNA binding proteins: central factors in genome stability. *Subcell Biochem* 2010;50:143–163.
45. Panagiotou OA, Ioannidis JP. What should the genome-wide significance threshold be? Empirical replication of borderline genetic associations. *Int J Epidemiol* 2012;41:273–286.
46. Kalluri R. Basement membranes: structure, assembly and role in tumour angiogenesis. *Nat Rev Cancer* 2003;3:422–433.
47. Turck N, Gross I, Gendry P, et al. Laminin isoforms: biological roles and effects on the intracellular distribution of nuclear pro-

- teins in intestinal epithelial cells. *Exp Cell Res* 2005;303:494–503.
48. Pouliot N, Saunders NA, Kaur P. Laminin 10/11: an alternative adhesive ligand for epidermal keratinocytes with a functional role in promoting proliferation and migration. *Exp Dermatol* 2002;11:387–397.
  49. Gudjonsson T, Ronnov-Jessen L, Villadsen R, et al. Normal and tumor-derived myoepithelial cells differ in their ability to interact with luminal breast epithelial cells for polarity and basement membrane deposition. *J Cell Sci* 2002;115:39–50.
  50. Patarroyo M, Tryggvason K, Virtanen I. Laminin isoforms in tumor invasion, angiogenesis and metastasis. *Semin Cancer Biol* 2002;12:197–207.
  51. Akhavan A, Griffith OL, Soroceanu L, et al. Loss of cell-surface laminin anchoring promotes tumor growth and is associated with poor clinical outcomes. *Cancer Res* 2012;72:2578–2588.
  52. Lukas J, Muller H, Bartkova J, et al. DNA tumor virus oncoproteins and retinoblastoma gene mutations share the ability to relieve the cell's requirement for cyclin D1 function in G1. *J Cell Biol* 1994;125:625–638.
  53. Matsushime H, Quelle DE, Shurtleff SA, et al. D-type cyclin-dependent kinase activity in mammalian cells. *Mol Cell Biol* 1994;14:2066–2076.
  54. Ragusa M, Statello L, Maugeri M, et al. Specific alterations of the microRNA transcriptome and global network structure in colorectal cancer after treatment with MAPK/ERK inhibitors. *J Mol Med (Berl)* 2012;90:1421–1438.
  55. Zhang P, Ma Y, Wang F, et al. Comprehensive gene and microRNA expression profiling reveals the crucial role of hsa-let-7i and its target genes in colorectal cancer metastasis. *Mol Biol Rep* 2012;39:1471–1478.
  56. Yang Y, Wang F, Shi C, et al. Cyclin D1 G870A polymorphism contributes to colorectal cancer susceptibility: evidence from a systematic review of 22 case-control studies. *PLoS One* 2012;7:e36813.
  57. Yang J, Zhang G, Chen J. CCND1 G870A polymorphism is associated with increased risk of colorectal cancer, especially for sporadic colorectal cancer and in Caucasians: a meta-analysis. *Clin Res Hepatol Gastroenterol* 2012;36:169–177.
  58. Turnbull C, Rapley EA, Seal S, et al. Variants near DMRT1, TERT and ATF7IP are associated with testicular germ cell cancer. *Nat Genet* 2010;42:604–607.
  59. Smith J, Mowla S, Prince S. Basal transcription of the human TBX3 gene, a key developmental regulator which is overexpressed in several cancers, requires functional NF-Y and Sp1 sites. *Gene* 2011;486:41–46.
  60. Washkowitz AJ, Gavrillov S, Begum S, et al. Diverse functional networks of Tbx3 in development and disease. *Wiley Interdiscip Rev Syst Biol Med* 2012;4:273–283.
  61. Renard CA, Labalette C, Armengol C, et al. Tbx3 is a downstream target of the Wnt/beta-catenin pathway and a critical mediator of beta-catenin survival functions in liver cancer. *Cancer Res* 2007;67:901–910.
  62. Morin PJ, Sparks AB, Korinek V, et al. Activation of beta-catenin-Tcf signaling in colon cancer by mutations in beta-catenin or APC. *Science* 1997;275:1787–1790.
  63. Yu J, Ma X, Cheung KF, et al. Epigenetic inactivation of T-box transcription factor 5, a novel tumor suppressor gene, is associated with colon cancer. *Oncogene* 2010;29:6464–6474.
  64. Garner C. Upward bias in odds ratio estimates from genome-wide association studies. *Genet Epidemiol* 2007;31:288–295.
  65. Setiawan VW, Haessler J, Schumacher F, et al. HNF1B and endometrial cancer risk: results from the PAGE study. *PLoS One* 2012;7:e30390.
  66. Lindstrom S, Schumacher FR, Campa D, et al. Replication of five prostate cancer loci identified in an Asian population—results from the NCI Breast and Prostate Cancer Cohort Consortium (BPC3). *Cancer Epidemiol Biomarkers Prev* 2012;21:212–216.
  67. Dumitrescu L, Carty CL, Taylor K, et al. Genetic determinants of lipid traits in diverse populations from the population architecture using genomics and epidemiology (PAGE) study. *PLoS Genet* 2011;7:e1002138.
  68. Buyske S, Wu Y, Carty CL, et al. Evaluation of the metabochip genotyping array in African Americans and implications for fine mapping of GWAS-identified loci: the PAGE study. *PLoS One* 2012;7:e35651.
  69. Pritchard JK, Przeworski M. Linkage disequilibrium in humans: models and data. *Am J Hum Genet* 2001;69:1–14.
  70. Newton-Cheh CN, Eijgelsheim M, Rice KM, et al. Common variants at ten loci influence QT interval duration in the QTGEN Study. *Nat Genet* 2009;41:399–406.

---

Author names in bold designate shared co-first authorship.  
Received August 10, 2012. Accepted December 14, 2012.

#### Reprint requests

Address requests for reprints to: Ulrike Peters, PhD, MPH, Fred Hutchinson Cancer Research Center, 1100 Fairview Avenue North, M4-B402, PO Box 19024, Seattle, Washington 98109-1024. e-mail: [upeters@fhcrc.org](mailto:upeters@fhcrc.org); fax: (206) 667-7850; or Li Hsu, PhD, Fred Hutchinson Cancer Research Center, 1100 Fairview Avenue North, M2-B500, Seattle, Washington 98109-1024. e-mail: [lih@fhcrc.org](mailto:lih@fhcrc.org); fax: (206) 667-7004.

#### Acknowledgments

The authors wish to thank the following:

**Asian Consortium:** The authors wish to thank the study participants and research staff for their contributions and commitment to this project, Regina Courtney for DNA preparation, and Jing He for data processing and analyses.

**The french Association Study Evaluating RISK for sporadic colorectal cancer:** The authors are very grateful to Dr Bruno Buecher without whom this project would not have existed. The authors also thank all those who agreed to participate in this study, including the patients and the healthy control persons, as well as all the physicians, technicians, and students.

**Darmkrebs: Chancen der Verhütung durch Screening:** The authors thank all participants and cooperating clinicians, and Ute Handte-Daub, Renate Hettler-Jensen, Utz Benschaid, Muhabbet Celik, and Ursula Eilber for excellent technical assistance.

**GECCO:** The authors would like to thank all those at the GECCO Coordinating Center for helping to bring together the data and people who made this project possible.

**Health Professionals Follow-up Study, Nurses' Health Study, and Physicians' Health Study:** The authors would like to acknowledge Patrice Soule and Hardeep Ranu of the Dana Farber Harvard Cancer Center High-Throughput Polymorphism Core who assisted in the genotyping for Nurses' Health Study, Health Professionals Follow-up Study, and Physician's Health Study under the supervision of Dr Immaculata Devivo and Dr David Hunter, Qin (Carolyn) Guo, and Lixue Zhu who assisted in programming for Nurses' Health Study and Health Professionals Follow-up Study, and Haiyan Zhang who assisted in programming for the Physicians' Health Study. The authors would like to thank the participants and staff of the Nurses' Health Study and the Health Professionals Follow-up Study for their valuable contributions as well as the following state cancer registries for their help: AL, AZ, AR, CA, CO, CT, DE, FL, GA, IA, ID, IL, IN, KY, LA, MA, MD, ME, MI, NC, ND, NE, NH, NJ, NY, OH, OK, OR, PA, RI, SC, TN, TX, VA, WA, and WY.

**Prostate, Lung, Colorectal Cancer, and Ovarian Cancer Screening Trial:** The authors thank Drs Christine Berg and Philip Prorok, Division of Cancer Prevention, National Cancer Institute, the Screening Center investigators and staff of the Prostate, Lung, Colorectal Cancer, and Ovarian Cancer Screening Trial, Mr Tom Riley and staff of Information Management Services, Inc, Ms Barbara O'Brien and staff of Westat, Inc, and Drs Bill Kopp, Wen Shao, and staff of SAIC-Frederick. Most importantly, the authors acknowledge

the study participants for their contributions to making this study possible.

**Postmenopausal Hormone study:** The authors would like to thank the study participants and staff of the Hormones and Colon Cancer study.

**Tennessee Colorectal Polyp Study:** The authors thank the study participants and the research staff for their contributions and commitment to this project, and Regina Courtney for DNA preparation.

**Women's Health Initiative:** The authors thank the Women's Health Initiative investigators and staff for their dedication, and the study participants for making the program possible. A full listing of Women's Health Initiative investigators can be found at: <https://cleo.whi.org/researchers/Documents%20%20Write%20a%20Paper/WHI%20Investigator%20Short%20List.pdf>.

#### Conflicts of interest

The authors disclose no conflicts.

#### Funding

The Genetics and Epidemiology of Colorectal Cancer Consortium study was supported by the National Cancer Institute, National Institutes of Health, and the US Department of Health and Human Services (U01 CA137088; R01 CA059045). The Asian Consortium was supported by a Grant-in-aid for Cancer Research, the Grant for the Third Term Comprehensive Control Research for Cancer and Grants-in-Aid for Scientific Research from the Japanese Ministry of Education, Culture, Sports, Science and Technology (17015018 and 221S0001). The french Association Study Evaluating RISK for sporadic colorectal cancer was supported by a Hospital Clinical Research Program (PHRC) and by the Regional Council of Pays de la Loire, the Groupement des Entreprises Françaises dans la Lutte Contre le Cancer, the Association Anne de Bretagne Génétique, and the Ligue Régionale Contre le Cancer. The Assessment of Risk in Colorectal Tumours in Canada study was supported by the National Institutes of Health through funding allocated to the Ontario Registry for Studies of Familial Colorectal Cancer (U01 CA074783; see the Colon Cancer Family Registry support section below); and by a GL2 grant from the Ontario Research Fund, the Canadian Institutes of Health Research, by a Cancer Risk Evaluation Program grant from the Canadian Cancer Society Research Institute, and by Senior Investigator Awards (T.J.H. and B.W.Z.) from the Ontario Institute for Cancer Research, through generous support from the Ontario Ministry of Economic Development and Innovation. The Hawaii Colorectal Cancer Studies 2 and 3 studies were supported by the National Institutes of Health (R01 CA60987). The Colon Cancer Family Registry was supported by the National Institutes of Health (RFA CA-95-011) and through cooperative agreements with members of the Colon Cancer Family Registry and P.I.s. This genome-wide scan was supported by the National Cancer Institute, National Institutes of Health (U01 CA122839). The content of this manuscript does not necessarily reflect the views or policies of the National Cancer Institute or any of the collaborating centers in the cancer family registries, nor does mention of trade names, commercial products, or organizations imply endorsement by the US Government or the cancer family registries. The following colon cancer family registries centers contributed data to this article and were supported by National Institutes of Health: the Australasian Colorectal Cancer Family Registry (U01 CA097735), Seattle Colorectal Cancer Family Registry (U01 CA074794), and the Ontario Registry for Studies of Familial Colorectal Cancer (U01 CA074783). The Darmkrebs: Chancen der Verhütung durch Screening study was supported by the German Research Council (Deutsche Forschungsgemeinschaft, BR 1704/6-1, BR 1704/6-3, BR 1704/6-4, and CH 117/1-1), and the German Federal Ministry of Education and Research (01KH0404 and

01ER0814). The Diet, Activity, and Lifestyle Study was supported by the National Institutes of Health (R01 CA48998 to M.L.S.); Guangzhou-1 was supported by the National Key Scientific and Technological Project (2011ZX09307-001-04) and the National Basic Research Program (2011CB504303) was supported by the People's Republic of China. The Health Professionals Follow-up Study was supported by the National Institutes of Health (P01 CA 055075, UM1 CA167552, R01 137178, and P50 CA 127003), the Nurses' Health Study was supported by the National Institutes of Health (R01 137178, P01 CA 087969, and P50 CA 127003), and the Physicians' Health Study was supported by the National Institutes of Health (CA42182). The Korean Cancer Prevention Study-II study was supported by the National R&D Program for cancer control (1220180), and the Seoul R&D Program (10526, Republic of Korea). The Multiethnic Cohort study was supported by the National Institutes of Health (R37 CA54281, P01 CA033619, and R01 CA63464). The Prostate, Lung, Colorectal Cancer, and Ovarian Cancer Screening Trial was supported by the Intramural Research Program of the Division of Cancer Epidemiology and Genetics, and supported by contracts from the Division of Cancer Prevention, National Cancer Institute, National Institutes of Health, Department of Health and Human Services. Control samples were genotyped as part of the Cancer Genetic Markers of Susceptibility prostate cancer scan, supported by the Intramural Research Program of the National Cancer Institute. The data sets used in this analysis were accessed with appropriate approval through the dbGaP online resource ([http://www.cgems.cancer.gov/data\\_access.html](http://www.cgems.cancer.gov/data_access.html)) through dbGaP accession number 000207v.1p1. Control samples also were genotyped as part of the GWAS of Lung Cancer and Smoking (Yeager, M et al. Nat Genet 2008;124:161-170). Support for this work was provided through the National Institutes of Health, Genes, Environment and Health Initiative (Z01 CP 010200). The human subjects participating in the genome-wide association study were derived from the Prostate, Lung, Colon, and Ovarian Screening Trial and the study was supported by intramural resources of the National Cancer Institute. Assistance with genotype cleaning, as well as with general study coordination, was provided by the Gene Environment Association Studies, Geneva Coordinating Center (U01 HG004446). Assistance with data cleaning was provided by the National Center for Biotechnology Information. Funding support for genotyping, which was performed at the Johns Hopkins University Center for Inherited Disease Research, was provided by the National Institutes of Health, Genes, Environment and Health Initiative (U01 HG 004438). The data sets used for the analyses described in this article were obtained from dbGaP at <http://www.ncbi.nlm.nih.gov/gap> through dbGaP accession number phs000093 v2.p2. The Postmenopausal Hormone Study was supported by the National Institutes of Health (R01 CA076366 to P.A.N.). The Shanghai-1 and Shanghai-2 studies were supported by the National Institutes of Health (R37CA070867, R01CA082729, R01CA124558, R01CA148667, and R01CA122364), as well as an Ingram Professorship and Research Reward funds from the Vanderbilt University School of Medicine. The Tennessee Colorectal Polyp Study was supported by the National Institutes of Health (P50CA95103 and R01CA121060) and was conducted by the Survey and Biospecimen Shared Resource, which was supported in part by the Vanderbilt-Ingram Cancer Center (P30 CA 68485). The VITamins And Lifestyle study was supported by the National Institutes of Health (K05 CA154337). The Women's Health Initiative program was funded by the National Heart, Lung, and Blood Institute, National Institutes of Health, US Department of Health and Human Services, through contracts HHSN268201100046C, HHSN268201100001C, HHSN268201100002C, HHSN268201100003C, HHSN268201100004C, and HHSN271201100004C.

## Supplementary Materials and Methods

### *Study Populations Included in GWAS and Follow-up Studies*

**GWAS in GECCO and CCFR.** We describe each study population used in the GWAS. For information on sample sizes and demographic factors please see Supplementary Table 1.

**Ontario Familial Colorectal Cancer Registry.** In GECCO, a subset of the Assessment of Risk in Colorectal Tumours in Canada from the OFCCR (Ontario Registry for Studies of Familial Colorectal Cancer) was used. Both the case-control study<sup>1</sup> and the OFCCR<sup>2</sup> have been described in detail previously, as have the GWAS results.<sup>3</sup> In brief, cases were confirmed incident colorectal cancer cases if they were ages 20 to 74 years, residents of Ontario, identified through comprehensive registry, and diagnosed between July 1997 and June 2000. Population-based controls were selected randomly among Ontario residents (random-digit dialing and listing of all Ontario residents), and matched by sex and 5-year age groups. A total of 1236 colorectal cancer cases and 1223 controls were genotyped successfully on at least one of the following: Illumina 1536 GoldenGate assay (Illumina, Inc, San Diego, CA), the Affymetrix GeneChip Human Mapping 100K and 500K Array Set (Affymetrix, Inc, Santa Clara, CA), or a 10K nonsynonymous SNP chip. Analysis was based on a set of unrelated subjects who were non-Hispanic, white by self-report, or by investigation of genetic ancestry. We further excluded subjects if there was a sample mix-up, if they were missing epidemiologic questionnaire data, if they were cases with a tumor in the appendix, or if they overlapped with the CCFR. In addition, only samples genotyped on the Affymetrix GeneChip 500K Array were used to avoid coverage issues in imputation.

**The french Association STudy Evaluating RISK for sporadic colorectal cancer.** Participants were recruited from the Pays de la Loire region in France between December 2002 and March 2006.<sup>4</sup> Eligibility criteria for cases included being Caucasian, age 40 years or older at diagnosis, and having no family history of colorectal cancer or polyps. Cases were patients with first primary colorectal cancer diagnosed in 1 of the 6 public hospitals and 5 clinics located in the Pays de la Loire region that participated in the study. Cases were confirmed based on medical and pathology reports. Controls were recruited at 2 Health Examination Centers of the Pays de la Loire region, and the recruitment of controls age 70 years and older was completed in the Departments of Internal Medicine and Hepatogastroenterology of the University Hospital Center of Nantes, located in the same region. Controls were eligible to participate if they were Caucasian, age 40 years or older, and had no family history of colorectal cancer or polyps. In the presence of the physician, each participant filled out a standardized

questionnaire on family information, medical history, lifestyle, and dietary intake. Cases and controls provided a blood sample.

**CCFR.** The CCFR is a National Cancer Institute-supported consortium consisting of 6 centers dedicated to the establishment of a comprehensive collaborative infrastructure for interdisciplinary studies in the genetic epidemiology of colorectal cancer.<sup>5</sup> The CCFR includes data from approximately 30,500 total subjects (10,500 probands and 20,000 unaffected and affected relatives and unrelated controls). Cases and controls, age 20–74 years, were recruited at the 6 participating centers beginning in 1998. CCFR implemented a standardized questionnaire that was administered to all participants and included established and suspected risk factors for colorectal cancer, which included questions on medical history and medication use, reproductive history (for female participants), family history, physical activity, demographics, alcohol and tobacco use, and dietary factors. The set 1 scan, which has been described previously,<sup>6</sup> included population-based cases and age-matched controls from the 3 population-based centers: Seattle, Toronto, and Australia. Cases were enriched genetically by oversampling those with a young age at onset or positive family history. Controls were matched to cases on age and sex. All cases and controls were self-reported as white, which was confirmed with genotype data. The set 2 scan included population-based cases and matched controls from all 6 colon CFR centers including the Mayo Clinic, Hawaii Cancer Registry, University of Southern California, Fred Hutchinson Cancer Research Center, Ontario Cancer Care, and University of Melbourne. As with set 1, cases were enriched genetically by oversampling those with a young age at onset or positive family history. Controls were same-generation family controls.

**Darmkrebs: Chancen der Verhütung durch Screening.** This study was initiated as a large population-based, case-control study in 2003 in the Rhine-Neckar-Odenwald region (southwest region of Germany) to assess the potential of endoscopic screening for reduction of colorectal cancer risk and to investigate etiologic determinants of disease, particularly lifestyle/environmental factors and genetic factors.<sup>7,8</sup> Cases with a first diagnosis of invasive colorectal cancer (International Classification of Diseases 10 codes C18-C20) who were at least 30 years of age (no upper age limit), German speaking, a resident in the study region, and mentally and physically able to participate in a 1-hour interview were recruited by their treating physicians either in the hospital a few days after surgery or by mail after discharge from the hospital. Cases were confirmed based on histologic reports and hospital discharge letters after diagnosis of colorectal cancer. All hospitals treating colorectal cancer patients in the study region participated. Based on estimates from population-based cancer registries, more than 50% of all potentially eligible patients with incident

colorectal cancer in the study region were included. Community-based controls were selected randomly from population registries, using frequency matching with respect to age (5-year groups), sex, and county of residence. Controls with a history of colorectal cancer were excluded. Controls were contacted by mail and follow-up telephone calls. The participation rate was 51%. During an in-person interview, data were collected on demographics, medical history, family history of colorectal cancer, and various lifestyle factors, as were blood and mouthwash samples. The set 1 scan consisted of a subset of participants recruited up until 2007, and samples were frequency matched on age and sex. The set 2 scan consisted of additional subjects who were recruited until 2010 as part of this ongoing study.

**Diet, Activity, and Lifestyle Study.** DALIS (Diet, Activity, and Lifestyle Study) was a population-based, case-control study of colon cancer. Participants were recruited between 1991 and 1994 from 3 locations: the Kaiser Permanente Medical Care Program of Northern California, an 8-county area in Utah, and the metropolitan Twin Cities area of Minnesota.<sup>9</sup> Eligibility criteria for cases included age at diagnosis between 30 and 79 years; diagnosis with first primary colon cancer (International Classification of Diseases for Oncology second edition codes 18.0 and 18.2–18.9) between October 1, 1991, and September 30, 1994; English speaking; and competency to complete the interview. Individuals with cancer of the rectosigmoid junction or rectum were excluded, as were those with a pathology report noting familial adenomatous polyposis, Crohn's disease, or ulcerative colitis. A rapid-reporting system was used to identify all incident cases of colon cancer, resulting in the majority of cases being interviewed within 4 months of diagnosis. Controls from the Kaiser Permanente Medical Care Program were selected randomly from membership lists. In Utah, controls younger than 65 years of age were selected randomly through random-digit dialing and driver's license lists. Controls, 65 years of age and older, were selected randomly from Health Care Financing Administration lists. In Minnesota, controls were identified from Minnesota driver's licenses or state identification lists. Controls were matched to cases by 5-year age groups and sex. The set 1 scan consisted of a subset of the study designed earlier, from Utah, Minnesota, and the Kaiser Permanente Medical Care Program, and was restricted to subjects who self-reported as white non-Hispanic. The set 2 scan consisted of subjects from Utah and Minnesota who were not genotyped in set 1. Set 2 was restricted to subjects who self-reported as white non-Hispanic and those who had appropriate consent to post data to dbGaP.

**Hawaii Colorectal Cancer Studies 2 and 3.** Patients with colorectal cancer were identified through the rapid reporting system of the Hawaii Surveillance, Epidemiology and End Results registry and consisted of all Japanese, Caucasian, and native Hawaiian residents of

Oahu who were newly diagnosed with an adenocarcinoma of the colon or rectum between January 1994 and August 1998.<sup>10</sup> Control subjects were selected from participants in an ongoing population-based health survey conducted by the Hawaii State Department of Health and from Health Care Financing Administration participants. Controls were matched to cases by sex, ethnicity, and age (within 2 years). Personal interviews were obtained from 768 matched pairs, resulting in a participation rate of 58.2% for cases and 53.2% for controls. A questionnaire, administered during an in-person interview, included questions about demographics, lifetime history of tobacco, alcohol use, aspirin use, physical activity, personal medical history, family history of colorectal cancer, height and weight, diet (Food Frequency Questionnaire), and postmenopausal hormone use. A blood sample was obtained from 548 (71%) interviewed cases and 662 (86%) interviewed controls. Surveillance, Epidemiology and End Results staging information was extracted from the Hawaii Tumor Registry. In GECCO, self-reported Caucasian subjects with DNA, and clinical and epidemiologic data, were selected for genotyping.

**Health Professionals Follow-up Study.** The HPFS (Health Professionals Follow-up Study) is a parallel prospective study to the NHS (Nurses' Health Study).<sup>11</sup> The HPFS cohort comprised 51,529 men who, in 1986, responded to a mailed questionnaire. The participants were US male dentists, optometrists, osteopaths, podiatrists, pharmacists, and veterinarians born between 1910 and 1946. Participants provided information on health-related exposures, including current and past smoking history, age, weight, height, diet, physical activity, aspirin use, and family history of colorectal cancer. Colorectal cancer and other outcomes were reported by participants or next-of-kin and were followed up through review of the medical and pathology record by physicians. Overall, more than 97% of self-reported colorectal cancers were confirmed by medical record review. Information was abstracted on histology and primary location. Incident cases were defined as those occurring after the subject provided the blood sample. Prevalent cases were defined as those occurring after enrollment in the study but before the subject provided the blood sample. Follow-up evaluation has been excellent, with 94% of the men responding to date. Colorectal cancer cases were ascertained through January 1, 2008. In 1993–1995, 18,825 men in the HPFS mailed blood samples by overnight courier, which were aliquoted into buffy coat and stored in liquid nitrogen. In 2001–2004, 13,956 men in the HPFS who had not provided a blood sample previously mailed in a swish-and-spit sample of buccal cells. Incident cases were defined as those occurring after the subject provided a blood or buccal sample. Prevalent cases were defined as those occurring after enrollment in the study in 1986, but before the subject provided either a blood or buccal sample. After excluding participants

with histories of cancer (except nonmelanoma skin cancer), ulcerative colitis, or familial polyposis, 2 case-control sets were constructed from which DNA was isolated from either buffy coat or buccal cells for genotyping, as follows: (1) a case-control set with cases of colorectal cancer matched to randomly selected controls who provided a blood sample and were free of colorectal cancer at the same time the colorectal cancer was diagnosed in the cases; and (2) a case-control set with cases of colorectal cancer matched to randomly selected controls who provided a buccal sample and were free of colorectal cancer at the same time the colorectal cancer was diagnosed in the case. For both case-control sets, matching criteria included year of birth (within 1 year) and month/year of blood or buccal cell sampling (within 6 months). Cases were pair-matched 1:1, 1:2, or 1:3 with a control participant(s).

In addition to colorectal cancer cases and controls, a set of adenoma cases and matched controls with available DNA from buffy coat were selected for genotyping. Over the follow-up period, data were collected on endoscopic screening practices and, if individuals had been diagnosed with a polyp, the polyps were confirmed to be adenomatous by medical record review. Adenoma cases were ascertained through January 1, 2008. A separate case-control set was constructed of participants diagnosed with advanced adenoma matched to control participants who underwent a lower endoscopy in the same time period and did not have an adenoma. Advanced adenoma was defined as an adenoma 1 cm or larger in diameter and/or with tubulovillous, villous, or high-grade dysplasia/carcinoma-in-situ histology. Matching criteria included year of birth (within 1 year) and month/year of blood sampling (within 6 months), the reason for their lower endoscopy (screening, family history, or symptoms), and the time period of any prior endoscopy (within 2 years). Controls matched to cases with a distal adenoma either had a negative sigmoidoscopy or colonoscopy examination, and controls matched to cases with proximal adenoma all had a negative colonoscopy.

**Multiethnic Cohort study.** The MEC (Multiethnic Cohort) was initiated in 1993 to investigate the impact of dietary and environmental factors on major chronic diseases, particularly cancer, in ethnically diverse populations in Hawaii and California.<sup>12</sup> The study recruited 96,810 men and 118,441 women aged 45–75 years between 1993 and 1996. Incident colorectal cancer cases occurring since January 1995 and controls were contacted for blood or saliva samples. The median interval between diagnosis and blood draw was 14 months (interquartile range, 10–19 mo) among cases and the participation rate was 74%. A sample of cohort participants was selected randomly to serve as controls at the onset of the nested case-control study (participation rate, 66%). The selection was stratified by sex, age, and race/ethnicity. Colorectal cancer cases were identified through the

Rapid Reporting System of the Hawaii Tumor Registry and through quarterly linkage to the Los Angeles County Cancer Surveillance Program. Both registries are members of Surveillance, Epidemiology and End Results. In GECCO, self-reported white subjects from the nested case-control study described earlier with DNA and clinical and epidemiologic data were selected for genotyping.

**Nurses' Health Study.** The NHS cohort began in 1976 when 121,700 married female registered nurses age 30–55 years returned the initial questionnaire that ascertained a variety of important health-related exposures.<sup>13</sup> Since 1976, follow-up questionnaires have been mailed every 2 years. Colorectal cancer and other outcomes were reported by participants or next-of-kin and followed up through review of the medical and pathology record by physicians. Overall, more than 97% of self-reported colorectal cancers were confirmed by medical-record review. Information was abstracted on histology and primary location. The rate of follow-up evaluation has been high: as a proportion of the total possible follow-up time, follow-up evaluation has been more than 92%. Colorectal cancer cases were ascertained through June 1, 2008. In 1989–1990, 32,826 women in NHS I mailed blood samples by overnight courier, which were aliquoted into buffy coat and stored in liquid nitrogen. In 2001–2004, 29,684 women in NHS I who did not previously provide a blood sample mailed a swish-and-spit sample of buccal cells. Incident cases were defined as those occurring after the subject provided a blood or buccal sample. Prevalent cases were defined as those occurring after enrollment in the study in 1976 but before the subject provided either a blood or buccal sample. After excluding participants with histories of cancer (except nonmelanoma skin cancer), ulcerative colitis, or familial polyposis, 2 case-control sets were constructed from which DNA was isolated from either buffy coat or buccal cells for genotyping: (1) a case-control set with cases of colorectal cancer matched to randomly selected controls who provided a blood sample and were free of colorectal cancer at the same time the colorectal cancer was diagnosed in the case; and (2) a case-control set with cases of colorectal cancer matched to randomly selected controls who provided a buccal sample and were free of colorectal cancer at the same time the colorectal cancer was diagnosed in the cases. For both case-control sets, matching criteria included year of birth (within 1 year) and month/year of blood or buccal cell sampling (within 6 months). Cases were pair matched 1:1, 1:2, or 1:3 with a control participant(s).

In addition to colorectal cancer cases and controls, a set of adenoma cases and matched controls with available DNA from buffy coat were selected for genotyping. Over the follow-up period, data were collected on endoscopic screening practices and, if individuals had been diagnosed with a polyp, the polyps were confirmed to be adenomatous by medical record review. Adenoma cases

were ascertained through June 1, 2008. A separate case-control set was constructed of participants diagnosed with advanced adenoma matched to control participants who underwent a lower endoscopy in the same time period and did not have an adenoma. Advanced adenoma was defined as an adenoma more than 1 cm in diameter and/or with tubulovillous, villous, or high-grade dysplasia/carcinoma-in-situ histology. Matching criteria included year of birth (within 1 year) and month/year of blood sampling (within 6 months), the reason for their lower endoscopy (screening, family history, or symptoms), and the time period of any prior endoscopy (within 2 years). Controls matched to cases with a distal adenoma either had a negative sigmoidoscopy or colonoscopy examination, and controls matched to cases with proximal adenoma all had a negative colonoscopy.

**Physicians' Health Study.** The PHS (Physicians' Health Study) was established as a randomized, double-blind, placebo-controlled trial of aspirin and  $\beta$ -carotene among 22,071 healthy US male physicians, between 40 and 84 years of age, in 1982.<sup>14,15</sup> Participants completed 2 mailed questionnaires before being assigned randomly, additional questionnaires at 6 and 12 months, and questionnaires annually thereafter. In addition, participants were sent postcards at 6 months to ascertain status. From August 1982 to December 1984, there were 14,916 baseline blood samples collected from the physicians during the run-in phase before randomization. When participants reported a diagnosis of cancer, medical records and pathology reports were reviewed by study physicians who were blinded to exposure data. Among those who provided baseline blood samples, colorectal cases were ascertained through March 31, 2008, and controls were matched on age (within 1 year for younger participants, up to 5 years for older participants) and smoking status (never, past, current). Cases were pair-matched 1:1, 1:2, or 1:3 with a control participant(s). Because of DNA availability, samples were genotyped in 2 batches on the same platform at the same genotyping center at different time points.

**Prostate, Lung, Colorectal, and Ovarian Cancer Screening Trial.** The PLCO (Prostate, Lung, Colorectal Cancer, and Ovarian Cancer Screening Trial) enrolled 154,934 participants (men and women, aged between 55 and 74 y) at 10 centers into a large, randomized, 2-arm trial to determine the effectiveness of screening to reduce cancer mortality. Sequential blood samples were collected from participants assigned to the screening arm. Participation was 93% at the baseline blood draw. In the observational (control) arm, buccal cells were collected via mail using the swish-and-spit protocol; the participation rate was 65%. Details of this study have been described previously<sup>16,17</sup> and are available online (<http://dcp.cancer.gov/plco>).

The set 1 scan included a subset of 577 colon cancer cases self-reported as being non-Hispanic white with

available DNA samples, questionnaire data, and appropriate consent for ancillary epidemiologic studies. Cases were excluded if they had a history of inflammatory bowel disease, polyps, polyposis syndrome, or cancer (excluding basal or squamous cell skin cancer). Controls originated from the Cancer Genetic Markers of Susceptibility prostate cancer scan<sup>18,19</sup> (all male) and the GWAS of Lung Cancer and Smoking<sup>20</sup> (enriched for smokers), along with an additional 92 non-Hispanic white female controls. For the set 2 scan, cases were individuals with colorectal cancer from both arms of the trial who were not already included in set 1. Samples were excluded if participants did not sign appropriate consent forms, if DNA was unavailable, if baseline questionnaire data with follow-up evaluation were unavailable, if they had a history of colon cancer before the trial, if they had a rare cancer, if they were already in a colon GWAS, or if they were a control in the prostate or lung populations. Controls were frequency-matched 1:1 to cases without replacement, and cases were not eligible to be controls. Matching criteria were age at enrollment (2-year blocks), enrollment date (2-year blocks), sex, race/ethnicity, trial arm, and study year of diagnosis (ie, controls must be cancer free into the case's year of diagnosis).

**Postmenopausal Hormones Supplementary Study to the CCFR.** Eligible case patients included all female residents, ages 50–74 years, residing in the 13 counties in Washington State, reporting to the Cancer Surveillance, Epidemiology and End Results program, who were newly diagnosed with invasive colorectal adenocarcinoma (ICD-O C18.0, C18.2–C18.9, C19.9, C20.0–C20.9) between October 1998 and February 2002.<sup>21</sup> Eligibility for all individuals was limited to those who were English speaking with available telephone numbers, through which they could be contacted. On average, cases were identified within 4 months of diagnosis. The overall response proportion of eligible cases identified was 73%. Community-based controls were selected randomly according to age distribution (in 5-year age intervals) of the eligible cases by using lists of licensed drivers from the Washington State Department of Licensing for individuals, ages 50–64 years, and rosters from the Health Care Financing Administration (now the Centers for Medicare and Medicaid), for individuals older than age 64. The overall response proportion of eligible controls was 66%. In GECCO, samples with sufficient DNA extracted from blood were genotyped. Only participants who were not part of the CCFR Seattle site were included in the sample set.

**VITamins And Lifestyle.** The VITAL (VITamins And Lifestyle) cohort comprised 77,721 Washington State men and women aged 50–76 years, recruited from 2000 to 2002, to investigate the association of supplement use and lifestyle factors with cancer risk. Subjects were recruited by mail, from October 2000 to December 2002, using names purchased from a commercial mailing

list. All subjects completed a 24-page questionnaire and buccal-cell specimens for DNA were self-collected by 70% of the participants. Subjects were followed up for cancer by linkage to the western Washington Surveillance, Epidemiology and End Results cancer registry and were censored when they moved out of the area covered by the registry or at time of death. Details of this study have been described previously.<sup>22</sup> In GECCO, a nested case-control set was genotyped. Samples included colorectal cancer cases with DNA, excluding subjects with colorectal cancer before baseline; in situ cases; (large cell) neuroendocrine carcinoma; squamous cell carcinoma; carcinoid tumor; Goblet-cell carcinoid; and any type of lymphoma, including non-Hodgkin, Mantle cell, large B-cell, or follicular lymphoma. Controls were matched on age at enrollment (within 1 year), enrollment date (within 1 year), sex, and race/ethnicity. One control was selected randomly per case among all controls who matched according to the 4 factors described earlier and for whom the control follow-up time was greater than the follow-up time of the case until diagnosis.

**Women's Health Initiative.** The WHI (Women's Health Initiative) is a long-term health study of 161,808 post-menopausal women aged 50–79 years at 40 clinical centers throughout the United States. WHI comprised a clinical trial arm, an observational study (OS) arm, and several extension studies. The details of WHI have been described previously<sup>23,24</sup> and are available online (<https://cleo.whi.org/SitePages/Home.aspx>). In GECCO, set 1 cases were selected from the September 12, 2005, database and comprised centrally adjudicated colon cancer cases from the OS arm who self-reported as white. Controls were first selected among controls previously genotyped as part of a hip fracture GWAS conducted within the WHI OS arm and matched to cases on age (within 3 years), enrollment date (within 365 days), hysterectomy status, and prevalent conditions at baseline. For 37 cases, there was no control match in the hip fracture GWAS. For these participants, we identified a matched control in the WHI OS arm based on the same criteria. In the set 2 scan, cases were selected from the August 2009 database and comprised centrally adjudicated colon and colorectal cancer cases from the OS and clinical trial arms who were not genotyped in set 1. In addition, case and control participants were subject to the following exclusion criteria: a prior history of colorectal cancer at baseline, institutional review board approval not available for data submission into dbGaP, and insufficient DNA available. Matching criteria included age (within years), race/ethnicity, WHI date (within 3 years), WHI Calcium and Vitamin D study date (within 3 years), and randomization arms (OS flag, hormone therapy assignments, dietary modification assignments, calcium/vitamin D assignments). In addition, they were matched by the 4 regions of randomization centers. Each case was matched with 1 control (1:1) who met the matching criteria ex-

actly. Control selection was performed in a time-forward manner, selecting one control for each case first from the risk set at the time of the case's event. The matching algorithm was allowed to select the closest match based on a criterion to minimize an overall distance measure.<sup>25</sup> Each matching factor was given the same weight. Additional available controls who were genotyped as part of the hip fracture GWAS were included to improve power.

### *Follow-up Studies*

In the following, we describe each study population used in the follow-up study. For information on sample sizes and demographic factors please see Supplementary Table 1.

**Asia Colorectal Cancer Consortium.** The study protocols were approved by relevant institutional review boards at all study sites, and all included subjects provided informed consent. Sample size, genotype platform, the number of SNPs used in imputation, and genomic inflation factors in each of the 5 studies are presented in Supplementary Table 8.

**Shanghai study 1 and 2.** Colorectal cancer cases were derived from the Shanghai Women's Health Study<sup>26</sup> and the Shanghai Men's Health Study,<sup>27</sup> both population-based cohort studies that are being conducted in urban Shanghai, China. A total of 777 pathologically diagnosed colorectal cancer cases with DNA available were identified in participants from the Shanghai Women's Health Study and Shanghai Men's Health Study and included in this study. A total of 758 cancer-free controls were derived from the Shanghai Women's Health Study/Shanghai Men's Health Study and frequency-matched to colorectal cancer cases by age and sex. To increase statistical power, we also included 2131 community female controls who were scanned using the Affymetrix Genome-Wide Human SNP Array 6.0 (Affymetrix 6.0) as part of an ongoing GWAS of breast cancer.<sup>28</sup> A total of 481 cases and 2632 controls were genotyped using Affymetrix 6.0 (Shanghai Study 1). A total of 296 cases and 257 controls were genotyped using Illumina HumanOmniExpress BeadChip (Illumina OmniExpress) (Shanghai Study 2).

**Guangzhou study 1.** This study contributed 694 cases and 972 controls. Histopathologically diagnosed colorectal cancer cases were recruited from the Sun Yat-Sen University Cancer Center between January 2002 and January 2012. Healthy controls were recruited from physical examination centers of several large general hospitals in Guangdong province communities.<sup>29</sup> At enrollment, controls reported no history of any cancer. All cases and controls were self-reported Han Chinese who lived in Guangdong Province at the time of recruitment. Blood samples from all cases and controls were obtained as the source of genomic DNA for the study.

**Aichi study 1.** This study is part of the Hospital-based Epidemiologic Research Program at Aichi Cancer

Center in Japan.<sup>30</sup> All first-visit outpatients 20–79 years of age at the Aichi Cancer Center during December 2000 to November 2005 were asked to participate in the Hospital-based Epidemiologic Research Program at Aichi Cancer Center. Of 29,736 eligible patients approached, 28,766 participated in the study, with a response rate of 96.7%. All participants completed self-administered questionnaires about their lifestyle and demographic characteristics and provided blood samples. Case status was confirmed via the Hospital-based Epidemiologic Research Program at Aichi Cancer Center database and the hospital-based cancer registry database. A total of 589 colorectal cancer cases were identified in this cohort and 497 were included in the GWAS. A total of 942 controls without any cancer at recruitment were selected randomly and frequency-matched to cases by age and sex.<sup>31</sup>

**Korean Cancer Prevention Study-II.** The Korean Cancer Prevention Study-II included 266,258 individuals, 20–77 years of age, who visited 16 health promotion centers nationwide from April 2004 to December 2008 in South Korea.<sup>32</sup> Subjects were interviewed at baseline to obtain exposure data. Cancer diagnoses were identified through 2008 using data from the national cancer registry and hospitalization records. For the study, we selected 325 colorectal cancer patients who provided a blood sample. Cancer-free cohort members ( $N = 977$ ) were selected randomly as controls.

**Tennessee Colorectal Polyp Study.** The Tennessee Colorectal Polyp Study was a colonoscopy-based, case-control study conducted in Nashville, Tennessee, from 2003 to 2011.<sup>33</sup> Eligible participants, aged 40–75 years old, were identified from patients at the Vanderbilt Gastroenterology Clinic and the Veteran's Affairs Tennessee Valley Health System Nashville Campus. Participants were excluded if they had genetic colorectal cancer syndromes, a prior history of inflammatory bowel disease, prevalent adenomatous polyps, or any cancer other than nonmelanoma skin cancer. Colonoscopic procedures were performed and reported using standard clinical protocols and all pathology diagnoses were determined by hospital pathologists. Participants provided DNA either before or after colonoscopy (blood and buccal samples were collected). The analysis included only participants of Caucasian race.

### *Genotyping, Quality Assurance/QC, and Imputation*

**GWAS in GECCO and CCFR.** We conducted a meta-analysis of GWAS from 13 studies within the GECCO consortium (10,729 cases and 13,328 controls) and additional GWAS within the CCFR (1967 cases and 1785 controls). The GWAS from CCFR, which consisted of participants from sites in the United States, Canada, and Australia, included a population-based, case-control set (CCFR set 1, 1171 cases and 983 controls) genotyped using Illumina Human1M or Human1M-Duo,<sup>6</sup> and a

sibling-pair set (CCFR set 2, 796 cases and 802 controls) genotyped using Illumina Omni1. The GECCO GWAS consisted of participants within The french Association Study Evaluating RISK for sporadic colorectal cancer; Hawaiian Colorectal Cancer Studies 2 and 3; DACHS [Darmkrebs: Chancen der Verhütung durch Screening]; DAL5; HPFS; MEC; NHS; OFCCR; PHS; Postmenopausal Hormone Study; PLCO; VITAL study; and the WHI. Phase one genotyping of a total of 1709 colon cancer cases and 4214 controls from PLCO, WHI, and DAL5 (PLCO set 1, WHI set 1, and DAL5 set 1) was performed using Illumina HumanHap 550K, 610K, or combined Illumina 300K and 240K, and has been described previously.<sup>34</sup> A total of 650 colorectal cancer cases and 522 controls from OFCCR were included in GECCO from previous genotyping using Affymetrix platforms.<sup>3</sup> A total of 5540 colorectal cancer cases and 5425 controls from the The french Association Study Evaluating RISK for sporadic colorectal cancer, Hawaiian Colorectal Cancer Studies 2 and 3, DACHS set 1, DAL5 set 2, the MEC, Postmenopausal Hormone study, PLCO set 2, VITAL study, and WHI set 2 were genotyped successfully using Illumina HumanCytoSNP. A total of 2004 colorectal cancer cases and 2244 controls from HPFS, NHS, PHS, and DACHS set 2, as well as a total of 826 advanced adenoma cases and 923 controls from HPFS and NHS were genotyped successfully using Illumina HumanOmniExpress.

DNA was extracted from blood samples or, for a subset of DACHS, HPFS, MEC, NHS, and PLCO samples, and for all VITAL samples, from buccal cells, using conventional methods. All studies included 1%–6% blinded duplicates to monitor the quality of the genotyping. All individual-level genotype data were managed and underwent quality assurance and QC at the University of Southern California (CCFR sets 1 and 2), the OFCCR, the University of Washington Genetics Coordinating Center (HPFS, NHS, PHS, and DACHS set 2), or the GECCO Coordinating Center at the Fred Hutchinson Cancer Research Center (all other studies). Details on the quality assurance/QC can be found in [Supplementary Table 2](#). In brief, samples were excluded based on call rate, heterozygosity, unexpected duplicates, gender discrepancy, and unexpectedly high identity-by-descent or unexpected genotype concordance ( $>65\%$ ) with another individual. All analyses were restricted to samples clustering with the Utah residents with northern and western European ancestry from the CEU population in principal component analysis, including the HapMap II populations as reference. SNPs were excluded if they were triallelic, not assigned an rs number, or were reported or observed as not performing consistently across platforms. In addition, genotyped SNPs were excluded based on call rate ( $<98\%$ ), lack of HWE in controls ( $P < 1 \times 10^{-4}$ ), and MAF ( $<5\%$  in set 1 for PLCO, WHI, DAL5, and OFCCR;  $<5/\text{number of samples}$  for remaining studies).

Because imputation of genotypes is established as standard practice in the analysis of genotype array data, all autosomal SNPs from all studies were imputed to the CEU population in HapMap II release 24, with the exception of OFCCR, which was imputed to HapMap II release 22. CCFR sets 1 and 2 were imputed using IMPUTE (available at: <https://mathgen.stats.ox.ac.uk/impute/impute.html>),<sup>35</sup> OFCCR was imputed using BEAGLE (available at: <http://faculty.washington.edu/browning/beagle/beagle.html>),<sup>36</sup> and all other studies were imputed using MACH (available at: <http://www.sph.umich.edu/csg/abecasis/MACH/tour/>).<sup>37</sup> Imputed data were merged with genotype data such that genotype data preferentially were selected if a SNP had both types of data, unless there was a difference in terms of reference allele frequency ( $>0.1$ ) or position ( $>100$  base pairs), in which case imputed data were used. Given the high agreement of imputation accuracy among MACH, IMPUTE, and BEAGLE,<sup>38</sup> the common practice to use different imputation programs is unlikely to cause heterogeneity<sup>39</sup> and it has become common practice to combine results across SNPs imputed using different programs. As a measurement of imputation accuracy we calculated  $R^2$ . Analyses of imputed data had different QC cut-off values than those for directly genotyped SNPs discussed earlier and were restricted to SNPs with either a MAF of 1% or greater or an  $R^2$  value greater than 0.3, with the exception of CCFR set 2, which was restricted to SNPs with both a MAF of 1% or greater and an  $R^2$  value of 0.3 or greater. After imputation and QC, a total of 2,708,280 SNPs were used in the meta-analysis of the GECCO and CCFR studies.

**Follow-up studies.** We selected the 10 most statistically significant regions (excluding known GWAS loci) based on the  $P$  value from the GECCO and CCFR meta-analyses for follow-up evaluation in colorectal cancer studies in Asian populations and adenoma studies in populations of European descent.

The Asian colorectal cancer follow-up study comprised a meta-analysis of 5 studies conducted in China, Japan, and South Korea, including 2293 colorectal cancer cases and 5780 controls. Cases and controls were genotyped using multiple SNP arrays, including Affymetrix Genome-Wide Human SNP Array 6.0, Affymetrix Genome-Wide Human SNP Array 5.0, Illumina Infinium HumanHap610 BeadChip, Illumina Human610-Quad BeadChip, and Illumina HumanOmniExpress BeadChip. Samples were excluded based on low call rate ( $<95\%$ ), heterozygosity, unexpected duplicates, gender discrepancy, and outlying population substructure. After quality control exclusions, 2098 cases and 5749 controls remained in the analysis. SNPs were excluded for low call rate ( $<95\%$ ), low genotype concordance ( $<95\%$ ) among positive QC samples, an MAF less than 5%, or an HWE  $P$  value less than  $1 \times 10^{-5}$  in controls. For each of the 5 studies, SNPs were imputed for autosomal SNPs that

were present in HapMap Japanese in Tokyo, Japan+Han Chinese individuals from Beijing, China Phase 2 release 22 using MACH.<sup>37</sup> SNPs with an  $R^2$  value greater than 0.5 were included in the analysis.

The colorectal adenoma follow-up study consisted of a US-based GWAS of 1049 cases and 987 controls.<sup>33</sup> DNA extracted from blood and buccal samples were genotyped using the Affymetrix Genome-Wide Human SNP Array 5.0. Samples were excluded based on low call rate ( $<95\%$ ), heterozygosity, unexpected duplicates, gender discrepancy, identity-by-descent, and outlying population substructure. After quality control exclusions, 958 cases and 909 controls remained in the analysis. SNPs were excluded for low call rate ( $<95\%$ ), MAF less than 1%, or HWE  $P$  value less than  $1 \times 10^{-6}$ . After quality control exclusions, a total of 402,326 SNPs remained in the analysis. Data were imputed to the 1000 Genomes Project and HapMap Phase 3 using IMPUTE.<sup>35</sup> SNPs with an  $R^2$  value greater than 0.5 were included in the analysis.

#### *Details on Functional Annotation Findings Using Bioinformatic Databases*

There are several bioinformatic tools available for the post-GWAS functional characterization of putative disease-causing loci through the University of California, Santa Cruz genome browser.<sup>40</sup> Annotation of non-protein-coding regions operates under the hypothesis that trait-associated alleles exert their effects by influencing transcriptional levels through multiple regulatory mechanisms. The University of California, Santa Cruz genome browser provides several tracks that can be used to annotate enhancers, promoters, insulators, and silencers<sup>40</sup> (for details see Supplementary Table 9). Such tools help expedite the discovery of causal variants by isolating a few likely culprits from a large background of variants in linkage disequilibrium with the surrogate marker (tag SNP). Because distal enhancers often facilitate cell-type-specific expression, it is helpful to look for evidence in a variety of cell lines in addition to those related to the trait. For example, the ENCODE (available at: <http://genome.ucsc.edu/ENCODE/>) transcription summary track assayed by RNA-sequencing can be displayed as an overlay of histograms denoting expression levels in various tissues marked by a specific color, thus allowing identification of cell-type-specific expression.

Similarly the histone modification tracks can provide additional evidence for cell-specific regulatory elements when displayed in this configuration. The methylation and acetylation of histone proteins changes chromatin accessibility for transcription and such marks can serve as a powerful tool for identifying both enhancer and promoter regions. There are 3 summary ENCODE tracks available to detect specific chemical modifications and were assayed in 7 different tissues using chromatin immunoprecipitation sequencing methodology. The H3K4me1 histone mark is associated with enhancers

downstream of transcription start sites. The H3k27Ac histone mark is similarly thought to enhance transcription and likely does so through the blocking of the repressive histone mark H3K27Me3. The last histone modification in the summary tracks, H3K4Me3, is associated with active promoters. Additional chemical modifications and cell lines are available under the Broad Institute histone modification track for further interrogation.

Regulatory regions are susceptible to DNase cutting and ENCODE has assayed this hypersensitivity in a large collection of cell types. The precision of the DNase cluster track is somewhat better than that of chromatin modifications. Identification of evolutionarily conserved segments, phylogenetic footprints, has been used to discover functionally important regions. However, histone marks and DNase hypersensitivity tracks are more robust tools for characterizing regulatory regions because these elements are not always constrained across vertebrate evolution. Functional hypotheses around regulatory regions can be strengthened with the ENCODE transcription factor track. By using the chromatin immunoprecipitation sequencing method, this track helps to identify the alteration of transcription factor binding sites, which potentially alter expression levels. As an example, CCCTC-binding factor is a transcription factor that assumes multiple forms and can act as an activator, a repressor/silencer, or an insulator. When binding chromatin insulators, it can prevent interactions between promoters and nearby enhancers or silencers. However, it also mediates long-range chromatin looping, which can bring enhancers in proximity of a gene's promoter. Combining the strengths and weaknesses of each of these tracks can provide *in silico* evidence for regulatory function, and enables selection of strong candidates for additional functional studies using reporter gene methods.

### Supplementary References

- Cotterchio M, Manno M, Klar N, et al. Colorectal screening is associated with reduced colorectal cancer risk: a case-control study within the population-based Ontario Familial Colorectal Cancer Registry. *Cancer Causes Control* 2005;16:865–875.
- Cotterchio M, Keown-Eyssen G, Sutherland H, et al. Ontario familial colon cancer registry: methods and first-year response rates. *Chronic Dis Can* 2000;21:81–86.
- Zanke BW, Greenwood CM, Rangrej J, et al. Genome-wide association scan identifies a colorectal cancer susceptibility locus on chromosome 8q24. *Nat Genet* 2007;39:989–994.
- Küry S, Buecher B, Robiou-du-Pont S, et al. Combinations of cytochrome P450 gene polymorphisms enhancing the risk for sporadic colorectal cancer related to red meat consumption. *Cancer Epidemiol Biomarkers Prev* 2007;16:1460–1467.
- Newcomb PA, Baron J, Cotterchio M, et al. Colon Cancer Family Registry: an international resource for studies of the genetic epidemiology of colon cancer. *Cancer Epidemiol Biomarkers Prev* 2007;16:2331–2343.
- Figueiredo JC, Lewinger JP, Song C, et al. Genotype-environment interactions in microsatellite stable/microsatellite instability-low colorectal cancer: results from a genome-wide association study. *Cancer Epidemiol Biomarkers Prev* 2011;20:758–766.
- Brenner H, Chang-Claude J, Seiler CM, et al. Protection from colorectal cancer after colonoscopy: population-based case-control study. *Ann Intern Med* 2011;154:22–30.
- Lilla C, Verla-Tebit E, Risch A, et al. Effect of NAT1 and NAT2 genetic polymorphisms on colorectal cancer risk associated with exposure to tobacco smoke and meat consumption. *Cancer Epidemiol Biomarkers Prev* 2006;15:99–107.
- Slattery ML, Potter J, Caan B, et al. Energy balance and colon cancer—beyond physical activity. *Cancer Res* 1997;57:75–80.
- Le Marchand L, Hankin JH, Wilkens LR, et al. Combined effects of well-done red meat, smoking, and rapid N-acetyltransferase 2 and CYP1A2 phenotypes in increasing colorectal cancer risk. *Cancer Epidemiol Biomarkers Prev* 2001;10:1259–1266.
- Rimm EB, Stampfer MJ, Colditz GA, et al. Validity of self-reported waist and hip circumferences in men and women. *Epidemiology* 1990;1:466–473.
- Kolonel LN, Henderson BE, Hankin JH, et al. A multiethnic cohort in Hawaii and Los Angeles: baseline characteristics. *Am J Epidemiol* 2000;151:346–357.
- Belanger CF, Hennekens CH, Rosner B, et al. The Nurses' Health Study. *Am J Nurs* 1978;78:1039–1040.
- Hennekens CH, Eberlein K. A randomized trial of aspirin and beta-carotene among U.S. physicians. *Prev Med* 1985;14:165–168.
- Christen WG, Gaziano JM, Hennekens CH. Design of Physicians' Health Study II—a randomized trial of beta-carotene, vitamins E and C, and multivitamins, in prevention of cancer, cardiovascular disease, and eye disease, and review of results of completed trials. *Ann Epidemiol* 2000;10:125–134.
- Prorok PC, Andriole GL, Bresalier RS, et al. Design of the Prostate, Lung, Colorectal and Ovarian (PLCO) cancer screening trial. *Control Clin Trials* 2000;21:273S–309S.
- Gohagan JK, Prorok PC, Hayes RB, et al. The Prostate, Lung, Colorectal and Ovarian (PLCO) cancer screening trial of the National Cancer Institute: history, organization, and status. *Control Clin Trials* 2000;21:251S–272S.
- National Cancer Institute. Cancer Genetic Markers of Susceptibility (CGEMS) data website. Available at: [http://cgems.cancer.gov/data\\_access.html](http://cgems.cancer.gov/data_access.html). CGEMS Data Accessed October 5, 2009.
- Yeager M, Chatterjee N, Ciampa J, et al. Identification of a new prostate cancer susceptibility locus on chromosome 8q24. *Nat Genet* 2009;41:1055–1057.
- Landi MT, Chatterjee N, Yu K, et al. A genome-wide association study of lung cancer identifies a region of chromosome 5p15 associated with risk for adenocarcinoma. *Am J Hum Genet* 2009;85:679–691.
- Newcomb PA, Zheng Y, Chia VM, et al. Estrogen plus progestin use, microsatellite instability, and the risk of colorectal cancer in women. *Cancer Res* 2007;67:7534–7539.
- White E, Patterson RE, Kristal AR, et al. ViTamins And Lifestyle cohort study: study design and characteristics of supplement users. *Am J Epidemiol* 2004;159:83–93.
- Hays J, Hunt JR, Hubbell FA, et al. The Women's Health Initiative recruitment methods and results. *Ann Epidemiol* 2003;13:S18–S77.
- The Women's Health Initiative Study Group. Design of the Women's Health Initiative clinical trial and observational study. *Control Clin Trials* 1998;19:61–109.
- Bergstralh EJ, Kosanke JL. Computerized matching of cases to controls. 56th ed. Rochester MN: Department of Health Sciences Research, Mayo Clinic, 1995.
- Zheng W, Chow WH, Yang G, et al. The Shanghai Women's Health Study: rationale, study design, and baseline characteristics. *Am J Epidemiol* 2005;162:1123–1131.
- Cai H, Zheng W, Xiang YB, et al. Dietary patterns and their correlates among middle-aged and elderly Chinese men: a report

- from the Shanghai Men's Health Study. *Br J Nutr* 2007;98:1006–1013.
28. Zheng W, Long J, Gao YT, et al. Genome-wide association study identifies a new breast cancer susceptibility locus at 6q25.1. *Nat Genet* 2009;41:324–328.
  29. Bei JX, Li Y, Jia WH, et al. A genome-wide association study of nasopharyngeal carcinoma identifies three new susceptibility loci. *Nat Genet* 2010;42:599–603.
  30. Matsuo K, Suzuki T, Ito H, et al. Association between an 8q24 locus and the risk of colorectal cancer in Japanese. *BMC Cancer* 2009;9:379.
  31. Nakata I, Yamashiro K, Yamada R, et al. Association between the SERPING1 gene and age-related macular degeneration and polypoidal choroidal vasculopathy in Japanese. *PLoS One* 2011;6:e19108.
  32. Jee SH, Sull JW, Lee JE, et al. Adiponectin concentrations: a genome-wide association study. *Am J Hum Genet* 2010;87:545–552.
  33. Edwards TL, Shrubsole MJ, Cai Q, et al. A study of prostaglandin pathway genes and interactions with current nonsteroidal anti-inflammatory drug use in colorectal adenoma. *Cancer Prev Res (Phila)* 2012;5:855–863.
  34. **Peters U, Hutter CM**, Hsu L, et al. Meta-analysis of new genome-wide association studies of colorectal cancer risk. *Hum Genet* 2012;131:217–234.
  35. Marchini J, Howie B, Myers S, et al. A new multipoint method for genome-wide association studies by imputation of genotypes. *Nat Genet* 2007;39:906–913.
  36. Browning SR, Browning BL. Rapid and accurate haplotype phasing and missing-data inference for whole-genome association studies by use of localized haplotype clustering. *Am J Hum Genet* 2007;81:1084–1097.
  37. Li Y, Willer CJ, Ding J, et al. MaCH: using sequence and genotype data to estimate haplotypes and unobserved genotypes. *Genet Epidemiol* 2010;34:816–834.
  38. Nothnagel M, Ellinghaus D, Schreiber S, et al. A comprehensive evaluation of SNP genotype imputation. *Hum Genet* 2009;125:163–171.
  39. Gogele M, Minelli C, Thakkestanian A, et al. Methods for meta-analyses of genome-wide association studies: critical assessment of empirical evidence. *Am J Epidemiol* 2012;175:739–749.
  40. Kent WJ, Sugnet CW, Furey TS, et al. The human genome browser at UCSC. *Genome Res* 2002;12:996–1006.
  41. Pruim RJ, Welch RP, Sanna S, et al. LocusZoom: regional visualization of genome-wide association scan results. *Bioinformatics* 2010;26:2336–2337.
- 

Author names in bold designate shared co-first authorship.

**Supplementary Table 1.** Descriptive Characteristics of Study Populations

Study name	Other name	Design	Country	Cases	Controls	Age range, y	Mean age, y	Female, %	Covariates used in analysis
GWAS									
Association Study Evaluating RISK for sporadic colorectal cancer	ASTERISK	Case-control	France	948	947	40–99	65.3	41.3	Age, sex, 3 PCAs, batch
Colorectal Cancer Studies 2&3	Hawaiian Colo2&3	Case-control	United States	87	125	38–86	65.2	44.8	Age, sex, 3 PCAs
Colon Cancer Family Registry*	CCFR	Case-control and sib-pair	United States, Canada, Australia	1967	1785	19–88	55.5	51.8	Age, sex, 3 PCAs, center
Darmkrebs: Chancen der Verhütung durch Screening	DACHS	Case-control	Germany	2376	2206	33–98	68.7	39.9	Age, sex, PCAs
Diet, Activity and Lifestyle Study	DALS	Case-control	United States	1116	1174	30–79	65.2	44.9	Age, sex, 3 PCAs, center
Health Professionals Follow-up Study	HPFS	Cohort	United States	403	402	48–83	65.2	0	age, 3PCAs
Multiethnic Cohort Study	MEC	Cohort	United States	328	346	45–76	63.0	46.4	Age, sex, 3 PCAs
Nurses' Health Study	NHS	Cohort	United States	553	955	44–69	59.8	100	Age, 3 PCAs
Ontario Familial Colorectal Cancer Registry	OFCCR	Case-control	Canada	650	522	31–79	64.1	52.0	Age, sex, 3 PCAs
Physicians' Health Study	PHS	Cohort	United States	382	389	40–84	58.4	0	Age, 3 PCAs, smoking
Postmenopausal Hormone study	PMH	Case-control	United States	280	122	50–75	64.8	100	Age, 3 PCAs
Prostate, Lung, Colorectal, and Ovarian Cancer Screening Trial	PLCO	Cohort	United States	1019	2391	55–75	64.0	30.8	Age, sex, 3 PCAs, center
VITamins And Lifestyle	VITAL	Cohort	United States	285	288	50–76	66.5	47.6	Age, sex, 3 PCAs
Women's Health Initiative	WHI	Cohort	United States	1476	2538	50–79	67.4	100	Age, 3 PCAs, region
Health Professionals' Follow-up Study, Adenoma Set	HPFS Ad	Cohort	United States	313	345	48–81	60.7	0	Age, 3 PCAs
Nurses' Health Study, Adenoma Set	NHS Ad	Cohort	United States	513	578	44–69	57.0	100	Age, 3 PCAs
Follow-up studies									
Asian Consortium, Colorectal Cancer				2098	5749				
Shanghai-1	Shanghai-1	Cohort	China	474	2628	25–75	53.22	91.62	Age, sex
Shanghai-2	Shanghai-2	Cohort	China	254	231	40–75	60.96	55.67	Age, sex
Guangzhou	Guangzhou	Case-control	China	641	972	14–85	50.36	30.81	Age, sex
Aichi	Aichi-1	Case-control	Japan	404	942	20–79	51.34	44.65	Age, sex
Korean Cancer Prevention Study-II	KCPS-II	Cohort	Korea	325	976	20–88	43.79	39.28	Age, sex
Tennessee Colorectal Polyp Study	TCPs	Case-control	United States	958	909	40–76	58.72	26.65	Age, sex

PCA, principal component analysis.

\*CCFR is a collaborating study with GECCO. The analysis of set 2 data did not adjust for PCs because of the sibling-pair study design.

**Supplementary Table 2.** Details on Genotyping Platform and Quality Assurance and Quality Control

Study	Genotyping platform <sup>a</sup>	Duplicate concordance, %	Mean sample call rate, %	SNP exclusions, <sup>b</sup> n	SNPs passing QC, n	Mean SNP call rate, %	Number of imputed SNPs by R <sup>2</sup>		
							<0.3	0.3–0.8	>0.8
ASTERISK	300K	100	99.97	30,446	252,176	99.95	76,043	443,302	1,856,490
Colo2&3	300K	100	99.95	40,390	258,978	99.96	71,487	445,613	1,854,778
DACHS Set 1	300K	99.9	99.93	33,588	255,208	99.90	70,989	434,295	1,869,458
DACHS Set 2	730K	100	99.84	32,159	609,115	99.85	18,551	154,813	1,865,294
DALS Set 1	550K, 610K	>97 <sup>c</sup>	99.69	34,644	516,631	99.82	20,173	180,322	1,912,832
DALS Set 2	300K	100	99.94	32,885	250,320	99.94	69,289	438,282	1,867,371
HPFS Set 1	730K	99.90	99.93	32,953	612,091	99.93	18,257	150,880	1,857,252
HPFS Set 2	730K	99.9	99.83	51,725	590,132	99.84	20,040	160,464	1,861,553
HPFS Ad	730K	100	99.86	61,201	597,470	99.86	18,610	155,527	1,861,220
MEC	300K	100	99.97	34,494	259,364	99.96	68,634	433,560	1,868,693
NHS Set 1	730K	100	99.93	47,295	628,541	99.93	17,142	147,723	1,855,814
NHS Set 2	730K	100	99.81	53,328	594,015	99.81	19,434	160,804	1,875,767
NHS Ad	730K	100	99.81	35,954	614,357	99.81	17,901	152,373	1,863,872
PHS Sets 1+2	730K	100	99.90	32,088	594,205	99.90	19,387	157,993	1,864,677
PLCO Set 1	300/240S and 610K	>97 <sup>c</sup>	99.65	33,342	503,351	99.85	20,855	184,854	1,921,986
PLCO Set 2	300K	99.90	99.80	38,655	253,702	99.90	68,059	434,769	1,870,311
PMH	300K	99.90	99.89	39,275	256,743	99.92	67,818	429,887	1,875,260
VITAL	300K	99.90	99.81	36,805	243,625	99.89	73,966	461,036	1,845,318
WHI set 1	550K, 550Kduo, 610K	>97 <sup>c</sup>	99.60	40,276	511,251	99.77	21,655	184,833	1,914,909
WHI set 2	300K	100	99.96	27,392	251,707	99.96	72,272	442,111	1,864,141

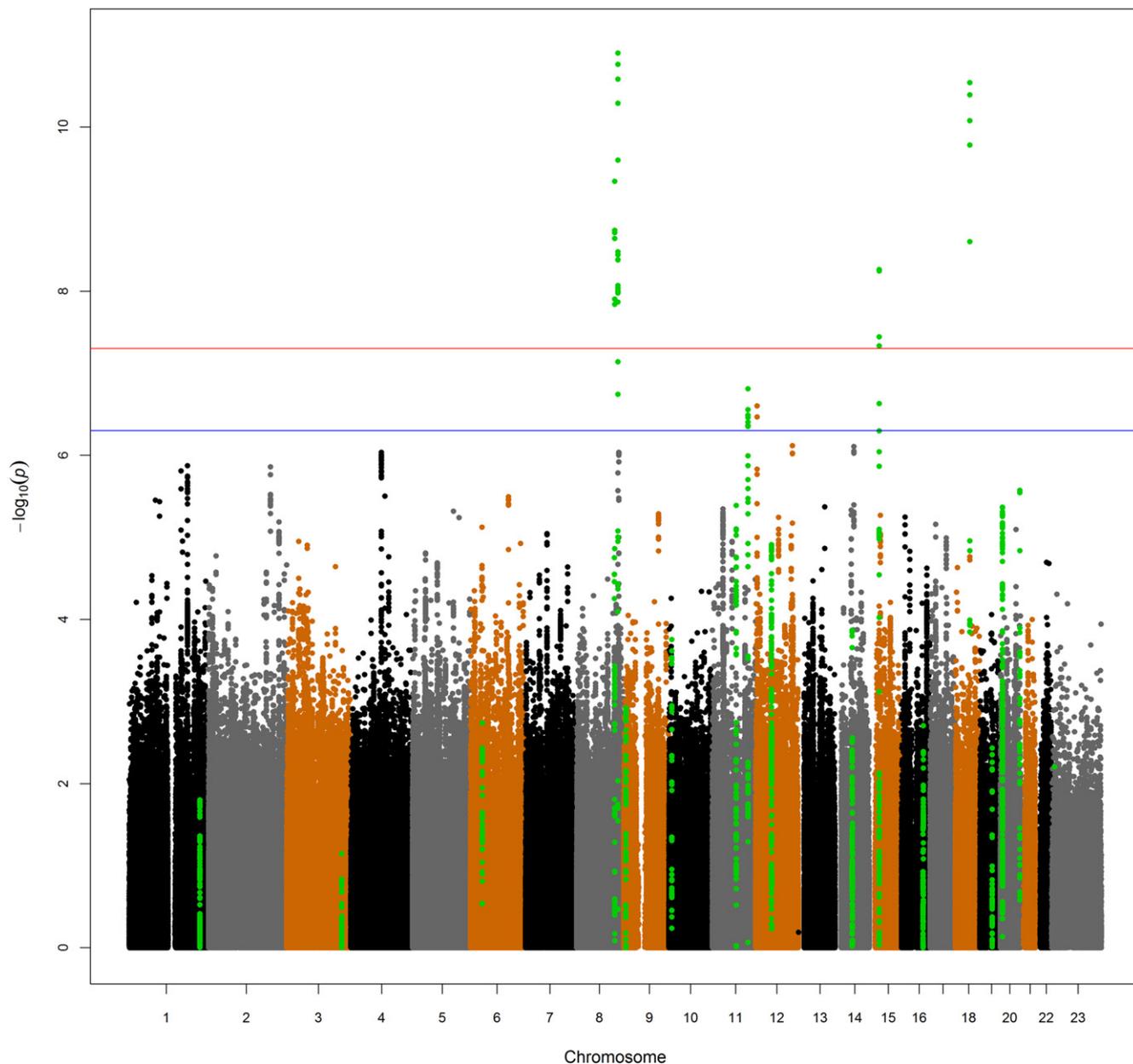
NOTE. CCFR and OFCCR had quality assurance/QC performed separately by OFCCR and CCFR investigators as documented by Zanke et al<sup>3</sup> and Figueiredo et al.<sup>6</sup>

ASTERISK, The french Association SStudy Evaluating RISK for sporadic colorectal cancer; Colo2&3, Hawaiian Colorectal Cancer Studies 2 and 3; DACHS, Darmkrebs: Chancen der Verhütung durch Screening; MEC, Multiethnic cohort; PMH, Postmenopausal Hormone study; VITAL, VITamins And Lifestyle.

<sup>a</sup>All platforms were Illumina assays, except for OFCCR, which was genotyped using Affymetrix platforms.

<sup>b</sup>Directly genotyped SNPs were excluded for a call rate less than 98%, HWE less than  $1 \times 10^{-4}$ , MAF less than 5 for WHI set 1, PLCO set 1, DALS set 1, and OFCCR set 1; MAF less than 5 per number of samples for remaining studies, and if SNPs reportedly did not perform consistently across platforms.

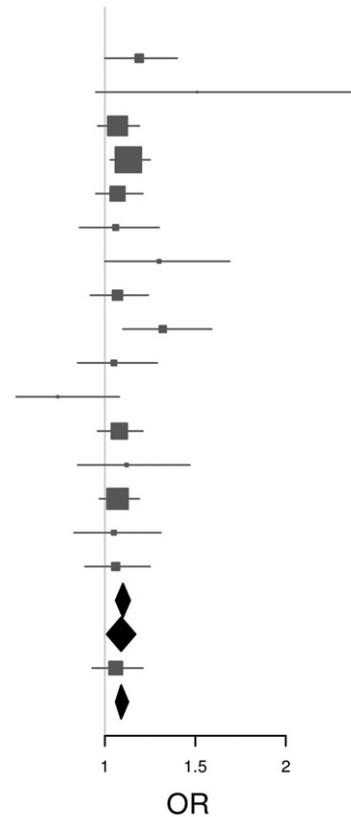
<sup>c</sup>Blinded duplicates were assessed across DALS set 1, PLCO set 1, and WHI set 1; exact concordance was not recorded, but all 98 pairs were identified as having concordance greater than 97%.



**Supplementary Figure 1.** Manhattan plot of the GWAS inverse-variance-weighted, fixed-effects meta-analysis, comprising 12,696 cases and 15,113 controls. The  $-\log_{10}$  of  $P$  values for 2,708,280 SNPs plotted against physical chromosomal positions. SNPs above the *blue line* represent those with a  $P$  value less than  $5 \times 10^{-7}$  whereas SNPs above the *red line* represent those with a  $P$  value less than  $5 \times 10^{-8}$ . The *green dots* represent previously identified loci as listed in Supplementary Table 4. Chromosome 23 is the X-chromosome. Because we do not have linkage disequilibrium (LD) information for SNPs on the X chromosome, we only show the result of the GWAS SNP on the X chromosome but not SNPs correlated with this GWAS SNP.

## rs10911251

Study	OR	95%CI	P
ASTERISK	1.19	(1.00–1.40)	4.52e-02
COLO23	1.51	(0.95–2.41)	8.46e-02
CCFR	1.07	(0.96–1.19)	2.22e-01
DACHS	1.13	(1.03–1.25)	7.97e-03
DALS	1.07	(0.95–1.21)	2.65e-01
HPFS	1.06	(0.86–1.30)	5.73e-01
MEC	1.30	(1.00–1.69)	5.09e-02
NHS	1.07	(0.92–1.24)	3.84e-01
OFCCR	1.32	(1.10–1.59)	2.76e-03
PHS	1.05	(0.85–1.29)	6.58e-01
PMH	0.74	(0.51–1.08)	1.16e-01
PLCO	1.08	(0.96–1.21)	2.19e-01
VITAL	1.12	(0.85–1.47)	4.37e-01
WHI	1.07	(0.97–1.19)	1.83e-01
HPFS Ad	1.05	(0.83–1.31)	6.96e-01
NHS Ad	1.06	(0.89–1.25)	5.28e-01
<b>GWAS</b>	<b>1.10</b>	<b>(1.06–1.14)</b>	<b>1.34e-06</b>
<b>Asian Follow-up</b>	<b>1.09</b>	<b>(1.01–1.17)</b>	<b>3.20e-02</b>
Adenoma Follow-up	1.06	(0.93–1.21)	3.66e-01
<b>GWAS+Follow-up</b>	<b>1.09</b>	<b>(1.06–1.13)</b>	<b>9.45e-08</b>

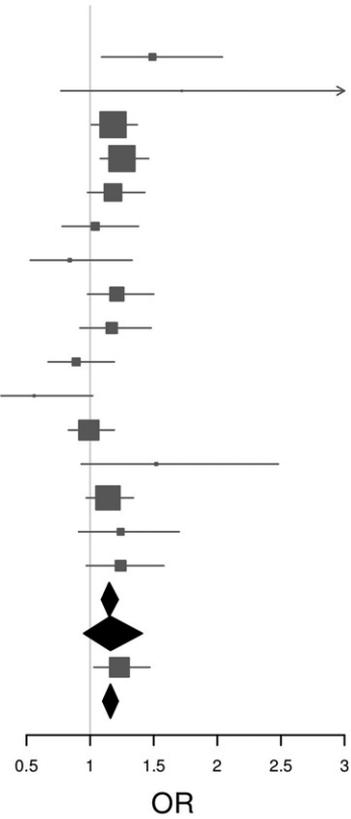


Het pv=0.687

**Supplementary Figure 2.** Forest plot for meta-analysis results for all new findings with a  $P$  value less than  $5 \times 10^{-7}$  in a combined analysis of GWAS and follow-up studies as listed in Table 1. ORs and 95% confidence intervals (95% CIs) are presented for each additional copy of the minor allele in the multiplicative model. The *grey boxes* are proportional in size to the inverse of the variance for each study, and the *lines* visually depict the confidence interval. Results from the fixed-effects meta-analysis are shown as *diamonds*. The width of the diamond represents the confidence interval.

**rs11903757**

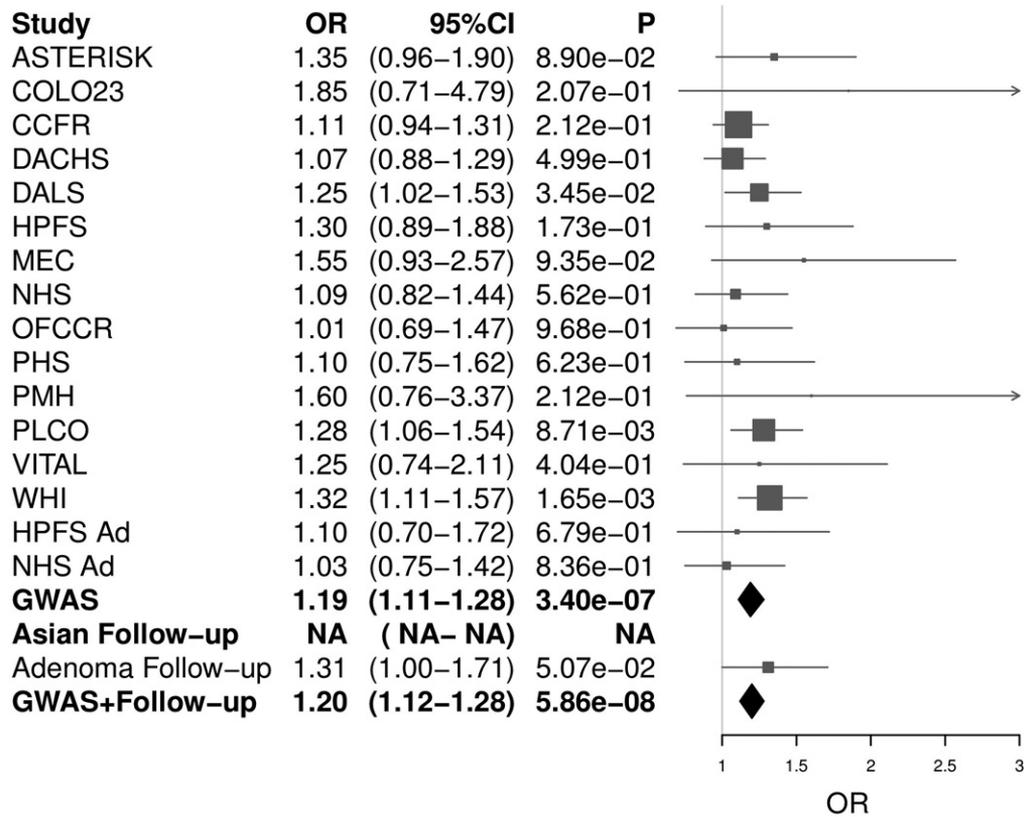
<b>Study</b>	<b>OR</b>	<b>95%CI</b>	<b>P</b>
ASTERISK	1.49	(1.09–2.04)	1.13e–02
COLO23	1.72	(0.77–3.83)	1.83e–01
CCFR	1.18	(1.01–1.37)	3.78e–02
DACHS	1.25	(1.08–1.46)	3.92e–03
DALS	1.18	(0.98–1.43)	8.53e–02
HPFS	1.04	(0.78–1.38)	7.83e–01
MEC	0.84	(0.53–1.33)	4.57e–01
NHS	1.21	(0.98–1.50)	7.23e–02
OFCCR	1.17	(0.92–1.48)	1.96e–01
PHS	0.89	(0.67–1.19)	4.45e–01
PMH	0.56	(0.30–1.02)	6.00e–02
PLCO	0.99	(0.83–1.19)	9.50e–01
VITAL	1.52	(0.93–2.48)	9.53e–02
WHI	1.14	(0.97–1.34)	1.04e–01
HPFS Ad	1.24	(0.91–1.70)	1.70e–01
NHS Ad	1.24	(0.97–1.58)	8.03e–02
<b>GWAS</b>	<b>1.15</b>	<b>(1.09–1.22)</b>	<b>1.38e–06</b>
<b>Asian Follow-up</b>	<b>1.16</b>	<b>(0.95–1.41)</b>	<b>1.34e–01</b>
Adenoma Follow-up	1.23	(1.03–1.47)	2.27e–02
<b>GWAS+Follow-up</b>	<b>1.16</b>	<b>(1.10–1.22)</b>	<b>3.71e–08</b>



Het pv=0.271

Supplementary Figure 2. Continued

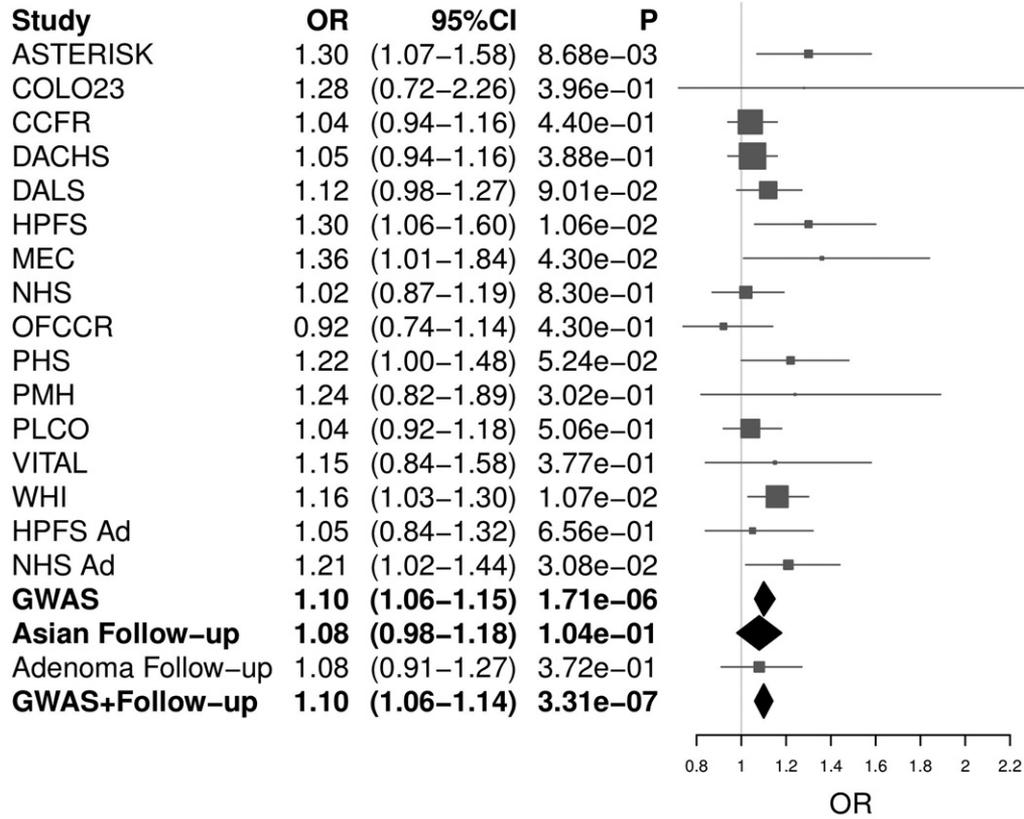
## rs3217810



Het pv=0.910

Supplementary Figure 2. Continued

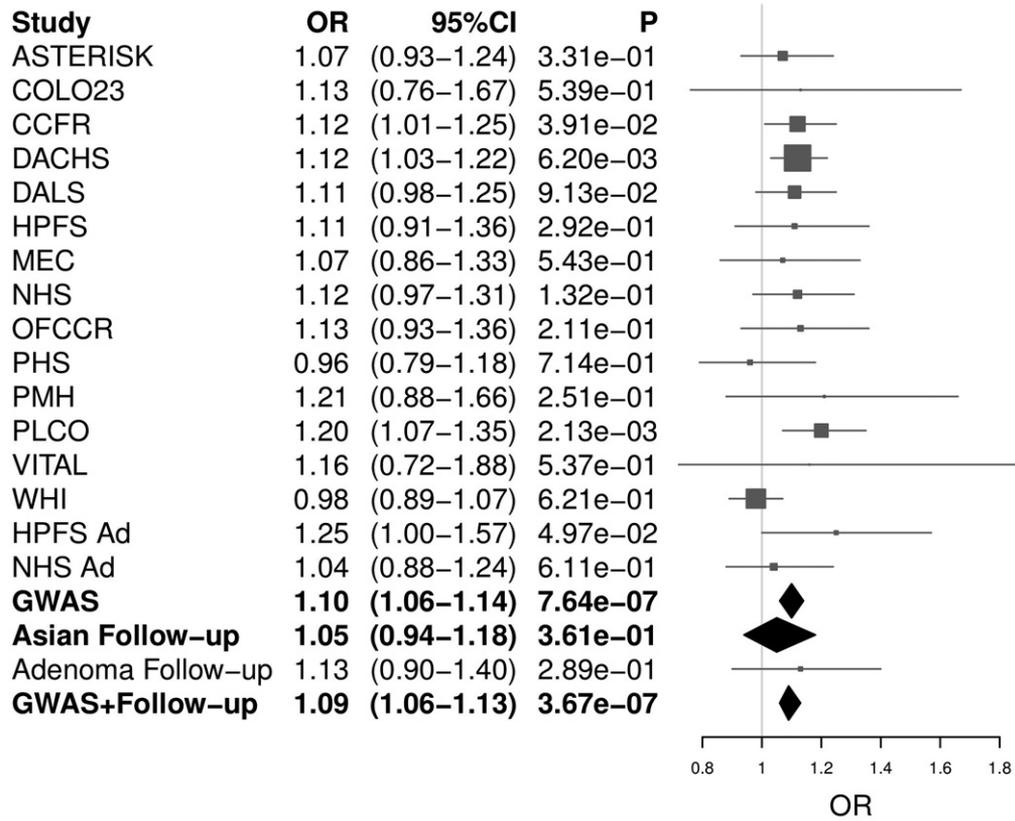
**rs3217901**



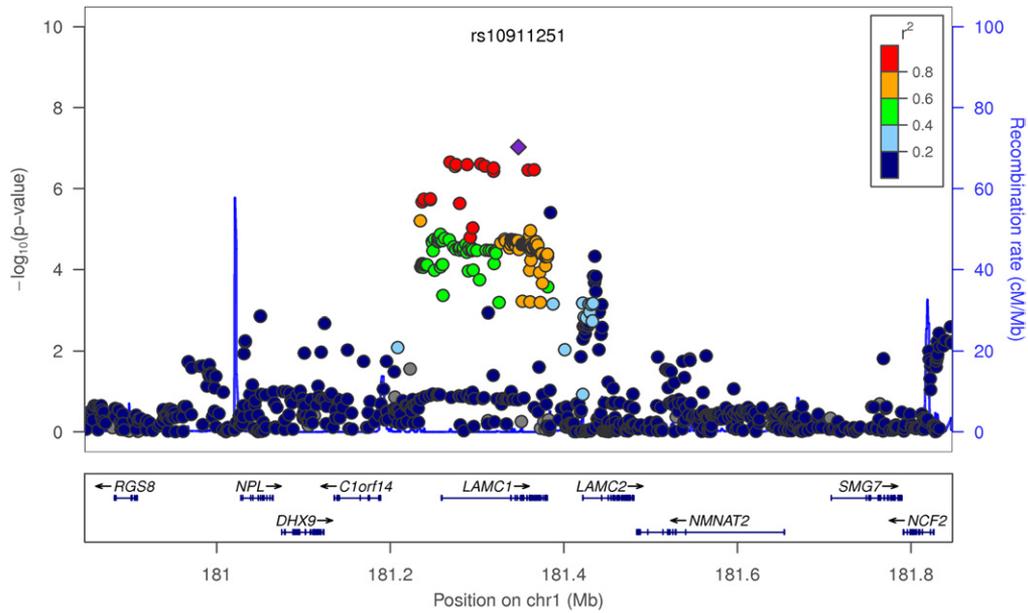
Het pv=0.507

Supplementary Figure 2. Continued

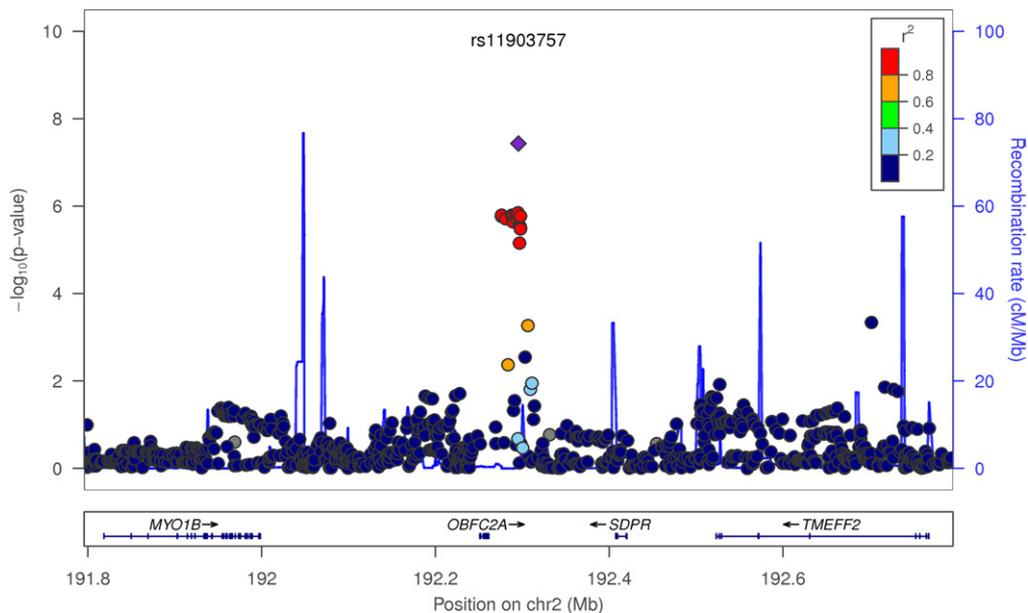
## rs59336

Het  $p_v=0.387$ 

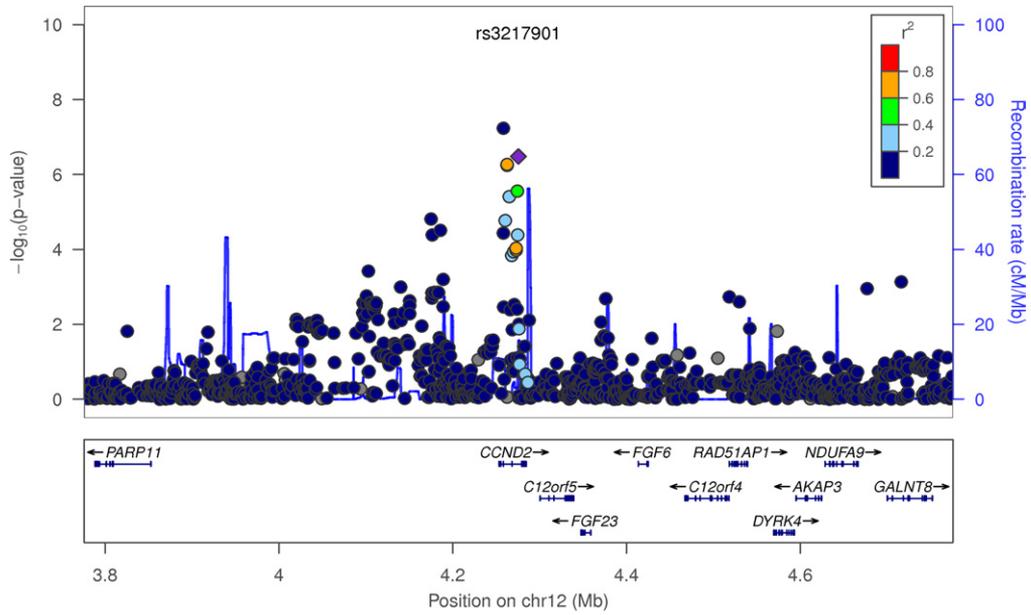
Supplementary Figure 2. Continued



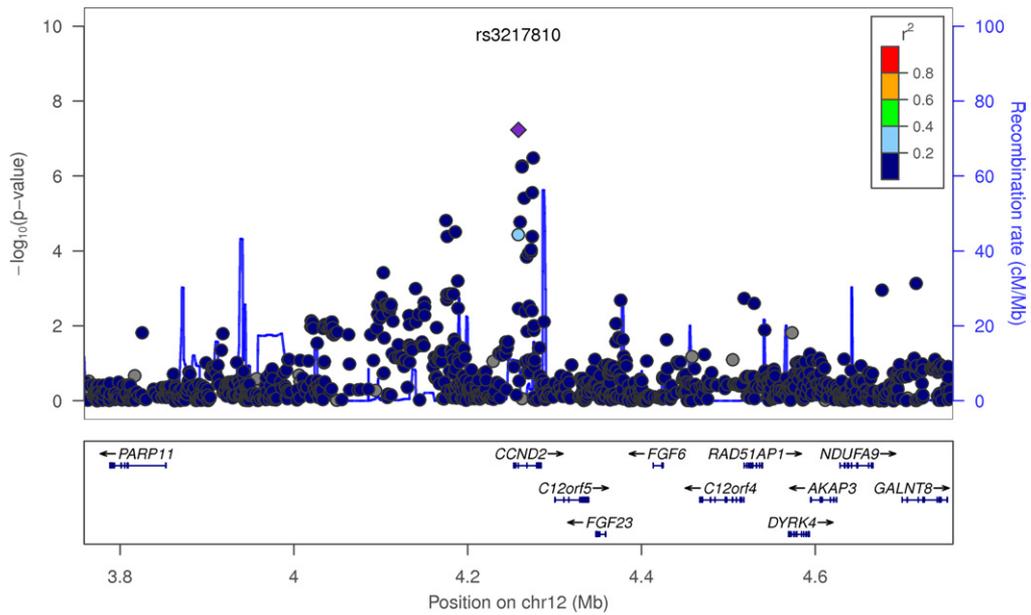
**Supplementary Figure 3.** Regional association results for all new findings with a  $P$  value less than  $5 \times 10^{-7}$ , as listed in Table 1. The *top half* of the figure shows the physical position of the SNP on the chromosome along the x-axis, and the  $-\log_{10}$  of the meta-analysis  $P$  value on the y-axis. Each *dot* on the plot represents the  $P$  value of the association for one SNP with colorectal cancer (allele test) across all studies. The most significant SNP in the region (index SNP) is marked as a *purple diamond*. The color scheme represents the pairwise correlation ( $r^2$ ) for the SNPs across the region with the index SNP. Correlation was calculated using the HapMap CEU data. The *bottom half* of the figure shows the position of the genes across the region. These regional association plots are also known as LocusZoom plots.<sup>41</sup>



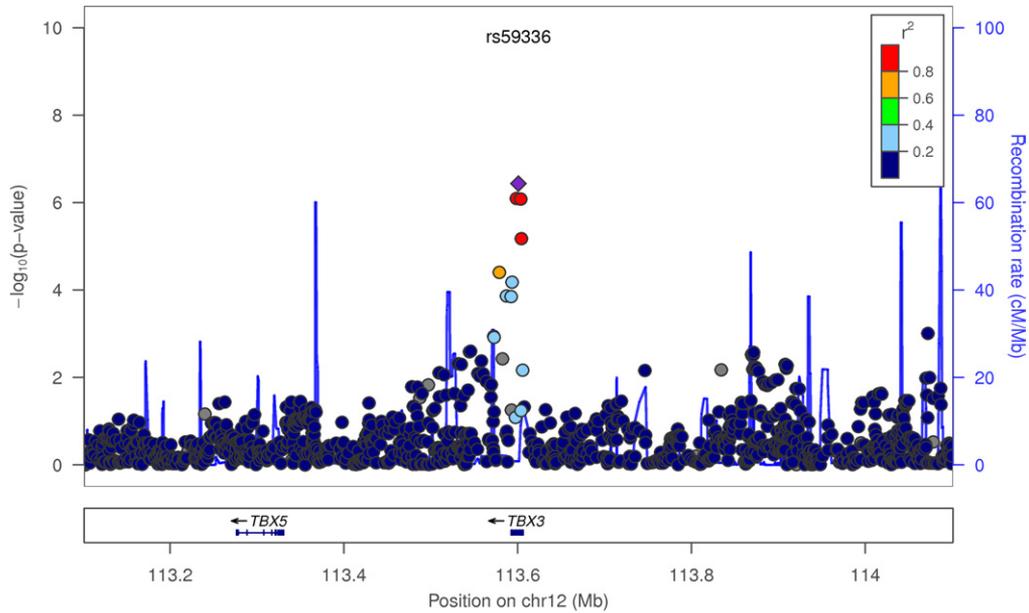
**Supplementary Figure 3.** Continued



Supplementary Figure 3. Continued



Supplementary Figure 3. Continued



Supplementary Figure 3. Continued

**Supplementary Table 5.** Risk Estimates for the 2 Top SNPs in 12p13.32/*CCND2* When Both Were Included Simultaneously in the Logistic Regression Analysis

SNP	OR (95% CI)	P value
Each SNP analyzed separately		
rs3217901	1.10 (1.06–1.15)	1.71E-06
rs3217810	1.19 (1.11–1.28)	3.40E-07
Both SNPs included simultaneously in the logistic regression analysis		
rs3217901	1.08 (1.03–1.13)	.0008
rs3217810	1.14 (1.06–1.23)	.004

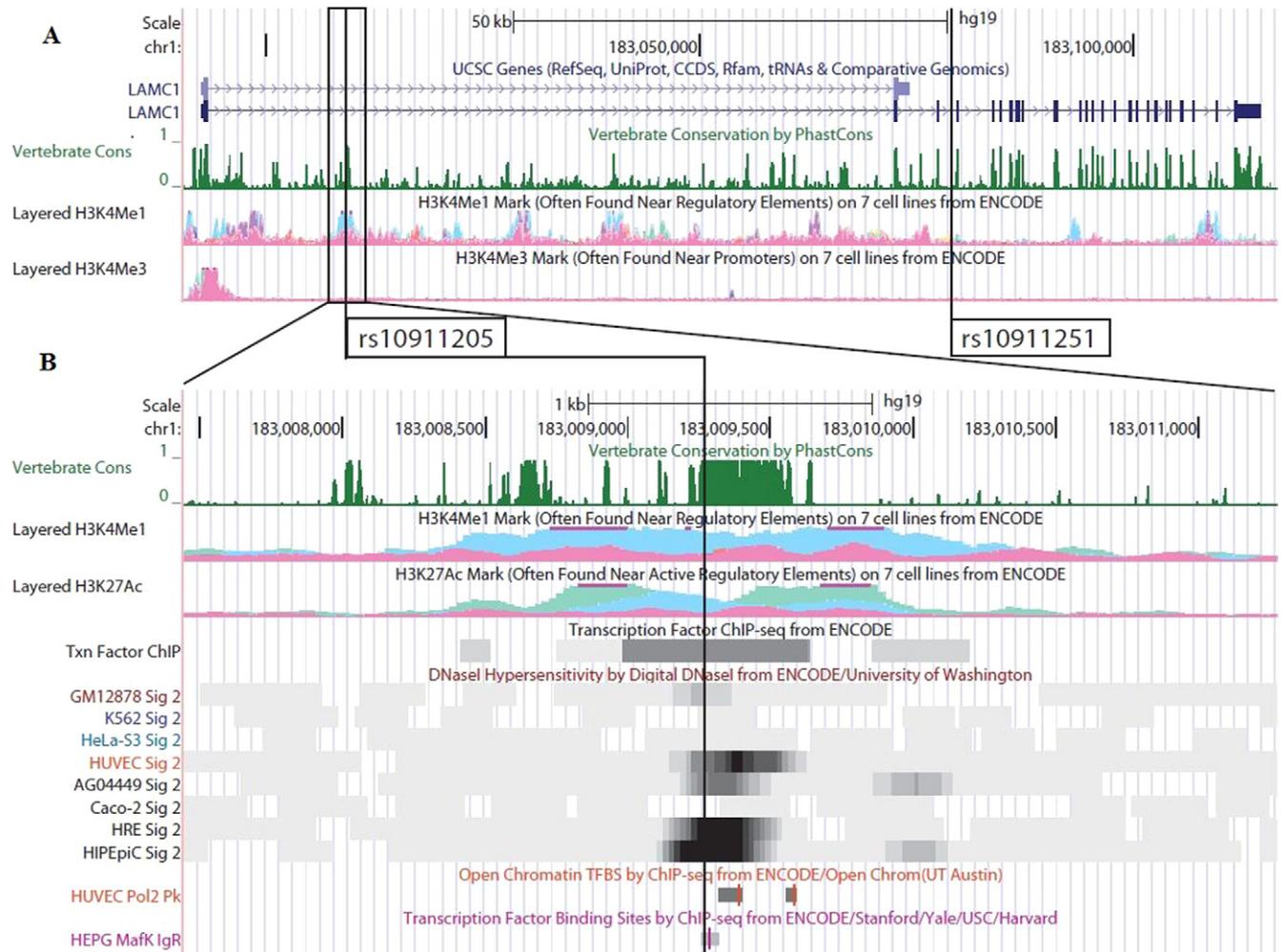
NOTE. Analysis was based on the log-additive model in GWAS of GECCO and CCFR only (12,696 cases and 15,113 controls). CI, confidence interval.

**Supplementary Table 7.** Risk Estimates for New Findings With  $P < 5 \times 10^{-7}$  Stratified by Colorectal Adenoma and Colorectal Cancer (Log-Additive Model)

SNP	Chromosome (gene)	Cancer/adenoma <sup>a</sup>	OR (95% CI)	P value	P heterogeneity
SNP with $P < 5 \times 10^{-8}$					
rs11903757	2q32.3 (NABP1)	Cancer	1.15 (1.08–1.21)	4.06E-06	.18
		Adenoma	1.24 (1.08–1.41)	1.46E-03	1.00
		Overall	1.16 (1.10–1.22)	3.71E-08	.27
SNPs with $P < 5 \times 10^{-7}$ and $P > 5 \times 10^{-8}$					
rs10911251	1q25.3 (LAMC1)	Cancer	1.10 (1.06–1.14)	1.41E-07	.55
		Adenoma	1.06 (0.96–1.16)	2.44E-01	.99
		Overall	1.09 (1.06–1.13)	9.45E-08	.69
rs3217810	12p13.32 (CCND2)	Cancer	1.20 (1.12–1.29)	2.34E-07	.87
		Adenoma	1.17 (0.97–1.41)	9.76E-02	.52
		Overall	1.20 (1.12–1.28)	5.86E-08	.91
rs3217901	12p13.32 (CCND2)	Cancer	1.10 (1.05–1.14)	2.98E-06	.41
		Adenoma	1.12 (1.01–1.25)	3.64E-02	.53
		Overall	1.10 (1.06–1.14)	3.31E-07	.51
rs59336	12q24.21 (TBX3)	Cancer	1.09 (1.05–1.13)	2.21E-06	.31
		Adenoma	1.12 (1.00–1.25)	5.73E-02	.44
		Overall	1.09 (1.06–1.13)	3.67E-07	.39

CI, confidence interval.

<sup>a</sup>Cancer (n = 13,968 cases and 19,939 controls, except for rs3217810, which has 11,870 cases and 14,190 controls); adenoma (1784 cases and 1832 controls); overall (15,752 cases and 21,771 controls, except for rs3217810, which has 13,654 cases and 16,022 controls).



**Supplementary Figure 4.** ENCODE integrate regulation tracks for *LAMC1*. (A) University of California, Santa Cruz (UCSC) genome browser position chr1:182,990,493–183,116,512 (build 37) containing the *LAMC1* protein coding gene. The University of California, Santa Cruz gene track shows 2 variant transcripts for *LAMC1*. Directly beneath the gene track is a histogram of multiple alignments of 46 vertebrate species indicating that there are multiple conserved elements in the gene, primarily concentrated near the 5' and 3' regulatory regions. Conservation can help unmask candidate variants that disrupt regulatory regions from other benign associations. The next 2 tracks are transparent overlays from 7 cell lines assayed by the ENCODE project showing the H3K4me1 mark and the H3K4me3 mark associated with active regulatory regions. Peaks in H3K4me3 mark are consistent with the promoter region of *LAMC1*, whereas H3K4me1 indicates additional enhancer regions in the first intron. The histone marks and pattern of transcription show coordinated, cell-type-specific activity increases in K562 (blue) and NHLF (pink) cells. (B) Focusing on the region containing rs10911205 (chr1:183,007,443–183,011,275), we find that the SNP lies within a strong evolutionarily constrained region. Below this track, evidence from the H3K4me1 and H3K27Ac marks are consistent with rs10911205 falling within a region of coordinated, cell-type-specific activity, most active in K562 (blue) cells and human skeletal muscle myoblasts (green) cells. The DNase and transcription factor chromatin immunoprecipitation sequencing (ChIP-seq) clusters shown in the last 2 tracks summarize data from a much wider range of cell lines and further supports tissue-specific accessibility for regulatory elements in the region surrounding rs10911205. The  $r^2$  value of rs10911205 with rs10911251 was 0.862. Taken together, evidence provided by the ENCODE integrated regulation tracks is consistent with rs10911205 being a strong functional candidate SNP for the strong rs10911251 association with colorectal cancer.

**Supplementary Table 8.** Sample Size and Genotyping Methods Used in Asian GWAS

Study	Genotyped		After quality control		Genotyping platform		Number of SNPs <sup>a</sup>	Inflation factor ( $\lambda$ ) <sup>b</sup>
	Cases	Controls	Cases	Controls	Cases	Controls		
Shanghai-1	481	2632	474	2628	Affymetrix 6.0	Affymetrix 6.0	502,145	1.03
Shanghai-2	296	257	254	231	Illumina	Illumina	515,701	1.03
Guangzhou-1	694	972	641	972	OmniExpress	OmniExpress	250,612	1.02
					Illumina	Illumina		
Aichi-1	497	942	404	942	Illumina	Illumina	232,426	1.04
					OmniExpress	HumanHap610		
KCPS-II	325	977	325	976	Affymetrix 5.0	Affymetrix 5.0	312,869	1.02
Overall	2293	5780	2098	5749				1.01

NOTE. Number of cases and controls differ from Supplementary Table 1 due to quality assurance/quality control exclusions.

<sup>a</sup>Number of SNPs in autosome used for imputation in GWAS.

<sup>b</sup>Genomic inflation factor ( $\lambda$ ) derived from 1,636,780 imputed SNPs with MAF >0.05 and high imputation quality (RSQR >0.50), adjusted with age, sex, and the first 10 principal components.

**Supplementary Table 9.** Tools for Functional Annotation of Noncoding Variants

UCSC genome browser	Genomic class	Description	Functional evidence
ENCODE transcription	Transcribed region	Transcription levels in 7 cell lines from ENCODE Assayed by high-throughput sequencing of polyadenylated RNA	Variable expression in different tissues provides evidence for cell-type-specific regulation when displayed as transparent overlay of each cell line
ENCODE layered H3K4Me1	Nonpromoter regulatory elements	Uses ChIP-seq method to identify regions of DNA that interact with the mono-methylation of lysine 4 of the H3 histone protein in 7 different cell lines Actual enhancer is likely a small portion of the broad region marked	Methylation of histone proteins changes chromatin accessibility for transcription H3K4Me1 is associated with enhancers downstream of the transcription start site
ENCODE layered H3K4Me3	Promoter regulatory element	Uses ChIP-seq method to identify regions of DNA that interact with the trimethylation of lysine 4 of the H3 histone protein in 7 different cell lines Actual regulatory element is likely a small portion of the broad region marked	H3K4Me3 is associated with promoters that are active or accessible for activation
ENCODE layered H3K27Ac	Nonpromoter regulatory elements	Uses ChIP-seq method to identify regions of DNA that interact with the acetylation of lysine 27 of the H3 histone protein in 7 different cell lines Actual regulatory element is likely a small portion of the broad region marked	H3K27Ac enhances transcription, possibly by blocking the spread of the repressive histone mark H3K27Me3 This mark often is found near active regulatory elements
ENCODE DNase clusters	Regulatory element	Measures digital DNaseI hypersensitivity clusters in a large collection of cell types from ENCODE Greater precision than histone modifications	Regulatory regions and promoters are susceptible DNase cutting Hypersensitivity is used to map chromatin accessibility
ENCODE Txn factor ChIP	Regulatory element	Transcription Factor ChIP-seq from ENCODE is assayed by chromatin immunoprecipitation using antibodies for specific transcription factors and sequencing the precipitated DNA	Marks regions where transcription factors bind DNA and exert specific functions Activators can recruit RNA polymerase, repressors suppress recruitment, and insulators block the activity of nearby activators or repressors
ENCODE UW CTCF binding (within the ENCODE transcription factor binding tracks)	Insulated element	CTCF binding sites are assayed by chromatin immunoprecipitation using antibodies for CTCF and sequencing the precipitated DNA	CTCF can function as a transcriptional activator, a repressor/silencer, or an insulator Binds chromatin insulators to prevent interaction between promoter and nearby enhancers or silencers Also mediates long-range chromatin looping, which can bring enhancers in proximity of a gene's promoter
Vertebrate multi-alignment and conservation (phastCons)	Conserved element	Multiple alignments of 46 vertebrate species Estimates the probability that each nucleotide belongs to a conserved element	Identification of evolutionarily conserved segments of homology, potentially identifying a functionally important region

ChIP-seq, chromatin immunoprecipitation sequencing; CTCF, CCCTC-binding factor; UCSC, University of California, Santa Cruz.