

# Discovery of common and rare genetic risk variants for colorectal cancer

**To further dissect the genetic architecture of colorectal cancer (CRC), we performed whole-genome sequencing of 1,439 cases and 720 controls, imputed discovered sequence variants and Haplotype Reference Consortium panel variants into genome-wide association study data, and tested for association in 34,869 cases and 29,051 controls. Findings were followed up in an additional 23,262 cases and 38,296 controls. We discovered a strongly protective 0.3% frequency variant signal at *CHD1*. In a combined meta-analysis of 125,478 individuals, we identified 40 new independent signals at  $P < 5 \times 10^{-8}$ , bringing the number of known independent signals for CRC to ~100. New signals implicate lower-frequency variants, Krüppel-like factors, Hedgehog signaling, Hippo-YAP signaling, long noncoding RNAs and somatic drivers, and support a role for immune function. Heritability analyses suggest that CRC risk is highly polygenic, and larger, more comprehensive studies enabling rare variant analysis will improve understanding of biology underlying this risk and influence personalized screening strategies and drug development.**

Colorectal cancer (CRC) is the fourth leading cancer-related cause of death worldwide<sup>1</sup> and presents a major public health burden. Up to 35% of interindividual variability in CRC risk has been attributed to genetic factors<sup>2,3</sup>. Family-based studies have identified rare high-penetrance mutations in at least a dozen genes, but collectively, these account for only a small fraction of familial risk<sup>4</sup>. Over the past decade, genome-wide association studies (GWASs) for sporadic CRC, which constitutes the majority of cases, have identified ~60 association signals at over 50 loci<sup>5–22</sup>. Yet most of the genetic factors contributing to CRC risk remain undefined. This severely hampers our understanding of biological processes underlying CRC. It also limits CRC precision prevention, including individualized preventive screening recommendations and development of cancer prevention drugs. The contribution of rare variation to sporadic CRC is particularly poorly understood.

To expand the catalog of CRC risk loci and improve our understanding of rare variants, genes and pathways influencing sporadic CRC risk and risk prediction, we performed the largest and most comprehensive whole-genome sequencing (WGS) study and GWAS meta-analysis for CRC so far, combining data from three consortia: the Genetics and Epidemiology of Colorectal Cancer Consortium (GECCO), the Colorectal Cancer Transdisciplinary Study (CORECT) and the Colon Cancer Family Registry (CCFR). Our study almost doubles the number of individuals analyzed, incorporating GWAS results from >125,000 individuals, and substantially expands and strengthens our understanding of biological processes underlying CRC risk.

## Results

**Study overview.** We performed WGS of 1,439 CRC cases and 720 controls of European ancestry at low sequencing depth (3.8–8.6×). We detected, called and estimated haplotype phase for 31.8 million genetic variants, including 1.7 million short insertion-deletion variants (indels) (Methods). These data include many rare variants not studied by GWAS. As other large-scale WGS studies have used a similar design, we expected to have near-complete ascertainment of single-nucleotide variants with minor allele count (MAC) > 5 (minor allele frequency (MAF) > 0.1%), and high accuracy at heterozygous genotypes<sup>23,24</sup>. We tested 14.4 million variants with  $\text{MAC} \geq 5$  for CRC association using logistic regression (Methods) but did not find any significant associations. To increase power to

detect associations with rare and low-frequency variants of modest effect, we imputed variants from the sequencing experiment into 34,869 cases and 29,051 controls of predominantly European (91.7%) and East Asian ancestry (8.3%) from 30 existing GWASs (Methods and Supplementary Table 1). By design, two-thirds of sequenced individuals were CRC cases, thereby enriching the panel for rare or low-frequency alleles that increase CRC risk. We contributed our sequencing data to the Haplotype Reference Consortium (HRC)<sup>25</sup> and imputed the 30 existing GWASs to the HRC panel, which comprises haplotypes for 32,488 individuals. Results of these GWAS meta-analyses (referred to as stage 1 meta-analysis; Methods) informed the design of a custom Illumina array comprising the OncoArray, a custom array to identify cancer risk loci<sup>26</sup>, and 15,802 additional variants selected based on stage 1 meta-analysis results. We genotyped 12,007 cases and 12,000 controls of European ancestry with this custom array, and combined them with an additional 11,255 cases and 26,296 controls with GWAS data, resulting in a stage 2 meta-analysis of 23,262 CRC cases and 38,296 controls (Methods, Supplementary Fig. 1 and Supplementary Table 1). Next, we performed a combined (stage 1 + stage 2) meta-analysis of up to 58,131 cases and 67,347 controls. This meta-analysis was based on the HRC panel-imputed data because, given its large size, this panel results in superior imputation quality and enables accurate imputation of variants with MAFs as low as 0.1% (ref. <sup>25</sup>). Here, we report new association signals discovered through our custom genotyping experiment and replicated in stage 2 at the Bonferroni significance threshold of  $P < 7.8 \times 10^{-6}$  (Methods), as well as distinct association signals passing the genome-wide significance (GWS) threshold of  $P < 5 \times 10^{-8}$  in the combined meta-analysis of up to 125,478 individuals.

**Colorectal cancer risk loci.** In the combined meta-analysis, we identified 30 new CRC risk loci reaching GWS and > 500 kilobases (kb) away from previously reported CRC risk variants (Table 1 and Supplementary Figs. 2 and 3). Some 22 of these were represented on our custom genotyping panel, either by the lead variant (15 loci) or by a variant in linkage disequilibrium (LD) (7 loci;  $r^2 > 0.7$ ). Of these 22 variants, eight attained the Bonferroni significance threshold in the stage 2 meta-analysis (Table 1).

Among these eight loci is the first rare variant signal identified for sporadic CRC, involving five 0.3% frequency variants at

5q21.1, near genes *CHD1* and *RGMB*. SNP rs145364999, which is intronic to *CHD1*, had high-quality genotyping (Supplementary Fig. 4). The variant was well imputed in the remaining sample sets (imputation quality  $r^2$  of 0.66–0.87; Supplementary Table 2) and there was no evidence of heterogeneity of effects (heterogeneity  $P=0.63$ ; Supplementary Table 2). The rare allele confers a strong protective effect (allelic odds ratio (OR)=0.52 in stage 2; 95% confidence interval (CI)=0.40–0.68). Chromatin remodeling factor *CHD1* provides an especially plausible candidate and is a synthetically essential gene<sup>27</sup> that is occasionally deleted in some cancers, but always retained in PTEN-deficient cancers<sup>28</sup>. The resulting mutually exclusive deletion pattern of *CHD1* and *PTEN* has been observed in prostate, breast and CRC The Cancer Genome Atlas data<sup>28</sup>. We hypothesize that the rare allele confers a protective effect by lowering *CHD1* expression, which is required for nuclear factor- $\kappa$ B (NF- $\kappa$ B) pathway activation and growth in cancer cells driven by loss of the tumor suppressor gene *PTEN*<sup>28</sup>. However, we cannot rule out involvement of nearby candidate gene *RGMB* that encodes a co-receptor for bone morphogenetic proteins BMP2 and BMP4, both of which are linked to CRC risk through GWAS<sup>9,11</sup>. Additionally, *RGMB* binds to PD-L2 (ref. 29), a known ligand of PD-1, an immune checkpoint blockade receptor targeted by cancer immunotherapy<sup>30</sup>.

The vast majority of new association signals involve common variants. We found associations near strong candidate genes for CRC risk in pathways or gene families not previously implicated by GWAS. Locus 13q22.1, represented by lead SNP rs78341008 (MAF 7.2%;  $P=3.2\times 10^{-10}$ ), is near *KLF5*, a known CRC oncogene that can be activated by somatic hotspot mutations or super-enhancer duplications<sup>31,32</sup>. *KLF5* encodes transcription factor Krüppel-like factor 5 (KLF5), which promotes cell proliferation and is highly expressed in intestinal crypt stem cells. We also found an association at locus 19p13.11, near *KLF2*. *KLF2* expression in endothelial cells is critical for normal blood vessel function<sup>33,34</sup>. Downregulated *KLF2* expression in colon tumor tissues contributes to structurally and functionally abnormal tumor blood vessels, leading to impaired blood flow and hypoxia in tumors<sup>35</sup>. Another locus at 9q31.1 is near *LPAR1*, which encodes a receptor for lysophosphatidic acid (LPA). LPA-induced expression of hypoxia-inducible factor 1 (HIF-1 $\alpha$ ), a key regulator of cellular adaptation to hypoxia and tumorigenesis, depends on *KLF5* (ref. 36). Additionally, LPA activates multiple signaling pathways and stimulates proliferation of colon cancer cells by activation of *KLF5* (ref. 37). Another locus (7p13) is near *SNHG15*, which encodes a long noncoding RNA (lncRNA) that epigenetically represses *KLF2* to promote pancreatic cancer proliferation<sup>38</sup>.

We found two loci near members of the Hedgehog (Hh) signaling pathway. Aberrant activation of this pathway, caused by somatic mutations or changes in expression, can drive tumorigenesis in many tumors<sup>39</sup>. Notably, downregulated stromal cell Hh signaling accelerates colonic tumorigenesis in mice<sup>40</sup>. Locus 3q13.2, represented by low-frequency lead SNP rs72942485 (MAF 2.2%;  $P=2.1\times 10^{-8}$ ), overlaps with *BOC*, which encodes a Hh coreceptor molecule. In medulloblastoma, upregulated *BOC* promotes Hh-driven tumor progression through cyclin D1-induced DNA damage<sup>41</sup>. In pancreatic cancer, a complex role for stromal *BOC* expression in tumorigenesis and angiogenesis has been reported<sup>42</sup>. Locus 4q31.21 is near *HHIP*, which encodes an inhibitor of Hh signaling. Notably, the Hh signaling pathway was also significantly enriched in our pathway analysis (described below).

Locus 11q22.1 is near *YAP1*, which encodes a critical downstream regulatory target in the Hippo signaling pathway that is gaining recognition as a pivotal player in organ size control and tumorigenesis<sup>43</sup>. *YAP1* is highly expressed in intestinal crypt stem cells, and in transgenic mice, its overexpression led to severe intestinal dysplasia and loss of differentiated cell types<sup>44</sup>, reminiscent of phenotypes observed in mice and humans with deleterious germline

*APC* mutations. Further, hypoxia-inducible factor 2 $\alpha$  (HIF-2 $\alpha$ ) promotes colon cancer growth by upregulating *YAP1* activity<sup>45</sup>.

We provide further evidence for a link between immune function and CRC pathogenesis, and implicate the major histocompatibility complex (MHC) in CRC risk. We identified a locus near genes *HLA-DRB1* and *HLA-DQA1* that is associated with immune-mediated diseases<sup>46</sup>.

We identified two new loci near known tumor suppressor genes. Locus 4q24 is near *TET2*, a chromatin-remodeling gene that is frequently somatically mutated in multiple cancers, including colon cancer<sup>47</sup>, and that overlaps with GWAS signals for multiple other cancers<sup>48–50</sup>. The *CDKN2B-CDKN2A-ANRIL* (*CDKN2B-AS1*) locus at 9p21.3 is a well-established hotspot of pleiotropic GWAS associations for many complex diseases, including coronary artery disease<sup>51</sup>, type 2 diabetes<sup>52</sup> and cancers<sup>50,53–56</sup>. Notably, lead variant rs1537372 is in high LD ( $r^2=0.82$ ) with variants associated with coronary artery disease<sup>51</sup> and endometriosis<sup>57</sup>, but not with the other cancer-associated variants. *CDKN2A* and *CDKN2B* encode cyclin-dependent kinase (CDK) inhibitors that regulate the cell cycle. *CDKN2A* is one of the most commonly inactivated genes in cancer, and is a high penetrance gene for melanoma<sup>58,59</sup>. *CDKN2B* activation is tightly controlled by the cytokine TGF- $\beta$ , further linking this signaling pathway with CRC tumorigenesis<sup>60</sup>.

Our findings implicate genes in pathways with established roles in CRC pathogenesis. We identified loci at *SMAD3* and *SMAD9*, members of the TGF- $\beta$  signaling pathway, which includes genes linked to familial CRC syndromes (for example, *SMAD4* and *BMPRIA*) and several GWAS-implicated genes (for example, *SMAD7*, *BMP2* and *BMP4*)<sup>61</sup>. We identified another locus near TGF- $\beta$  receptor 1 (*TGFBR1*). Nearby gene *GALNT12* harbors inactivating germline and somatic mutations in human colon cancers<sup>62</sup> and therefore could also be the regulated effector gene. We identified a locus at 14q23.1 near *DACT1*, a member of the Wnt- $\beta$ -catenin pathway which includes genes previously linked to familial CRC syndromes (*APC*)<sup>63</sup>, and several GWAS-implicated genes (for example, *CTNBN1*<sup>18</sup> and *TCF7L2*)<sup>17</sup>. Genes related to telomere biology were linked by other GWASs: *TERC*<sup>10</sup> and *TERT*<sup>22</sup>, which encode the RNA and protein subunits of telomerase, respectively, and *FEN1* (ref. 17), which is involved in telomere stability<sup>64</sup>. A new locus at 20q13.33 harbors another gene related to telomere biology, *RTEL1*. This gene is involved in DNA double-strand-break repair, and overlaps with GWAS signals for cancers<sup>55,65</sup> and inflammation-related phenotypes, including inflammatory bowel disease<sup>66</sup> and atopic dermatitis<sup>67</sup>.

Of 61 signals at 56 loci previously associated with CRC at GWS, 42 showed association evidence at  $P<5\times 10^{-8}$  in the combined meta-analysis, and 55 at  $P<0.05$  in the independent stage 2 meta-analysis (Supplementary Table 3). Notably, the association of rs755229494 at locus 5q22.2 ( $P=2.1\times 10^{-12}$ ) was driven by studies of predominantly subjects with Ashkenazi Jewish ancestry, and this SNP is in perfect LD with known missense SNP rs1801155 in the *APC* gene (p.Ile1307Lys), the minor allele of which is enriched in this population (MAF 6%), but rare in other populations<sup>68,69</sup>.

**Delineating distinct association signals at CRC risk loci.** To identify additional independent association signals at known or new CRC risk loci, we conducted conditional analysis using individual-level data of 125,478 participants (Methods). At nine loci, we observed ten new independent association signals that attained a  $P$  value in a joint multiple-variant analysis ( $P_j$ )  $<5\times 10^{-8}$  (Table 2, Supplementary Table 4 and Supplementary Fig. 5). Because this analysis focused on  $<5\%$  of the genome, we also report signals at  $P_j<1\times 10^{-5}$  in Supplementary Table 5. At 22 loci, we observed 25 new suggestive associations with  $P_j<1\times 10^{-5}$ .

At locus 11q13.4, near *POLD3* and *CHRD12*, we identified a new low-frequency variant (lead SNP rs61389091, MAF 3.94%)

**Table 1 | New CRC risk loci reaching genome-wide significance ( $P < 5 \times 10^{-8}$ ) in the combined (stage 1 and stage 2) meta-analysis**

Locus	Nearby gene(s)	rsID lead variant	Chr.	Position (build 37)	Alleles (risk/other)	RAF (%)	Stage 1 meta-analysis: up to 34,869 cases and 29,051 controls			Stage 2 meta-analysis: up to 23,262 cases and 38,296 controls			Combined meta-analysis: up to 58,131 cases and 67,347 controls		
							OR	95% CI	P	OR	95% CI	P	OR	95% CI	P
Rare variants															
5q21.1	<i>RGMB; CHD1</i>	rs145364999 <sup>a</sup>	5	98,206,082	T/A	99.69	1.57	1.20–2.05	$9.0 \times 10^{-4}$	1.93	1.48–2.52	$1.0 \times 10^{-6}$	1.74	1.45–2.10	$6.3 \times 10^{-9}$
Low-frequency variants															
3q13.2	<i>BOC</i>	rs72942485	3	112,999,560	G/A	98.02	1.16	1.07–1.26	$2.5 \times 10^{-4}$	1.23	1.12–1.35	$1.5 \times 10^{-5}$	1.19	1.12–1.26	$2.1 \times 10^{-8}$
Common variants															
1p34.3	<i>FHL3</i>	rs4360494 <sup>b</sup>	1	38,455,891	G/C	45.39	1.05	1.03–1.08	$2.9 \times 10^{-5}$	1.06	1.03–1.08	$3.3 \times 10^{-5}$	1.05	1.04–1.07	$3.8 \times 10^{-9}$
1p32.3	<i>TTC22; PCSK9</i>	rs12144319 <sup>a</sup>	1	55,246,035	C/T	25.48	1.07	1.04–1.10	$1.4 \times 10^{-6}$	1.07	1.04–1.10	$5.5 \times 10^{-6}$	1.07	1.05–1.09	$3.3 \times 10^{-11}$
2q24.2	<i>MARCH7; TANC1</i>	rs448513 <sup>b</sup>	2	159,964,552	C/T	32.60	1.06	1.03–1.08	$1.9 \times 10^{-5}$	1.05	1.02–1.08	$5.8 \times 10^{-4}$	1.05	1.03–1.07	$4.4 \times 10^{-8}$
2q33.1	<i>SATB2</i>	rs983402 <sup>a</sup>	2	199,781,586	T/C	33.12	1.05	1.03–1.08	$7.2 \times 10^{-5}$	1.08	1.05–1.11	$1.0 \times 10^{-8}$	1.07	1.05–1.09	$7.7 \times 10^{-12}$
3q22.2	<i>SLCO2A1</i>	rs10049390 <sup>b</sup>	3	133,701,119	A/G	73.53	1.06	1.03–1.09	$4.9 \times 10^{-5}$	1.07	1.04–1.10	$1.8 \times 10^{-5}$	1.06	1.04–1.08	$3.8 \times 10^{-9}$
4q24	<i>TET2</i>	rs1391441	4	106,128,760	A/G	67.20	1.05	1.02–1.07	$1.5 \times 10^{-4}$	1.06	1.03–1.09	$2.3 \times 10^{-5}$	1.05	1.03–1.07	$1.6 \times 10^{-8}$
4q31.21	<i>HHIP</i>	rs11727676	4	145,659,064	C/T	9.80	1.08	1.03–1.13	$4.5 \times 10^{-4}$	1.10	1.05–1.14	$1.5 \times 10^{-5}$	1.09	1.06–1.12	$2.9 \times 10^{-8}$
6p21.32	<i>HLA-DRB1; HLA-DQA1</i>	rs9271695 <sup>a</sup>	6	32,593,080	G/A	79.54	1.09	1.06–1.13	$1.3 \times 10^{-7}$	1.09	1.05–1.12	$1.7 \times 10^{-7}$	1.09	1.07–1.12	$1.1 \times 10^{-13}$
7p13	<i>MYO1G; SNHG15; CCM2; TBRG4</i>	rs12672022 <sup>b</sup>	7	45,136,423	T/C	83.45	1.07	1.04–1.11	$1.6 \times 10^{-5}$	1.06	1.03–1.10	$4.4 \times 10^{-4}$	1.07	1.04–1.09	$2.8 \times 10^{-8}$
9p21.3	<i>ANRIL; CDKN2A; CDKN2B</i>	rs1537372 <sup>b</sup>	9	22,103,183	G/T	56.92	1.05	1.02–1.07	$1.4 \times 10^{-4}$	1.06	1.03–1.08	$2.4 \times 10^{-5}$	1.05	1.03–1.07	$1.4 \times 10^{-8}$
9q22.33	<i>GALNT12; TGFBRI</i>	rs34405347 <sup>b</sup>	9	101,679,752	T/G	90.34	1.08	1.04–1.13	$5.5 \times 10^{-5}$	1.09	1.04–1.13	$1.5 \times 10^{-4}$	1.09	1.05–1.12	$3.1 \times 10^{-8}$
9q31.3	<i>LPAR1</i>	rs10980628	9	113,671,403	C/T	21.06	1.05	1.02–1.09	$3.1 \times 10^{-4}$	1.08	1.05–1.11	$1.3 \times 10^{-6}$	1.07	1.04–1.09	$2.8 \times 10^{-9}$
11q22.1	<i>YAP1</i>	rs2186607	11	101,656,397	T/A	51.78	1.05	1.03–1.08	$1.1 \times 10^{-5}$	1.05	1.03–1.08	$3.3 \times 10^{-5}$	1.05	1.04–1.07	$1.5 \times 10^{-9}$
12q12	<i>PRICKLE1; YAF2</i>	rs11610543 <sup>b</sup>	12	43,134,191	G/A	50.13	1.05	1.03–1.08	$1.1 \times 10^{-5}$	1.06	1.03–1.08	$2.8 \times 10^{-5}$	1.05	1.04–1.07	$1.3 \times 10^{-9}$
12q13.3	<i>STAT6; LRP1; NAB2</i>	rs4759277	12	57,533,690	A/C	35.46	1.07	1.04–1.09	$8.4 \times 10^{-7}$	1.04	1.02–1.07	$1.6 \times 10^{-3}$	1.05	1.04–1.07	$9.4 \times 10^{-9}$
13q13.3	<i>SMAD9</i>	rs7333607 <sup>a</sup>	13	37,462,010	G/A	23.50	1.09	1.06–1.12	$2.5 \times 10^{-8}$	1.07	1.04–1.10	$4.4 \times 10^{-6}$	1.08	1.06–1.10	$6.3 \times 10^{-13}$
13q22.1	<i>KLF5</i>	rs78341008 <sup>b</sup>	13	73,791,554	C/T	7.19	1.13	1.07–1.18	$1.4 \times 10^{-6}$	1.11	1.05–1.16	$4.8 \times 10^{-5}$	1.12	1.08–1.16	$3.2 \times 10^{-10}$
13q34	<i>COL4A2; COL4A1; RAB20</i>	rs8000189	13	111,075,881	T/C	64.01	1.05	1.02–1.07	$2.1 \times 10^{-4}$	1.07	1.04–1.10	$1.3 \times 10^{-6}$	1.06	1.04–1.08	$1.8 \times 10^{-9}$
14q23.1	<i>DACT1</i>	rs17094983 <sup>b</sup>	14	59,189,361	G/A	87.73	1.10	1.07–1.15	$8.4 \times 10^{-8}$	1.08	1.04–1.12	$9.0 \times 10^{-5}$	1.09	1.06–1.12	$4.6 \times 10^{-11}$
15q22.33	<i>SMAD3</i>	rs56324967 <sup>a</sup>	15	67,402,824	C/T	67.57	1.07	1.04–1.10	$2.2 \times 10^{-7}$	1.08	1.05–1.11	$9.8 \times 10^{-8}$	1.07	1.05–1.09	$1.1 \times 10^{-13}$
16q23.2	<i>MAF</i>	rs9930005 <sup>b</sup>	16	80,043,258	C/A	43.03	1.05	1.03–1.08	$1.3 \times 10^{-5}$	1.05	1.02–1.07	$4.0 \times 10^{-4}$	1.05	1.03–1.07	$2.1 \times 10^{-8}$
17p12	<i>LINC00675</i>	rs1078643 <sup>a</sup>	17	10,707,241	A/G	76.36	1.07	1.04–1.10	$9.2 \times 10^{-6}$	1.09	1.05–1.12	$1.1 \times 10^{-7}$	1.08	1.05–1.10	$6.6 \times 10^{-12}$
17q24.3	<i>LINC00673</i>	rs983318 <sup>a</sup>	17	70,413,253	A/G	25.26	1.07	1.04–1.10	$1.2 \times 10^{-6}$	1.05	1.02–1.08	$8.0 \times 10^{-4}$	1.06	1.04–1.08	$5.6 \times 10^{-9}$
17q25.3	<i>RAB40B; METRNL</i>	rs75954926 <sup>a</sup>	17	81,061,048	G/A	65.68	1.10	1.07–1.13	$9.4 \times 10^{-11}$	1.09	1.06–1.12	$4.8 \times 10^{-9}$	1.09	1.07–1.11	$3.0 \times 10^{-18}$
19p13.11	<i>KLF2</i>	rs34797592 <sup>b</sup>	19	16,417,198	T/C	11.82	1.09	1.05–1.13	$8.2 \times 10^{-6}$	1.09	1.05–1.13	$1.2 \times 10^{-5}$	1.09	1.06–1.12	$4.2 \times 10^{-10}$
19q13.43	<i>TRIM28</i>	rs73068325	19	59,079,096	T/C	18.26	1.06	1.03–1.09	$2.1 \times 10^{-4}$	1.07	1.04–1.11	$5.0 \times 10^{-5}$	1.07	1.04–1.09	$4.2 \times 10^{-8}$
20q13.12	<i>TOX2; HNF4A</i>	rs6031311 <sup>b</sup>	20	42,666,475	T/C	75.91	1.07	1.04–1.10	$1.7 \times 10^{-6}$	1.05	1.02–1.08	$7.6 \times 10^{-4}$	1.06	1.04–1.08	$6.8 \times 10^{-9}$
20q13.33	<i>TNFRSF6B; RTEL1</i>	rs2738783 <sup>b,c</sup>	20	62,308,612	T/G	20.29	1.07	1.04–1.10	$2.6 \times 10^{-6}$	1.05	1.02–1.08	$3.3 \times 10^{-3}$	1.06	1.04–1.08	$5.3 \times 10^{-8}$

Lead variant is the most associated variant at the locus. Reference SNP cluster ID (rsID) based on NCBI dbSNP Build 150. Alleles are on the + strand. Chr., chromosome; RAF, risk allele frequency, based on stage 2 data; OR, odds ratio estimate for the risk allele. All P values reported in this table are based on fixed-effects inverse variance-weighted meta-analysis. <sup>a</sup>Variant or LD proxy ( $r^2 > 0.7$ ) was selected for our custom genotyping panel and formally replicates in the stage 2 meta-analysis at a Bonferroni significance threshold of  $P < 7.8 \times 10^{-6}$ . <sup>b</sup>Variant or LD proxy ( $r^2 > 0.7$ ) was selected for our custom genotyping panel but did not attain Bonferroni significance in the stage 2 meta-analysis. <sup>c</sup>This SNP reached genome-wide significance in the combined (stage 1 + stage 2) sample-size weighted meta-analysis based on likelihood-ratio test results ( $P = 4.9 \times 10^{-9}$ ).

**Table 2 | Additional new conditionally independent association signals at known and newly identified CRC risk loci that reach genome-wide significance ( $P < 5 \times 10^{-8}$ ) in the combined meta-analysis of up to 125,478 individuals**

Locus	Nearby gene(s)	rsID lead variant	Chr.	Position (build 37)	Alleles (risk/other)	RAF (%)	OR <sub>unconditional</sub>	95% CI	P <sub>unconditional</sub>	Joint multiple-variant analysis			
										Conditioning variant(s)	OR <sub>conditional</sub>	95% CI	P <sub>conditional</sub>
Low-frequency variants													
11q13.4	<i>POLD3</i>	rs61389091	11	74,427,921	C/T	96.06	1.23	1.18–1.29	$1.2 \times 10^{-18}$	rs7121958 <sup>a</sup> , rs7946853	1.21	1.16–1.27	$3.7 \times 10^{-16}$
Common variants													
2q33.1	<i>SATB2</i>	rs11884596	2	199,612,407	C/T	38.23	1.06	1.04–1.08	$1.1 \times 10^{-9}$	rs983402	1.06	1.04–1.07	$3.6 \times 10^{-9}$
5p15.33	<i>TERT</i> ; <i>CLPTM1L</i>	rs78368589	5	1,240,204	T/C	5.97	1.14	1.10–1.18	$9.4 \times 10^{-12}$	rs2735940 <sup>a</sup>	1.12	1.08–1.16	$4.1 \times 10^{-9}$
5p13.1	<i>LINC00603</i> ; <i>PTGER4</i>	rs7708610	5	40,102,443	A/G	35.64	1.04	1.02–1.06	$1.5 \times 10^{-5}$	rs12514517 <sup>a</sup>	1.06	1.04–1.08	$3.8 \times 10^{-9}$
6p21.32	<i>HLA-B</i> ; <i>MICA</i> ; <i>MICB</i> ; <i>NFKB1I</i> ; <i>TNF</i>	rs2516420	6	31,449,620	C/T	92.63	1.10	1.06–1.13	$1.3 \times 10^{-7}$	rs9271695, rs116685461, rs116353863	1.12	1.08–1.16	$2.0 \times 10^{-10}$
8q24.21	<i>MYC</i>	rs4313119	8	128,571,855	G/T	74.86	1.06	1.04–1.08	$1.0 \times 10^{-9}$	rs6983267 <sup>a</sup> , rs7013278	1.06	1.04–1.08	$2.1 \times 10^{-9}$
12p13.32	<i>CCND2</i>	rs3217874	12	4,400,808	T/C	42.82	1.08	1.06–1.10	$1.2 \times 10^{-17}$	rs3217810 <sup>a</sup> , rs35808169 <sup>a</sup>	1.06	1.04–1.08	$2.4 \times 10^{-9}$
15q13.3	<i>GREM1</i>	rs17816465	15	33,156,386	A/G	20.55	1.07	1.04–1.09	$6.8 \times 10^{-9}$	rs2293581 <sup>a</sup> , rs12708491 <sup>a</sup>	1.07	1.05–1.10	$1.4 \times 10^{-10}$
20p12.3	<i>BMP2</i>	rs28488	20	6,762,221	T/C	63.88	1.06	1.04–1.08	$2.6 \times 10^{-11}$	rs189583 <sup>a</sup> , rs4813802 <sup>a</sup> , rs994308	1.07	1.05–1.09	$2.6 \times 10^{-14}$
20p12.3	<i>BMP2</i>	rs994308	20	6,603,622	C/T	59.39	1.08	1.06–1.10	$4.8 \times 10^{-18}$	rs189583 <sup>a</sup> , rs4813802 <sup>a</sup> , rs28488	1.06	1.05–1.08	$8.6 \times 10^{-12}$

Lead variant is the most associated variant at the locus in the conditional analysis. Reference SNP cluster ID (rsID) based on NCBI dbSNP Build 150. Alleles are on the + strand. Chr., chromosome; RAF, risk allele frequency, based on stage 2 data; OR, odds ratio estimates are for the risk allele. Conditioning variants are the lead variant of other conditionally independent association signals with  $P < 1 \times 10^{-5}$  within 1 Mb of the new association signal. Because of extensive LD, we used a 2-Mb distance for the MHC region (6p21.32). All lead variants for the new association signals are in linkage equilibrium with any previously reported CRC risk variants at the locus ( $r^2 < 0.10$ ). <sup>a</sup>Conditioning variant is either the index variant, or a variant in LD with the index variant reported in previous GWAS. Details and full results are in Supplementary Table 5.

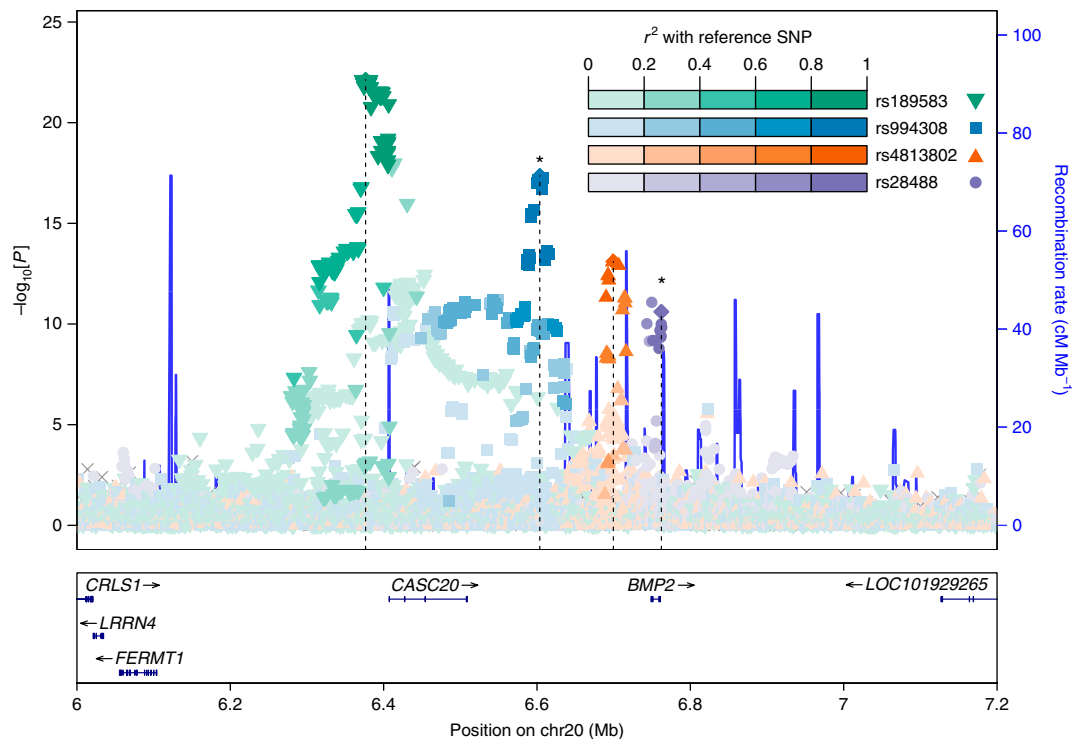
separated by a recombination hotspot from the known common variant signal<sup>12</sup> (LD  $r^2$  between lead SNPs  $< 0.01$ ). At 5p15.33, we identified another lower-frequency variant association (lead SNP rs78368589, MAF 5.97%), which was independent from the previously reported common variant signal 56 kb away, near *TERT* and *CLPTM1L* (LD  $r^2$  with lead SNP rs2735940  $< 0.01$ )<sup>22</sup>. Variants in this region were linked to many cancer types, including lung, prostate, breast and ovarian cancer<sup>70</sup>.

The remaining eight new signals involved common variants. At new locus 2q33.1, near genes *PLCL1* and *SATB2*, two statistically independent associations (LD  $r^2$  between two lead SNPs  $< 0.01$ ) are separated by a recombination hotspot (Supplementary Fig. 5). In the MHC region, we identified a conditionally independent signal near genes involved in NF- $\kappa$ B signaling, including the gene encoding tumor necrosis factor- $\alpha$ , genes for the stress-signaling proteins *MICA* and *MICB*, and *HLA-B*. Locus 20p12.3, near *BMP2*, harbored four distinct association signals (Fig. 1), two of which have been reported<sup>10,11</sup> (Supplementary Table 5). All four SNPs selected in the model were in pairwise linkage equilibrium (maximum LD  $r^2 = 0.039$ , between rs189583 and rs994308). Our conditional analysis further confirmed that the signal ~1-Mb centromeric of *BMP2*, near gene *HAOI*, is independent. At locus 8q24.21 near *MYC*, the locus showing the second strongest statistical evidence of association in the combined meta-analysis (lead SNP rs6983267;  $P = 3.4 \times 10^{-64}$ ), we identified a second independent signal (lead SNP rs4313119,  $P_1 = 2.1 \times 10^{-9}$ ; LD  $r^2$  with rs6983267  $< 0.001$ ). At the recently reported locus 5p13.1 (ref. <sup>22</sup>), near the noncoding

RNA gene *LINC00603*, we identified an additional signal (lead SNP rs7708610) that was partly masked by the reported signal in the single-variant analysis due to the negative correlation between rs7708610 and rs12514517 ( $r = -0.18$ ;  $r^2 = 0.03$ ). This caused significance for both SNPs to increase markedly when fitted jointly (rs7708610, unconditional  $P = 1.5 \times 10^{-5}$  and  $P_1 = 3.8 \times 10^{-9}$ ). At locus 12p13.32 near *CCND2*, we identified a new signal (lead SNP rs3217874,  $P_1 = 2.4 \times 10^{-9}$ ) and confirmed two previously associated signals<sup>13–15</sup> (Supplementary Note). At the *GREM1* locus on 15q13.3, two independent signals have been described<sup>11</sup>. Our analyses suggest that this locus harbors three signals. A new signal represented by SNP rs17816465 is conditionally independent from the other two signals ( $P_1 = 1.4 \times 10^{-10}$ , conditioned on rs2293581 and rs12708491; LD with conditioning SNPs  $r^2 < 0.01$ ; Supplementary Note).

Additionally, signals with  $P_1$  values approaching GWS were observed at new locus 3q13.2 near *BOC* (rs13086367, unconditional  $P = 6.7 \times 10^{-8}$ ,  $P_1 = 6.9 \times 10^{-8}$ , MAF = 47.4%), 96 kb from the low-frequency signal represented by rs72942485 (unconditional  $P = 2.1 \times 10^{-8}$ ,  $P_1 = 1.3 \times 10^{-8}$ , MAF = 2.2%); at known locus 10q22.3 near *ZMIZ1* (rs1250567, unconditional  $P = 3.1 \times 10^{-8}$ ,  $P_1 = 7.2 \times 10^{-8}$ , MAF = 45.1%); and at new locus 13q22.1 near *KLF5* (rs45597035, unconditional  $P = 2.7 \times 10^{-9}$ ,  $P_1 = 8.1 \times 10^{-8}$ , MAF = 34.4%) (Supplementary Table 5). Furthermore, we clarified previously reported independent association signals (Supplementary Note).

**Associations of CRC risk variants with other traits.** Some 19 of the GWS association signals for CRC were in high LD ( $r^2 > 0.7$ )



**Fig. 1 | Conditionally independent association signals at the *BMP2* locus.** Regional association plot showing the unconditional  $-\log_{10}[P]$  for the association with CRC risk in the combined meta-analysis of up to 125,478 individuals, as a function of genomic position (build 37) for each variant in the chromosome 20 (chr20) region. The lead variants are indicated by a diamond symbol, and their positions are indicated by dashed vertical lines. The color labeling and shape of all other variants indicate the lead variant with which they are in strongest LD. The two new genome-wide significant signals are indicated by asterisks.

with at least one SNP in the NHGRI-EBI GWAS Catalog<sup>46</sup> that has significant association in GWASs of other traits. Notable overlap included SNPs associated with other cancers, immune-related traits (for example, tonsillectomy, inflammatory bowel disease and circulating white blood cell traits), obesity traits, blood pressure and other cardiometabolic traits (Supplementary Table 6).

**Mechanisms underlying CRC association signals.** To further localize variants driving the 40 newly identified signals, we used association evidence to define credible sets of variants that are 99% probable to contain the causal variant (Methods). The 99% credible set size for new loci ranged from one (17p12) to 93 (2q33.1). For 11 distinct association signals, the set included ten or fewer variants (Supplementary Table 7). At locus 17p12, we narrowed the candidate variant to rs1078643, located in exon 1 of the lncRNA *LINC00675*, which is primarily expressed in gastrointestinal tissues. Small credible sets were observed for locus 4q31.21 (two variants, indexed by synonymous SNP rs11727676 in *HHIP*), and signals at known loci near *GREM1* (one variant) and *CCND2* (two variants).

We performed functional annotation of credible set variants to nominate putative causal variants. Eight sets contained coding variants, but only the synonymous SNP in *HHIP* had a high posterior probability of driving the association (Supplementary Table 8). Next, we examined overlap of credible sets with regulatory genomic annotations from 51 existing CRC-relevant data sets to examine noncoding functions (Methods). Also, to better refine regulatory elements in active enhancers, we performed the assay for transposase-accessible chromatin using sequencing (ATAC-seq) to measure chromatin accessibility in four colonic crypts and used the resulting data to annotate GWAS signals.

Of the 40 sets, 36 overlapped with active enhancers identified by histone mark H3K27ac, which was measured in normal colonic

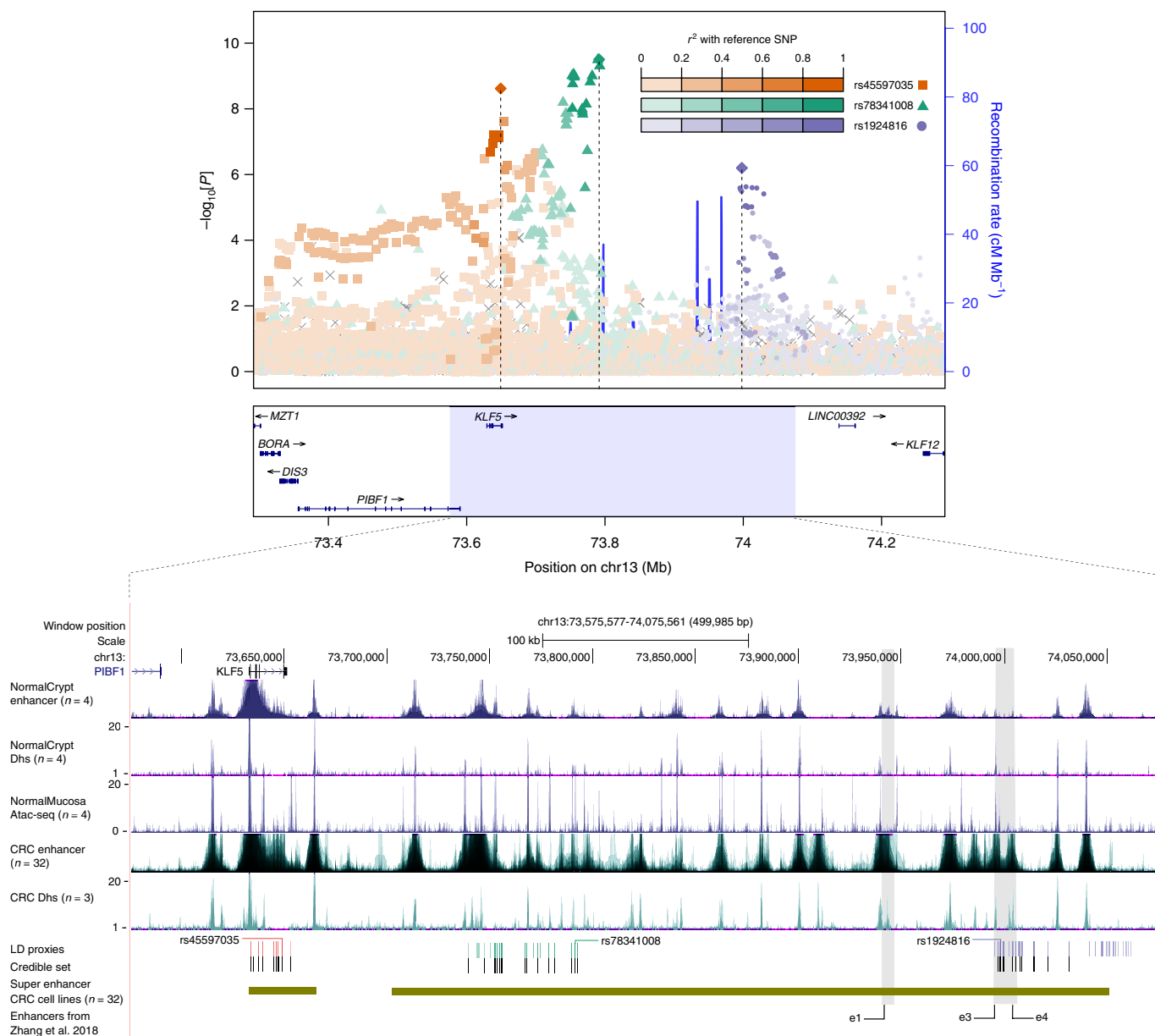
crypt epithelium, CRC cell lines or CRC tissue (Supplementary Table 8 and Supplementary Fig. 6). Twenty of these 36 sets overlapped with super-enhancers. Notably, when compared with epigenomics data from normal colonic crypt epithelium, all 36 sets overlapped with enhancers that gained or lost activity in one or more CRC specimens. Eleven of these sets overlapped with enhancers that recurrently gained or lost activity in  $\geq 20$  CRC cell lines.

The locus at GWAS hotspot 9p21 overlaps with a super-enhancer, and the credible set is entirely intronic to *ANRIL*, alias *CDKN2B-AS1*. The Genotype-Tissue Expression (GTEx) data show that the antisense lncRNA *ANRIL* is exclusively expressed in transverse colon and small intestine. Notably, ANRIL recruits SUZ12 and EH22 to epigenetically silence tumor suppressor genes *CDKN2A* and *CDKN2B*<sup>71</sup>.

Noncoding somatic driver mutations or focal amplifications have been reported in regions regulating expression of *MYC*<sup>72</sup>, *TERT*<sup>73</sup> and *KLF5* (ref. <sup>31</sup>), which are now implicated by GWAS for CRC. We checked whether GWAS-identified association signals colocalize with these regions and found that the *KLF5* signal overlaps with the somatically amplified super-enhancer flanked by *KLF5* and *KLF12* (Fig. 2). Also, the previously reported signal in the *TERT* promoter region<sup>22</sup> overlaps with the recurrent somatically mutated region in multiple cancers<sup>73</sup>.

To test whether CRC associations are nonrandomly distributed across genomic features, we used GARFIELD<sup>74</sup>. Focusing on DNase I hypersensitive site peaks that identify open chromatin, we observed significant enrichment across many cell types, particularly fetal tissues, with the strongest enrichment observed in fetal gastrointestinal tissues, CD20<sup>+</sup> primary cells (B cells) and embryonic stem cells (Supplementary Fig. 7 and Supplementary Table 9).

We used MAGENTA<sup>75</sup> to identify pathways or gene sets enriched for associations with CRC, assessing two gene *P*-value cutoffs: the

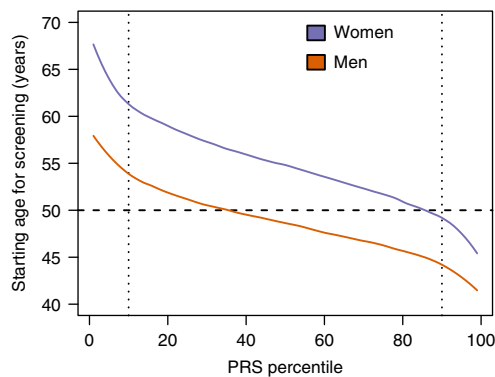


**Fig. 2 | Functional genomic annotation of new CRC risk locus overlapping with *KLF5* super-enhancer.** Top, regional association plot showing the unconditional  $-\log_{10}[P]$  for the association with CRC risk in the combined meta-analysis of up to 125,478 individuals, as a function of genomic position (build 37) for each variant in the chromosome 13 (chr13) region. Lead variants are indicated by diamonds and their positions are indicated by dashed vertical lines. The color labeling and shape of all other variants indicate the lead variant with which they are in strongest LD. Bottom, UCSC genome browser annotations for region overlapping with the super-enhancer flanked by *KLF5* and *KLF12*, spanning variants in LD with rs78341008, and with two conditionally independent association signals indexed by rs45597035 and rs1924816. The region is annotated with the following tracks (top to bottom): UCSC gene annotations; epigenetic profiles showing MACS2 peak calls as transparent overlays for different samples taken from nondiseased colonic crypt cells (NormalCrypt) or colon mucosa tissue (NormalMucosa) (purple) and from different primary CRC cell lines or tumor samples (teal); position of the lead variants and variants in LD with the lead; variants in the 99% credible set; the union of super-enhancers called using the ROSE package; targeted enhancers (e1, e3 and e4) shown by Zhang *et al.*<sup>31</sup> to have combinatorial effects on *KLF5* expression (gray bars). ATAC-seq data newly generated for this study show high-resolution annotation of putative binding regions within the active super-enhancer, further fine-mapping putative causal variants at each of the three signals. Dhs: DNase I hypersensitive sites.

95th and 75th percentiles. At the 75th percentile, we observed enrichment of multiple KEGG cancer pathways at a false discovery rate (FDR) of 0.05. This was not observed for the 95th percentile cutoff and suggests that many more loci that are shared with other cancer types remain to be identified in larger studies. Using the 75th (95th) percentile cutoff, at FDR 0.05 and 0.20, we found enrichment of 7 (5) and 53 (24) gene sets, respectively. Established pathways related to TGF- $\beta$ /SMAD and BMP signaling were among the top

enriched pathways. Other notable enriched pathways included Hedgehog signaling, basal cell carcinoma, melanogenesis, cell cycle, S phase and telomere maintenance (Supplementary Table 10).

**Polygenicity of CRC and contribution of rare variants.** To estimate the contribution of rare variants ( $MAF \leq 1\%$ ) to CRC heritability, we used the LD- and MAF-stratified component GREML (GREML-LDMS) method implemented in GCTA<sup>76</sup> (Methods). Assuming a



**Fig. 3 | Recommended age to start CRC screening based on a polygenic risk score.** The PRS was constructed using the 95 known and newly discovered variants. The horizontal line represents the recommended age for first endoscopy for a person with average risk using the current screening guidelines for CRC. The risk threshold to determine the age for the first screening was set as the average of 10-year CRC risks for a 50-year-old man (1.25%) and woman (0.68%) who have not previously received an endoscopy. Details are in Methods.

lifetime risk of 4.3%, we estimated that all imputed autosomal variants explain 21.6% (95% CI = 17.5–25.7%) of the variation in liability for CRC, with almost half of this contributed by rare variants ( $h_g^2 = 9.7\%$ , 95% CI = 6.2–13.3%; likelihood ratio test  $P = 0.003$ ); the estimated liability-scale heritability for variants with MAF > 1% is 11.8% (95% CI = 8.9–14.7%). Our overall estimate falls within the range of heritability reported by large twin studies<sup>2</sup>. Because heritability estimates for rare variants are sensitive to potential biases due to technical effects or population stratification<sup>77</sup>, and the contribution of rare variants is probably underestimated due to limitations of genotype imputation, the results should be interpreted with caution. Overall, the findings suggest that missing heritability is not large, but that many rare and common variants have yet to be identified.

**Familial relative risk explained by GWAS-identified variants.** Adjusting for the winner's curse<sup>78</sup>, the familial relative risk (RR) to first-degree relatives ( $\lambda_0$ ) attributable to GWAS-identified variants rose from 1.072 for the 55 previously described autosomal risk variants that showed evidence for replication at  $P < 0.05$ , to 1.092 after inclusion of 40 new signals, and increased further to 1.098 when we included 25 suggestive association signals reported in Supplementary Table 5 (Methods). Assuming a  $\lambda_0$  of 2.2, the 55 established signals account for 8.8% of familial RR explained (95% CI = 8.1–9.4). Established signals combined with 40 newly discovered signals account for 11.2% (95% CI = 10.5–12.0), and adding 25 suggestive signals increases this to 11.9% (95% CI = 11.1–12.7).

**Implications for stratified screening prevention.** We demonstrate how using a polygenic risk score (PRS) derived from 95 independent association signals could affect clinical guidelines for preventive screening. The difference in recommended starting age for screening for those in the highest 1% (and 10%) percentiles of risk compared with lowest percentiles is 18 years (and 10 years) for men, and 24 years (and 12 years) for women (Fig. 3; Methods). Supplementary Table 11 gives risk allele frequency estimates in different populations for variants included in the PRS. As expected, risk allele frequencies vary across populations. Furthermore, differences in LD between tagging and true causal variants across populations can lead to less prediction accuracy and subsequent lower predictive power of the PRS in non-European populations. Accordingly, it will be important to develop ancestry-specific PRSs that incorporate detailed fine-mapping results for each GWAS signal.

## Discussion

To further define the genetic architecture of sporadic CRC, we performed low-coverage WGS and imputation into a large set of GWAS data. We discovered 40 new CRC signals and replicated 55 previously reported signals. We found the first rare variant signal for sporadic CRC, which represents the strongest protective rare allelic effect identified so far. Our analyses highlight new genes and pathways contributing to underlying CRC risk and suggest roles for Krüppel-like factors, Hedgehog signaling, Hippo-YAP signaling and immune function. Multiple loci provide new evidence for an important role of lncRNAs in CRC tumorigenesis<sup>79</sup>. Functional genomic annotations support that most of the sporadic CRC genetic risk lies in noncoding genomic regions. We further show how newly discovered variants can lead to improved risk prediction.

This study underscores the critical importance of large-scale GWAS collaboration. Although discovery of the rare variant signal was possible only through increased coverage and improved imputation accuracy enabled by imputation panels, sample size was pivotal for discovery of new CRC loci. The results suggest that CRC has a highly polygenic architecture, much of which remains undefined. This also suggests that continued GWAS efforts, together with increasingly comprehensive imputation panels that allow for improved low-frequency and rare genetic variant imputation, will uncover more CRC risk variants. In addition, to investigate sites that are not imputable, large-scale deep sequencing will be needed. Importantly, the prevailing European bias in CRC GWASs limits the generalizability of findings and the application of PRSs in non-European (especially African) populations<sup>80</sup>. Therefore, a broader representation of ancestries in CRC GWASs is necessary.

Studies of somatic genomic alterations in cancer have mostly focused on the coding genome, and identification of noncoding drivers has been challenging<sup>73</sup>. Yet, noncoding somatic driver mutations or focal amplifications in regulatory regions affecting expression have been reported for *MYC*<sup>72</sup>, *TERT*<sup>73</sup> and *KLF5* (ref. 31). The observed overlap between GWAS-identified CRC risk loci and somatic driver regions strongly suggests that expanding the search of somatic driver mutations to noncoding regulatory elements will yield additional discoveries and that searches for somatic drivers can be guided by GWAS findings.

Additionally, we found loci near proposed drug targets, including *CHD1*, which is implicated by the rare variant signal, and *KLF5*. So far, cancer drug target discovery research has almost exclusively focused on properties of cancer cells, yielding drugs targeting proteins that are either highly expressed or expressed in a mutant form due to frequent recurrent somatic missense mutations (for example, *BRAF* p.Val600Glu) or gene fusion events. In stark contrast with other common complex diseases, cancer GWAS results are not being used extensively to inform drug target selection. It has been estimated that selecting targets supported by GWAS could double the success rate in clinical development<sup>81</sup>. Our discoveries corroborate that GWAS results can considerably inform drug discovery, not only for treating cancers but also for chemoprevention in high-risk individuals.

In summary, in the largest genome-wide scan for sporadic CRC risk thus far, we identified the first rare variant signal for sporadic CRC, and almost doubled the number of known association signals. Our findings provide a substantial number of new leads that may spur downstream investigation into the biology of CRC risk, and that will influence drug development and clinical guidelines, such as personalized screening decisions.

**URLs.** EPACTS, <https://genome.sph.umich.edu/wiki/EPACTS>; ENCODE data processing pipeline, [https://github.com/kundajelab/atac\\_dnase\\_pipelines](https://github.com/kundajelab/atac_dnase_pipelines); WashU Epigenome Browser, <https://epigenomegateway.wustl.edu>.

## Online content

Any methods, additional references, Nature Research reporting summaries, source data, statements of data availability and associated accession codes are available at <https://doi.org/10.1038/s41588-018-0286-6>.

Received: 11 May 2018; Accepted: 22 October 2018;  
Published online: 3 December 2018

## References

- Ferlay, J. et al. Cancer incidence and mortality worldwide: sources, methods and major patterns in GLOBOCAN 2012. *Int. J. Cancer* **136**, E359–E386 (2015).
- Lichtenstein, P. 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.* **343**, 78–85 (2000).
- Czene, K., Lichtenstein, P. & Hemminki, K. Environmental and heritable causes of cancer among 9.6 million individuals in the Swedish Family-Cancer Database. *Int. J. Cancer* **99**, 260–266 (2002).
- Sud, A., Kinnersley, B. & Houlston, R. S. Genome-wide association studies of cancer: current insights and future perspectives. *Nat. Rev. Cancer* **17**, 692–704 (2017).
- Tomlinson, I. et al. A genome-wide association scan of tag SNPs identifies a susceptibility variant for colorectal cancer at 8q24.21. *Nat. Genet.* **39**, 984–988 (2007).
- Broderick, P. et al. A genome-wide association study shows that common alleles of *SMAD7* influence colorectal cancer risk. *Nat. Genet.* **39**, 1315–1317 (2007).
- Tomlinson, I. P. M. et al. A genome-wide association study identifies colorectal cancer susceptibility loci on chromosomes 10p14 and 8q23.3. *Nat. Genet.* **40**, 623–630 (2008).
- Tenesa, A. et al. Genome-wide association scan identifies a colorectal cancer susceptibility locus on 11q23 and replicates risk loci at 8q24 and 18q21. *Nat. Genet.* **40**, 631–637 (2008).
- COGENT Study et al. Meta-analysis of genome-wide association data identifies four new susceptibility loci for colorectal cancer. *Nat. Genet.* **40**, 1426–1435 (2008).
- Houlston, R. S. 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.* **42**, 973–977 (2010).
- Tomlinson, I. P. M. 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.* **7**, e1002105 (2011).
- Dunlop, M. G. et al. Common variation near *CDKN1A*, *POLD3* and *SHROOM2* influences colorectal cancer risk. *Nat. Genet.* **44**, 770–776 (2012).
- Peters, U. et al. Identification of genetic susceptibility loci for colorectal tumors in a genome-wide meta-analysis. *Gastroenterology* **144**, 799–807.e24 (2013).
- Jia, W.-H. et al. Genome-wide association analyses in East Asians identify new susceptibility loci for colorectal cancer. *Nat. Genet.* **45**, 191–196 (2013).
- Whiffin, N. et al. Identification of susceptibility loci for colorectal cancer in a genome-wide meta-analysis. *Hum. Mol. Genet.* **23**, 4729–4737 (2014).
- Wang, H. et al. Trans-ethnic genome-wide association study of colorectal cancer identifies a new susceptibility locus in *VTG1A*. *Nat. Commun.* **5**, 4613 (2014).
- Zhang, B. et al. Large-scale genetic study in East Asians identifies six new loci associated with colorectal cancer risk. *Nat. Genet.* **46**, 533–542 (2014).
- Schumacher, F. R. et al. Genome-wide association study of colorectal cancer identifies six new susceptibility loci. *Nat. Commun.* **6**, 7138 (2015).
- Al-Tassan, N. A. et al. A new GWAS and meta-analysis with 1000Genomes imputation identifies novel risk variants for colorectal cancer. *Sci. Rep.* **5**, 10442 (2015).
- Orlando, G. et al. Variation at 2q35 (*PNKD* and *TMBIM1*) influences colorectal cancer risk and identifies a pleiotropic effect with inflammatory bowel disease. *Hum. Mol. Genet.* **25**, 2349–2359 (2016).
- Zeng, C. et al. Identification of susceptibility loci and genes for colorectal cancer risk. *Gastroenterology* **150**, 1633–1645 (2016).
- Schmit, S. L. et al. Novel common genetic susceptibility loci for colorectal cancer. *J. Natl. Cancer Inst.* <https://doi.org/10.1093/jnci/djy099> (2018).
- Fuchsberger, C. et al. The genetic architecture of type 2 diabetes. *Nature* **536**, 41–47 (2016).
- 1000 Genomes Project Consortium. et al. A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
- McCarthy, S. et al. A reference panel of 64,976 haplotypes for genotype imputation. *Nat. Genet.* **48**, 1279–1283 (2016).
- Amos, C. I. et al. The Oncoarray Consortium: a network for understanding the genetic architecture of common cancers. *Cancer Epidemiol. Biomarkers. Prev.* **26**, 126–135 (2017).
- Zhao, D. & DePinho, R. A. Synthetic essentiality: Targeting tumor suppressor deficiencies in cancer. *Bioessays* **39**, (2017).
- Zhao, D. et al. Synthetic essentiality of chromatin remodelling factor CHD1 in PTEN-deficient cancer. *Nature* **542**, 484–488 (2017).
- Xiao, Y. et al. RGMb is a novel binding partner for PD-L2 and its engagement with PD-L2 promotes respiratory tolerance. *J. Exp. Med.* **211**, 943–959 (2014).
- Topalian, S. L. et al. Safety, activity, and immune correlates of anti-PD-1 antibody in cancer. *N. Engl. J. Med.* **366**, 2443–2454 (2012).
- Zhang, X. et al. Somatic superenhancer duplications and hotspot mutations lead to oncogenic activation of the KLF5 transcription factor. *Cancer Discov.* **8**, 108–125 (2018).
- Giannakis, M. et al. Genomic correlates of immune-cell infiltrates in colorectal carcinoma. *Cell Rep.* **15**, 857–865 (2016).
- Dekker, R. J. et al. KLF2 provokes a gene expression pattern that establishes functional quiescent differentiation of the endothelium. *Blood* **107**, 4354–4363 (2006).
- Boon, R. A. et al. KLF2 suppresses TGF-beta signaling in endothelium through induction of Smad7 and inhibition of AP-1. *Arterioscler. Thromb. Vasc. Biol.* **27**, 532–539 (2007).
- Chakraborty, D. et al. Dopamine stabilizes tumor blood vessels by up-regulating angiopoietin 1 expression in pericytes and Kruppel-like factor-2 expression in tumor endothelial cells. *Proc. Natl Acad. Sci. USA* **108**, 20730–20735 (2011).
- Lee, S.-J. et al. Regulation of hypoxia-inducible factor 1 $\alpha$  (HIF-1 $\alpha$ ) by lysophosphatidic acid is dependent on interplay between p53 and Krüppel-like factor 5. *J. Biol. Chem.* **288**, 25244–25253 (2013).
- Zhang, H. et al. Lysophosphatidic acid facilitates proliferation of colon cancer cells via induction of Krüppel-like factor 5. *J. Biol. Chem.* **282**, 15541–15549 (2007).
- Ma, Z. et al. Long non-coding RNA SNHG15 inhibits P15 and KLF2 expression to promote pancreatic cancer proliferation through EZH2-mediated H3K27me3. *Oncotarget* **8**, 84153–84167 (2017).
- Evangelista, M., Tian, H. & de Sauvage, F. J. The hedgehog signaling pathway in cancer. *Clin. Cancer Res.* **12**, 5924–5928 (2006).
- Gerling, M. et al. Stromal Hedgehog signalling is downregulated in colon cancer and its restoration restrains tumour growth. *Nat. Commun.* **7**, 12321 (2016).
- Mille, F. et al. The Shh receptor Boc promotes progression of early medulloblastoma to advanced tumors. *Dev. Cell.* **31**, 34–47 (2014).
- Mathew, E. et al. Dosage-dependent regulation of pancreatic cancer growth and angiogenesis by hedgehog signaling. *Cell Rep.* **9**, 484–494 (2014).
- Zhao, B., Li, L., Lei, Q. & Guan, K.-L. The Hippo-YAP pathway in organ size control and tumorigenesis: an updated version. *Genes Dev.* **24**, 862–874 (2010).
- Camargo, F. D. et al. YAP1 increases organ size and expands undifferentiated progenitor cells. *Curr. Biol.* **17**, 2054–2060 (2007).
- Ma, X., Zhang, H., Xue, X. & Shah, Y. M. Hypoxia-inducible factor 2 $\alpha$  (HIF-2 $\alpha$ ) promotes colon cancer growth by potentiating Yes-associated protein 1 (YAP1) activity. *J. Biol. Chem.* **292**, 17046–17056 (2017).
- MacArthur, J. et al. The new NHGRI-EBI Catalog of published genome-wide association studies (GWAS Catalog). *Nucleic Acids Res.* **45**, D896–D901 (2017).
- Seshagiri, S. et al. Recurrent R-spondin fusions in colon cancer. *Nature* **488**, 660–664 (2012).
- Song, F. et al. Identification of a melanoma susceptibility locus and somatic mutation in *TET2*. *Carcinogenesis* **35**, 2097–2101 (2014).
- Eeles, R. A. et al. Identification of seven new prostate cancer susceptibility loci through a genome-wide association study. *Nat. Genet.* **41**, 1116–1121 (2009).
- Michailidou, K. et al. Association analysis identifies 65 new breast cancer risk loci. *Nature* **551**, 92–94 (2017).
- Schunkert, H. et al. Large-scale association analysis identifies 13 new susceptibility loci for coronary artery disease. *Nat. Genet.* **43**, 333–338 (2011).
- Scott, L. J. et al. A genome-wide association study of type 2 diabetes in Finns detects multiple susceptibility variants. *Science* **316**, 1341–1345 (2007).
- Al Olama, A. A. et al. A meta-analysis of 87,040 individuals identifies 23 new susceptibility loci for prostate cancer. *Nat. Genet.* **46**, 1103–1109 (2014).
- Timofeeva, M. N. et al. Influence of common genetic variation on lung cancer risk: meta-analysis of 14 900 cases and 29 485 controls. *Hum. Mol. Genet.* **21**, 4980–4995 (2012).
- Shete, S. et al. Genome-wide association study identifies five susceptibility loci for glioma. *Nat. Genet.* **41**, 899–904 (2009).



56. Bishop, D. T. et al. Genome-wide association study identifies three loci associated with melanoma risk. *Nat. Genet.* **41**, 920–925 (2009).
57. Sapkota, Y. et al. Meta-analysis identifies five novel loci associated with endometriosis highlighting key genes involved in hormone metabolism. *Nat. Commun.* **8**, 15539 (2017).
58. Cannon-Albright, L. A. et al. Assignment of a locus for familial melanoma, MLM, to chromosome 9p13-p22. *Science* **258**, 1148–1152 (1992).
59. Hussussian, C. J. et al. Germline p16 mutations in familial melanoma. *Nat. Genet.* **8**, 15–21 (1994).
60. Seoane, J. et al. TGFbeta influences Myc, Miz-1 and Smad to control the CDK inhibitor p15INK4b. *Nat. Cell Biol.* **3**, 400–408 (2001).
61. Jung, B., Staudacher, J. J. & Beauchamp, D. Transforming growth factor  $\beta$  superfamily signaling in development of colorectal cancer. *Gastroenterology* **152**, 36–52 (2017).
62. Guda, K. et al. Inactivating germ-line and somatic mutations in polypeptide N-acetylgalactosaminyltransferase 12 in human colon cancers. *Proc. Natl Acad. Sci. USA* **106**, 12921–12925 (2009).
63. Groden, J. et al. Identification and characterization of the familial adenomatous polyposis coli gene. *Cell* **66**, 589–600 (1991).
64. Saharia, A. et al. FEN1 ensures telomere stability by facilitating replication fork re-initiation. *J. Biol. Chem.* **285**, 27057–27066 (2010).
65. Eeles, R. A. et al. Identification of 23 new prostate cancer susceptibility loci using the iCOGS custom genotyping array. *Nat. Genet.* **45**, 385–391 (2013).
66. Liu, J. Z. et al. Association analyses identify 38 susceptibility loci for inflammatory bowel disease and highlight shared genetic risk across populations. *Nat. Genet.* **47**, 979–986 (2015).
67. Paternoster, L. et al. Multi-ancestry genome-wide association study of 21,000 cases and 95,000 controls identifies new risk loci for atopic dermatitis. *Nat. Genet.* **47**, 1449–1456 (2015).
68. Laken, S. J. et al. Familial colorectal cancer in Ashkenazim due to a hypermutable tract in APC. *Nat. Genet.* **17**, 79–83 (1997).
69. Niell, B. L., Long, J. C., Rennert, G. & Gruber, S. B. Genetic anthropology of the colorectal cancer-susceptibility allele APC I1307K: evidence of genetic drift within the Ashkenazim. *Am. J. Hum. Genet.* **73**, 1250–1260 (2003).
70. Karami, S. et al. Telomere structure and maintenance gene variants and risk of five cancer types. *Int. J. Cancer* **139**, 2655–2670 (2016).
71. Congrains, A., Kamide, K., Ohishi, M. & Rakugi, H. ANRIL: molecular mechanisms and implications in human health. *Int. J. Mol. Sci.* **14**, 1278–1292 (2013).
72. Zhang, X. et al. Identification of focally amplified lineage-specific super-enhancers in human epithelial cancers. *Nat. Genet.* **48**, 176–182 (2016).
73. Rheinbay, E. et al. Discovery and characterization of coding and non-coding driver mutations in more than 2,500 whole cancer genomes. Preprint at <https://www.biorxiv.org/content/early/2017/12/23/237313> (2017).
74. Iotchkova, V. et al. GARFIELD - GWAS analysis of regulatory or functional information enrichment with LD correction. Preprint at <https://www.biorxiv.org/content/early/2016/11/07/085738> (2016).
75. Segrè, A. V. et al. Common inherited variation in mitochondrial genes is not enriched for associations with type 2 diabetes or related glycemic traits. *PLoS Genet.* **6**, e1001058 (2010).
76. Yang, J. et al. Genetic variance estimation with imputed variants finds negligible missing heritability for human height and body mass index. *Nat. Genet.* **47**, 1114–1120 (2015).
77. Bhatia, G. et al. Subtle stratification confounds estimates of heritability from rare variants. Preprint at <https://www.biorxiv.org/content/early/2016/04/12/048181> (2016).
78. Zhong, H. & Prentice, R. L. Bias-reduced estimators and confidence intervals for odds ratios in genome-wide association studies. *Biostatistics* **9**, 621–634 (2008).
79. Cheetham, S. W., Gruhl, F., Mattick, J. S. & Dinger, M. E. Long noncoding RNAs and the genetics of cancer. *Br. J. Cancer* **108**, 2419–2425 (2013).
80. Popejoy, A. B. & Fullerton, S. M. Genomics is failing on diversity. *Nature* **538**, 161–164 (2016).
81. Nelson, M. R. et al. The support of human genetic evidence for approved drug indications. *Nat. Genet.* **47**, 856–860 (2015).

## Acknowledgements

A full list of acknowledgements appears in the Supplementary Note.

## Author contributions

J.R.H., S.A.B., T.A.H., H.M.K., D.V.C., M.W., F.R.S., J.D.S., D.A., M.H.A., K.A., C.A.-C., V.A., C.B., J.A.B., S.I.B., S.B., D.T.B., J.B., H.Boeing, H.Brenner, S.Brezina, S.Buch, D.D.B., A.B.-H., K.B., B.J.C., P.T.C., S.C.-B., A.T.C., J.C.-C., S.J.C., M.-D.C., S.H.C., A.J.C., K.C., A.d.l.C., D.F.E., S.G.E., F.E., D.R.E., E.J.M.F., J.C.F., D.F., S.G., G.G.G., E.G., P.J.G., J.S.G., A.G., M.J.G., R.W.H., J.H., H.H., R.B.H., P.H., M.H., J.L.H., W.-Y.H., T.J.H., D.J.H., R.J., E.J.J., M.A.J., T.O.K., T.J.K., H.R.K., L.N.K., C.K., S.K., S.-S.K., L.L.M., S.C.L., C.I.L., L.L., A.L., N.M.L., S.M., S.D.M., V.M., G.M., M.M., R.L.M., L.M., R.M., A.N., P.A.N., K.O., N.C.O.-M., B.P., P.S.P., R.P., V.P., P.D.P.P., E.A.P., R.L.P., G.R., H.S.R., E.R., M.R.-B., C.S., R.E.S., D.S., M.-H.S., S.S., M.L.S., C.M.T., S.N.T., A.T., C.M.U., F.J.B.v.D., B.V.G., H.v.K., J.V., K.V., P.V., L.V., V.V., E.W., C.R.W., A.W., M.O.W., A.H.W., B.W.Z., W.Z., P.C.S., J.D.P., M.C.B., G.C., V.M., G.R.A., D.A.N., S.B.G., L.H. and U.P. conceived and designed the experiments. T.A.H., M.W., J.D.S., K.F.D., D.D., R.I., E.K., H.L., C.E.M., E.P., J.R., T.S., S.S.T., D.J.V.D.B., M.C.B. and D.A.N. performed the experiments. J.R.H., H.M.K., S.C., S.L.S., D.V.C., C.Q., J.J., C.K.E., P.G., F.R.S., D.M.L., S.C.N., N.A.S.-A., C.A.L., M.L., T.L.L., Y.-R.S., A.K., G.R.A. and L.H. performed statistical analysis. J.R.H., S.A.B., T.A.H., H.M.K., S.C., S.L.S., D.V.C., C.Q., J.J., C.K.E., P.G., M.W., F.R.S., D.M.L., S.C.N., N.A.S.-A., B.L.B., C.S.C., C.M.C., K.R.C., J.G., W.-L.H., C.A.L., S.M.L., M.L., Y.L., T.L.L., M.S., Y.-R.S., A.K., G.R.A., L.H. and U.P. analyzed the data. H.M.K., C.K.E., D.A., M.H.A., K.A., C.A.-C., V.A., C.B., J.A.B., S.I.B., S.B., D.T.B., J.B., H.Boeing, H.Brenner, S.Brezina, S.Buch, D.D.B., A.B.-H., K.B., B.J.C., P.T.C., S.C.-B., A.T.C., J.C.-C., S.J.C., M.-D.C., S.H.C., A.J.C., K.C., A.d.l.C., D.F.E., S.G.E., F.E., D.R.E., E.J.M.F., J.C.F., R.F., L.M.F., D.F., M.G., S.G., W.J.G., G.G.G., P.J.G., W.M.G., J.S.G., A.G., M.J.G., R.W.H., J.H., H.H., S.H., R.B.H., P.H., M.H., J.L.H., W.-Y.H., T.J.H., D.J.H., G.I.-S., G.E.I., R.J., E.J.J., M.A.J., A.D.J., C.E.J., T.O.K., T.J.K., H.R.K., L.N.K., C.K., T.K., S.K., S.-S.K., S.C.L., L.L.M., S.C.L., F.L., C.I.L., L.L., W.L., A.L., N.M.L., S.M., S.D.M., V.M., G.M., M.M., R.L.M., L.M., N.M., R.M., A.N., P.A.N., K.O., S.O., N.C.O.-M., B.P., P.S.P., R.P., V.P., P.D.P.P., M.P., E.A.P., R.L.P., L.R., G.R., H.S.R., E.R., M.R.-B., L.C.S., C.S., R.E.S., M.S., M.-H.S., K.S., S.S., M.L.S., M.C.S., Z.K.S., C.S., C.M.T., S.N.T., D.C.T., A.E.T., A.T., C.M.U., F.J.B.v.D., B.V.G., H.v.K., J.V., K.V., P.V., L.V., V.V., K.W., S.J.W., E.W., A.K.W., C.R.W., A.W., M.O.W., A.H.W., S.H.Z., B.W.Z., Q.Z., W.Z., P.C.S., J.D.P., M.C.B., A.K., G.C., V.M., G.R.A., S.B.G. and U.P. contributed reagents, materials and analysis tools. J.R.H., S.A.B., T.A.H., J.J., L.H. and U.P. wrote the paper.

## Competing interests

G.R.A. has received compensation from 23andMe and Helix. He is currently an employee of Regeneron Pharmaceuticals. H.H. performs collaborative research with Ambray Genetics, InVitaie Genetics, and Myriad Genetic Laboratories, is on the scientific advisory board for InVitaie Genetics and Genome Medical, and has stock in Genome Medical. R.P. has participated in collaborative funded research with Myriad Genetics Laboratories and InVitaie Genetics but has no financial competitive interest.

## Additional information

**Supplementary information** is available for this paper at <https://doi.org/10.1038/s41588-018-0286-6>.

**Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints).

**Correspondence and requests for materials** should be addressed to U.P.

**Publisher's note:** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

© This is a U.S. government work and not under copyright protection in the U.S.; foreign copyright protection may apply 2018

Jeroen R. Huyghe <sup>1,140</sup>, Stephanie A. Bien <sup>1,140</sup>, Tabitha A. Harrison <sup>1,140</sup>, Hyun Min Kang<sup>2</sup>, Sai Chen <sup>2</sup>, Stephanie L. Schmit <sup>3</sup>, David V. Conti<sup>4</sup>, Conghui Qu<sup>1</sup>, Jihyoun Jeon <sup>5</sup>, Christopher K. Edlund<sup>4</sup>, Peyton Greenside<sup>6</sup>, Michael Wainberg<sup>7</sup>, Fredrick R. Schumacher <sup>8</sup>, Joshua D. Smith<sup>9</sup>, David M. Levine<sup>10</sup>, Sarah C. Nelson <sup>10</sup>, Nasa A. Sinnott-Armstrong <sup>11</sup>, Demetrius Albanes<sup>12</sup>, M. Henar Alonso <sup>13,14,15</sup>, Kristin Anderson<sup>16</sup>, Coral Arnau-Collell<sup>17</sup>, Volker Arndt <sup>18</sup>, Christina Bamia<sup>19,20</sup>, Barbara L. Banbury<sup>1</sup>, John A. Baron<sup>21</sup>, Sonja I. Berndt<sup>12</sup>, Stéphane Bézieau <sup>22</sup>, D. Timothy Bishop <sup>23</sup>, Juergen Boehm<sup>24</sup>, Heiner Boeing<sup>25</sup>, Hermann Brenner<sup>18,26,27</sup>, Stefanie Brezina <sup>28</sup>, Stephan Buch<sup>29</sup>, Daniel D. Buchanan <sup>30,31,32</sup>, Andrea Burnett-Hartman<sup>33</sup>, Katja Butterbach<sup>18</sup>, Bette J. Caan<sup>34</sup>, Peter T. Campbell<sup>35</sup>, Christopher S. Carlson<sup>1,36</sup>, Sergi Castellví-Bel<sup>17</sup>, Andrew T. Chan<sup>37,38,39,40,41,42</sup>, Jenny Chang-Claude<sup>43,44</sup>, Stephen J. Chanock <sup>12</sup>, Maria-Dolores Chirlaque<sup>14,45</sup>, Sang Hee Cho<sup>46</sup>, Charles M. Connolly<sup>1</sup>, Amanda J. Cross<sup>47,48</sup>, Katarina Cuk<sup>18</sup>, Keith R. Curtis<sup>1</sup>, Albert de la Chapelle <sup>49</sup>, Kimberly F. Doheny<sup>50</sup>, David Duggan<sup>51</sup>, Douglas F. Easton <sup>52,53</sup>, Sjoerd G. Elias<sup>54</sup>, Faye Elliott<sup>23</sup>, Dallas R. English<sup>55,56</sup>, Edith J. M. Feskens<sup>57</sup>, Jane C. Figueiredo<sup>58,59</sup>, Rocky Fischer<sup>60</sup>, Liesel M. FitzGerald<sup>56,61</sup>, David Forman<sup>62</sup>, Manish Gala<sup>37,39</sup>, Steven Gallinger<sup>63</sup>, W. James Gauderman<sup>4</sup>, Graham G. Giles<sup>55,56</sup>, Elizabeth Gillanders<sup>64</sup>, Jian Gong<sup>1</sup>, Phyllis J. Goodman<sup>65</sup>, William M. Grady<sup>66</sup>, John S. Grove<sup>67</sup>, Andrea Gsur<sup>28</sup>, Marc J. Gunter<sup>68</sup>, Robert W. Haile<sup>69</sup>, Jochen Hampe <sup>29</sup>, Heather Hampel<sup>70</sup>, Sophia Harlid<sup>71</sup>, Richard B. Hayes <sup>72</sup>, Philipp Hofer <sup>28</sup>, Michael Hoffmeister<sup>18</sup>, John L. Hopper<sup>55,73</sup>, Wan-Ling Hsu<sup>10</sup>, Wen-Yi Huang <sup>12</sup>, Thomas J. Hudson <sup>74</sup>, David J. Hunter<sup>41,75</sup>, Gemma Ibañez-Sanz<sup>13,76,77</sup>, Gregory E. Idos<sup>4</sup>, Roxann Ingersoll<sup>50</sup>, Rebecca D. Jackson<sup>78</sup>, Eric J. Jacobs <sup>35</sup>, Mark A. Jenkins <sup>55</sup>, Amit D. Joshi <sup>39,41</sup>, Corinne E. Joshi<sup>79</sup>, Temitope O. Keku<sup>80</sup>, Timothy J. Key<sup>81</sup>, Hyeong Rok Kim<sup>82</sup>, Emiko Kobayashi<sup>1</sup>, Laurence N. Kolonel<sup>83</sup>, Charles Kooperberg<sup>1</sup>, Tilman Kühn<sup>43</sup>, Sébastien Küry <sup>22</sup>, Sun-Seog Kweon<sup>84,85</sup>, Susanna C. Larsson <sup>86</sup>, Cecelia A. Laurie <sup>10</sup>, Loic Le Marchand<sup>67</sup>, Suzanne M. Leal<sup>87</sup>, Soo Chin Lee<sup>88,89</sup>, Flavio Lejbkovicz<sup>90,91,92</sup>, Mathieu Lemire<sup>74</sup>, Christopher I. Li<sup>1</sup>, Li Li<sup>93</sup>, Wolfgang Lieb<sup>94</sup>, Yi Lin<sup>1</sup>, Annika Lindblom<sup>95,96</sup>, Noralane M. Lindor<sup>97</sup>, Hua Ling<sup>50</sup>, Tin L. Louie<sup>10</sup>, Satu Männistö<sup>98</sup>, Sanford D. Markowitz<sup>99</sup>, Vicente Martín <sup>14,100</sup>, Giovanna Masala <sup>101</sup>, Caroline E. McNeil<sup>102</sup>, Marilena Melas<sup>4</sup>, Roger L. Milne <sup>55,56</sup>, Lorena Moreno<sup>17</sup>, Neil Murphy<sup>68</sup>, Robin Myte<sup>71</sup>, Alessio Naccarati <sup>103,104</sup>, Polly A. Newcomb <sup>1,36</sup>, Kenneth Offit<sup>105,106</sup>, Shuji Ogino<sup>40,41,107,108</sup>, N. Charlotte Onland-Moret<sup>54</sup>, Barbara Pardini <sup>104,109</sup>, Patrick S. Parfrey<sup>110</sup>, Rachel Pearlman<sup>70</sup>, Vittorio Perduca <sup>111,112</sup>, Paul D. P. Pharoah <sup>52</sup>, Mila Pinchev<sup>91</sup>, Elizabeth A. Platz<sup>79</sup>, Ross L. Prentice<sup>1</sup>, Elizabeth Pugh<sup>50</sup>, Leon Raskin <sup>113</sup>, Gad Rennert <sup>91,92,114</sup>, Hedy S. Rennert<sup>91,92,114</sup>, Elio Riboli<sup>115</sup>, Miguel Rodríguez-Barranco<sup>14,116</sup>, Jane Romm<sup>50</sup>, Lori C. Sakoda<sup>1,117</sup>, Clemens Schafmayer<sup>118</sup>, Robert E. Schoen<sup>119</sup>, Daniela Seminara<sup>64</sup>, Mitul Shah<sup>53</sup>, Tameka Shelford<sup>50</sup>, Min-Ho Shin <sup>84</sup>, Katerina Shulman<sup>120</sup>, Sabina Sieri<sup>121</sup>, Martha L. Slattery<sup>122</sup>, Melissa C. Southey<sup>123</sup>, Zsafia K. Stadler<sup>124</sup>, Christa Stegmaier<sup>125</sup>, Yu-Ru Su<sup>1</sup>, Catherine M. Tangen<sup>65</sup>, Stephen N. Thibodeau<sup>126</sup>, Duncan C. Thomas<sup>4</sup>, Sushma S. Thomas<sup>1</sup>, Amanda E. Toland <sup>127</sup>, Antonia Trichopoulou<sup>19,20</sup>, Cornelia M. Ulrich<sup>24</sup>, David J. Van Den Berg<sup>4</sup>, Franzel J. B. van Duijnhoven<sup>57</sup>, Bethany Van Guelpen<sup>71</sup>, Henk van Kranen<sup>128</sup>, Joseph Vijai <sup>124</sup>, Kala Visvanathan<sup>79</sup>, Pavel Vodicka<sup>103,129,130</sup>, Ludmila Vodickova<sup>103,129,130</sup>, Veronika Vymetalkova<sup>103,129,130</sup>, Korbinian Weigl <sup>18,27,131</sup>, Stephanie J. Weinstein<sup>12</sup>, Emily White<sup>1</sup>, Aung Ko Win <sup>32,55</sup>, C. Roland Wolf<sup>132</sup>, Alicja Wolk <sup>86,133</sup>, Michael O. Woods<sup>134</sup>, Anna H. Wu<sup>4</sup>, Syed H. Zaidi<sup>74</sup>, Brent W. Zanke<sup>135</sup>, Qing Zhang<sup>136</sup>, Wei Zheng <sup>137</sup>, Peter C. Scacheri<sup>138</sup>, John D. Potter<sup>1</sup>, Michael C. Bassik <sup>11</sup>, Anshul Kundaje<sup>7,11</sup>, Graham Casey<sup>139</sup>, Victor Moreno <sup>13,14,15,77</sup>, Goncalo R. Abecasis<sup>2</sup>, Deborah A. Nickerson<sup>9,141</sup>, Stephen B. Gruber<sup>4,141</sup>, Li Hsu<sup>1,10,141</sup> and Ulrike Peters <sup>1,36,141\*</sup>

<sup>1</sup>Public Health Sciences Division, Fred Hutchinson Cancer Research Center, Seattle, WA, USA. <sup>2</sup>Department of Biostatistics and Center for Statistical Genetics, University of Michigan, Ann Arbor, MI, USA. <sup>3</sup>Department of Cancer Epidemiology, H. Lee Moffitt Cancer Center and Research Institute, Tampa, FL, USA. <sup>4</sup>Department of Preventive Medicine, USC Norris Comprehensive Cancer Center, Keck School of Medicine, University of Southern California, Los Angeles, CA, USA. <sup>5</sup>Department of Epidemiology, University of Michigan, Ann Arbor, MI, USA. <sup>6</sup>Biomedical Informatics Program, Stanford University, Stanford, CA, USA. <sup>7</sup>Department of Computer Science, Stanford University, Stanford, CA, USA. <sup>8</sup>Department of Population and Quantitative Health Sciences, Case Western Reserve University, Cleveland, OH, USA. <sup>9</sup>Department of Genome Sciences, University of Washington, Seattle, WA, USA. <sup>10</sup>Department of Biostatistics, University of Washington, Seattle, WA, USA. <sup>11</sup>Department of Genetics, Stanford University, Stanford, CA, USA. <sup>12</sup>Division of Cancer Epidemiology and Genetics, National Cancer Institute, National Institutes of Health, Bethesda, MD, USA. <sup>13</sup>Cancer Prevention and Control Program, Catalan Institute of Oncology-IDIBELL, L'Hospitalet de Llobregat, Barcelona, Spain. <sup>14</sup>CIBER de Epidemiología y Salud Pública (CIBERESP), Madrid, Spain. <sup>15</sup>Department of Clinical Sciences, Faculty of Medicine, University of Barcelona, Barcelona, Spain. <sup>16</sup>Division of Epidemiology and Community Health, University of Minnesota, Minneapolis, MN, USA. <sup>17</sup>Gastroenterology Department, Hospital Clínic, Institut d'Investigacions Biomèdiques August Pi i Sunyer (IDIBAPS), Centro de Investigación Biomédica en Red de Enfermedades Hepáticas y Digestivas (CIBEREHD), University of Barcelona, Barcelona, Spain. <sup>18</sup>Division of Clinical Epidemiology and Aging Research, German Cancer Research Center (DKFZ), Heidelberg, Germany. <sup>19</sup>Hellenic Health Foundation, Athens, Greece. <sup>20</sup>WHO Collaborating Center for Nutrition and Health, Unit of Nutritional Epidemiology and Nutrition in Public Health, Department of Hygiene, Epidemiology and Medical Statistics, School of Medicine, National and Kapodistrian University of Athens, Athens, Greece. <sup>21</sup>Department of Medicine, University of North Carolina School of Medicine, Chapel Hill, NC, USA. <sup>22</sup>Service de Génétique Médicale, Centre Hospitalier Universitaire (CHU) Nantes, Nantes, France. <sup>23</sup>Leeds Institute of Medical Research at St James's, University of Leeds, Leeds, UK. <sup>24</sup>Huntsman Cancer Institute and Department of Population Health Sciences, University of Utah, Salt Lake City, UT, USA. <sup>25</sup>Department of Epidemiology, German Institute of Human Nutrition (DIfE), Potsdam-Rehbrücke, Germany. <sup>26</sup>Division of Preventive Oncology, German Cancer Research Center (DKFZ) and National Center for Tumor Diseases (NCT), Heidelberg, Germany. <sup>27</sup>German Cancer Consortium (DKTK), German Cancer Research Center (DKFZ), Heidelberg, Germany. <sup>28</sup>Institute of Cancer Research, Department of Medicine I, Medical University of Vienna, Vienna, Austria. <sup>29</sup>Department of Medicine I, University Hospital Dresden, Technische Universität Dresden (TU Dresden), Dresden, Germany. <sup>30</sup>Colorectal Oncogenomics Group, Department of Clinical Pathology, The University of Melbourne, Parkville, Victoria, Australia. <sup>31</sup>University of Melbourne Centre for Cancer Research, Victorian Comprehensive Cancer Centre, Parkville, Victoria, Australia. <sup>32</sup>Genomic Medicine and Family Cancer Clinic, Royal Melbourne Hospital, Parkville, Victoria, Australia. <sup>33</sup>Institute for Health Research, Kaiser Permanente Colorado, Denver, CO, USA. <sup>34</sup>Division of Research, Kaiser Permanente Medical Care Program, Oakland, CA, USA. <sup>35</sup>Behavioral and Epidemiology Research Group, American Cancer Society, Atlanta, GA, USA. <sup>36</sup>Department of Epidemiology, University of Washington, Seattle, WA, USA. <sup>37</sup>Division of Gastroenterology, Massachusetts General Hospital and Harvard Medical School, Boston, MA, USA. <sup>38</sup>Channing Division of Network Medicine, Brigham and Women's Hospital and Harvard Medical School, Boston, MA, USA. <sup>39</sup>Clinical and Translational Epidemiology Unit, Massachusetts General Hospital and Harvard Medical School, Boston, MA, USA. <sup>40</sup>Broad Institute of Harvard and MIT, Cambridge, MA, USA. <sup>41</sup>Department of Epidemiology, Harvard T.H. Chan School of Public Health, Harvard University, Boston, MA, USA. <sup>42</sup>Department of Immunology and Infectious Diseases, Harvard T.H. Chan School of Public Health, Harvard University, Boston, MA, USA. <sup>43</sup>Division of Cancer Epidemiology, German Cancer Research Center (DKFZ), Heidelberg, Germany. <sup>44</sup>Cancer Epidemiology Group, University Medical Centre Hamburg-Eppendorf, University Cancer Centre Hamburg (UCC), Hamburg, Germany. <sup>45</sup>Department of Epidemiology, Regional Health Council, IMIB-Arrixaca, Murcia University, Murcia, Spain. <sup>46</sup>Department of Hematology-Oncology, Chonnam National University Hospital, Hwasun, South Korea. <sup>47</sup>Department of Epidemiology and Biostatistics, Imperial College London, London, UK. <sup>48</sup>Department of Surgery and Cancer, Imperial College London, London, UK. <sup>49</sup>Department of Cancer Biology and Genetics and the Comprehensive Cancer Center, The Ohio State University, Columbus, OH, USA. <sup>50</sup>Center for Inherited Disease Research (CIDR), Institute of Genetic Medicine, Johns Hopkins University, Baltimore, MD, USA. <sup>51</sup>Translational Genomics Research Institute - An Affiliate of City of Hope, Phoenix, AZ, USA. <sup>52</sup>Department of Public Health and Primary Care, University of Cambridge, Cambridge, UK. <sup>53</sup>Centre for Cancer Genetic Epidemiology, Department of Oncology, University of Cambridge, Cambridge, UK. <sup>54</sup>Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht University, Utrecht, The Netherlands. <sup>55</sup>Centre for Epidemiology and Biostatistics, Melbourne School of Population and Global Health, The University of Melbourne, Melbourne, Victoria, Australia. <sup>56</sup>Cancer Epidemiology and Intelligence Division, Cancer Council Victoria, Melbourne, Victoria, Australia. <sup>57</sup>Division of Human Nutrition and Health, Wageningen University and Research, Wageningen, The Netherlands. <sup>58</sup>Department of Medicine, Samuel Oschin Comprehensive Cancer Institute, Cedars-Sinai Medical Center, Los Angeles, CA, USA. <sup>59</sup>Department of Preventive Medicine, Keck School of Medicine, University of Southern California, Los Angeles, CA, USA. <sup>60</sup>University of Michigan Comprehensive Cancer Center, Ann Arbor, MI, USA. <sup>61</sup>Menzies Institute for Medical Research, University of Tasmania, Hobart, Tasmania, Australia. <sup>62</sup>International Agency for Research on Cancer, World Health Organization, Lyon, France. <sup>63</sup>Lunenfeld Tanenbaum Research Institute, Mount Sinai Hospital, University of Toronto, Toronto, Ontario, Canada. <sup>64</sup>Division of Cancer Control and Population Sciences, National Cancer Institute, Bethesda, MD, USA. <sup>65</sup>SWOG Statistical Center, Fred Hutchinson Cancer Research Center, Seattle, WA, USA. <sup>66</sup>Clinical Research Division, Fred Hutchinson Cancer Research Center, Seattle, WA, USA. <sup>67</sup>University of Hawaii Cancer Research Center, Honolulu, HI, USA. <sup>68</sup>Nutrition and Metabolism Section, International Agency for Research on Cancer, World Health Organization, Lyon, France. <sup>69</sup>Division of Oncology, Department of Medicine, Stanford University, Stanford, CA, USA. <sup>70</sup>Division of Human Genetics, Department of Internal Medicine, The Ohio State University Comprehensive Cancer Center, Columbus, OH, USA. <sup>71</sup>Department of Radiation Sciences, Oncology Unit, Umeå University, Umeå, Sweden. <sup>72</sup>Division of Epidemiology, Department of Population Health, New York University School of Medicine, New York, NY, USA. <sup>73</sup>Department of Epidemiology, School of Public Health and Institute of Health and Environment, Seoul National University, Seoul, South Korea. <sup>74</sup>Ontario Institute for Cancer Research, Toronto, Ontario, Canada. <sup>75</sup>Nuffield Department of Population Health, University of Oxford, Oxford, UK. <sup>76</sup>Gastroenterology Department, Bellvitge University Hospital, L'Hospitalet de Llobregat, Barcelona, Spain. <sup>77</sup>Colorectal Cancer Group, ONCOBELL Program, Bellvitge Biomedical Research Institute-IDIBELL, Hospitalet de Llobregat, Barcelona, Spain. <sup>78</sup>Department of Medicine, Division of Endocrinology, Diabetes and Metabolism, The Ohio State University, Columbus, OH, USA. <sup>79</sup>Department of Epidemiology, Johns Hopkins Bloomberg School of Public Health, Johns Hopkins University, Baltimore, MD, USA. <sup>80</sup>Center for Gastrointestinal Biology and Disease, University of North Carolina, Chapel Hill, NC, USA. <sup>81</sup>Cancer Epidemiology Unit, Nuffield Department of Population Health, University of Oxford, Oxford, UK. <sup>82</sup>Department of Surgery, Chonnam National University Hwasun Hospital and Medical School, Hwasun, Korea. <sup>83</sup>Office of Public Health Studies, University of Hawaii Manoa, Honolulu, HI, USA. <sup>84</sup>Department of Preventive Medicine, Chonnam National University Medical School, Gwangju, Korea. <sup>85</sup>Jeonnam Regional Cancer Center, Chonnam National University Hwasun Hospital, Hwasun, Korea. <sup>86</sup>Institute of Environmental Medicine, Karolinska Institutet, Stockholm, Sweden. <sup>87</sup>Center for Statistical Genetics, Department of Molecular and Human Genetics, Baylor College of Medicine, Houston, TX, USA. <sup>88</sup>Department of Haematology-Oncology, National University Cancer Institute, Singapore, Singapore. <sup>89</sup>Cancer Science Institute of Singapore, National University of Singapore, Singapore, Singapore. <sup>90</sup>The Clalit Health Services, Personalized Genomic Service, Carmel, Haifa, Israel. <sup>91</sup>Department of Community Medicine and Epidemiology, Lady Davis Carmel Medical Center, Haifa, Israel. <sup>92</sup>Clalit National Cancer Control Center, Haifa, Israel. <sup>93</sup>Center for Community Health Integration and Case Comprehensive Cancer Center, Case Western Reserve University, Cleveland, OH, USA. <sup>94</sup>Institute of Epidemiology, PopGen Biobank, Christian-Albrechts-University Kiel, Kiel, Germany. <sup>95</sup>Department of Clinical Genetics, Karolinska University Hospital, Stockholm, Sweden. <sup>96</sup>Department of Molecular Medicine and Surgery, Karolinska Institutet, Stockholm, Sweden. <sup>97</sup>Department of Health Science Research, Mayo Clinic, Scottsdale, AZ, USA. <sup>98</sup>Department of Public Health Solutions, National Institute for Health and

Welfare, Helsinki, Finland. <sup>99</sup>Departments of Medicine and Genetics, Case Comprehensive Cancer Center, Case Western Reserve University, and University Hospitals of Cleveland, Cleveland, OH, USA. <sup>100</sup>Biomedicine Institute (BIOMED), University of León, León, Spain. <sup>101</sup>Cancer Risk Factors and Life-Style Epidemiology Unit, Institute of Cancer Research, Prevention and Clinical Network - ISPRO, Florence, Italy. <sup>102</sup>USC Norris Comprehensive Cancer Center, University of Southern California, Los Angeles, CA, USA. <sup>103</sup>Department of Molecular Biology of Cancer, Institute of Experimental Medicine of the Czech Academy of Sciences, Prague, Czech Republic. <sup>104</sup>Italian Institute for Genomic Medicine (IIGM), Turin, Italy. <sup>105</sup>Clinical Genetics Service, Department of Medicine, Memorial Sloan Kettering Cancer Center, New York, NY, USA. <sup>106</sup>Department of Medicine, Weill Cornell Medical College, New York, NY, USA. <sup>107</sup>Program in MPE Molecular Pathological Epidemiology, Department of Pathology, Brigham and Women's Hospital, Harvard Medical School, Boston, MA, USA. <sup>108</sup>Department of Oncologic Pathology, Dana-Farber Cancer Institute, Boston, MA, USA. <sup>109</sup>Department of Medical Sciences, University of Turin, Turin, Italy. <sup>110</sup>The Clinical Epidemiology Unit, Memorial University Medical School, Newfoundland, Canada. <sup>111</sup>Laboratoire de Mathématiques Appliquées MAP5 (UMR CNRS 8145), Université Paris Descartes, Paris, France. <sup>112</sup>CESP (Inserm U1018), Facultés de Médecine Université Paris-Sud, UVSQ, Université Paris-Saclay, Gustave Roussy, Villejuif, France. <sup>113</sup>Division of Epidemiology, Vanderbilt Epidemiology Center, Vanderbilt University School of Medicine, Nashville, TN, USA. <sup>114</sup>Ruth and Bruce Rappaport Faculty of Medicine, Technion-Israel Institute of Technology, Haifa, Israel. <sup>115</sup>School of Public Health, Imperial College London, London, UK. <sup>116</sup>Escuela Andaluza de Salud Pública. Instituto de Investigación Biosanitaria ibs.GRANADA, Hospitales Universitarios de Granada, Universidad de Granada, Granada, Spain. <sup>117</sup>Division of Research, Kaiser Permanente Northern California, Oakland, CA, USA. <sup>118</sup>Department of General and Thoracic Surgery, University Hospital Schleswig-Holstein, Campus Kiel, Kiel, Germany. <sup>119</sup>Department of Medicine and Epidemiology, University of Pittsburgh Medical Center, Pittsburgh, PA, USA. <sup>120</sup>Oncology Unit, Hillel Yaffe Medical Center, Hadera, Israel. <sup>121</sup>Epidemiology and Prevention Unit, Fondazione IRCCS Istituto Nazionale dei Tumori, Milan, Italy. <sup>122</sup>Department of Internal Medicine, University of Utah, Salt Lake City, UT, USA. <sup>123</sup>Genetic Epidemiology Laboratory, Department of Pathology, The University of Melbourne, Melbourne, Australia. <sup>124</sup>Department of Medicine, Memorial Sloan Kettering Cancer Center, New York, NY, USA. <sup>125</sup>Saarland Cancer Registry, Saarbrücken, Germany. <sup>126</sup>Division of Laboratory Genetics, Department of Laboratory Medicine and Pathology, Mayo Clinic, Rochester, MN, USA. <sup>127</sup>Departments of Cancer Biology and Genetics and Internal Medicine, Comprehensive Cancer Center, The Ohio State University, Columbus, OH, USA. <sup>128</sup>National Institute for Public Health and the Environment (RIVM), Bilthoven, The Netherlands. <sup>129</sup>Institute of Biology and Medical Genetics, First Faculty of Medicine, Charles University, Prague, Czech Republic. <sup>130</sup>Faculty of Medicine and Biomedical Center in Pilsen, Charles University, Pilsen, Czech Republic. <sup>131</sup>Medical Faculty, University of Heidelberg, Heidelberg, Germany. <sup>132</sup>School of Medicine, University of Dundee, Dundee, Scotland, UK. <sup>133</sup>Department of Surgical Sciences, Uppsala University, Uppsala, Sweden. <sup>134</sup>Memorial University of Newfoundland, Discipline of Genetics, St. John's, Newfoundland, Canada. <sup>135</sup>Division of Hematology, University of Toronto, Toronto, Ontario, Canada. <sup>136</sup>Genomics Shared Resource, Fred Hutchinson Cancer Research Center, Seattle, WA, USA. <sup>137</sup>Division of Epidemiology, Department of Medicine, Vanderbilt-Ingram Cancer Center, Vanderbilt Epidemiology Center, Vanderbilt University School of Medicine, Nashville, TN, USA. <sup>138</sup>Department of Genetics and Genome Sciences, Case Western Reserve University School of Medicine, Case Comprehensive Cancer Center, Cleveland, OH, USA. <sup>139</sup>Center for Public Health Genomics, University of Virginia, Charlottesville, VA, USA. <sup>140</sup>These authors contributed equally: J. R. Huyghe, S. A. Bien, T. A. Harrison. <sup>141</sup>These authors jointly supervised this work: D. A. Nickerson, S. B. Gruber, L. Hsu, U. Peters \*e-mail: [upeters@fredhutch.org](mailto:upeters@fredhutch.org)

## Methods

**Study samples.** After quality control (QC), this study included WGS data for 1,439 CRC cases and 720 controls from 5 studies, and GWAS array data for 58,131 CRC or advanced adenoma cases (3,674; 6.3% of cases), and 67,347 controls from 45 studies from GECCO, CORECT and CCFR. The stage 1 meta-analysis comprised existing genotyping data from 30 studies that were included in previously published CRC GWAS<sup>13,18,22</sup>. After QC, the stage 1 meta-analysis included 34,869 cases and 29,051 controls. Study participants were predominantly of European ancestry (31,843 cases and 26,783 controls; 91.7% of participants). Because the vast majority of known CRC risk variants are shared between Europeans and East Asians<sup>57</sup>, we included 3,026 cases and 2,268 controls of East Asian ancestry to increase power for discovery. The stage 2 meta-analysis comprised newly generated genotype data involving four genotyping projects and 22 studies. After QC, the stage 2 meta-analysis included 23,262 cases and 38,296 controls, which were all of European ancestry. Studies, sample selection and matching are described in Supplementary Note. Supplementary Table 1 provides details on sample numbers and demographic characteristics of study participants. All participants provided written informed consent, and each study was approved by the relevant research ethics committee or institutional review board. Four normal colon mucosa biopsies for ATAC-seq were obtained from patients with a normal colon at colonoscopy at the Institut d'Investigació Biomèdica de Bellvitge (IDIBELL), Spain. Patients signed informed consent, and the protocol was approved by the Bellvitge Hospital Ethics Committee (Colscreen protocol PR084/16).

**Whole-genome sequencing.** We performed low-pass WGS of 2,192 samples from five studies at the University of Washington Northwest Genomics Center (Seattle, WA, USA). Cases and controls were processed and sequenced together. Libraries were prepared with ThruPLEX DNA-seq kits (Rubicon Genomics) and paired-end sequencing was performed using Illumina HiSeq 2500 sequencers. Reads were mapped to human reference genome (GRCh37 assembly) using Burrows-Wheeler aligner BWA v0.6.2 (ref. <sup>83</sup>). Fold genomic coverage averaged  $5.3\times$  (range:  $3.8\text{--}8.6\times$ ). We used the GotCloud population-based multisample variant calling pipeline<sup>83</sup> for post-processing of BAM files with initial alignments, and to detect and call single-nucleotide variants and short indels. After removing duplicated reads and recalibrating base quality scores, QC checks included sample contamination detection. Variants were jointly called across all samples. To identify high-quality sites, the GotCloud pipeline performs a two-step filtering process. First, lower quality variants are identified by applying individual variant quality statistic filters. Next, variants failing multiple filters are used as negative examples to train a support vector machine classifier. Finally, we performed a haplotype-aware genotype refinement step via Beagle<sup>84</sup> and ThunderVCF<sup>85</sup> on the support vector machine-filtered VCF files. After further sample QC, we excluded samples with estimated DNA contamination  $>3\%$  ( $n=16$ ), duplicated samples ( $n=5$ ) or related individuals ( $n=1$ ), sex discrepancies ( $n=0$ ), and samples with low concordance with GWAS array data ( $n=11$ ). We checked for ancestry outliers by performing principal component analysis (PCA) after merging in data for shared, LD-pruned single-nucleotide variants for 1,092 individuals from the 1000 Genomes Project<sup>86</sup>. After QC, sequences were available for 1,439 CRC cases and 720 controls of European ancestry.

**GWAS genotype data and quality control.** Details of genotyping and QC for studies included in the stage 1 meta-analysis are described elsewhere<sup>13,18,22</sup>. Supplementary Table 1 provides details of genotyping platforms used. Before association analysis, we pooled individual-level genotype data of all stage 1 studies for a subset of SNPs to enable identification of unexpected duplicates and close relatives. We calculated identity by descent for each pair of samples using KING-robust<sup>87</sup> and excluded duplicates and individuals that were second-degree or more closely related. As part of stage 2, 28,805 individuals from 19 studies were newly genotyped on a custom Illumina array based on the Infinium OncoArray-500K<sup>26</sup> and a panel of 15,802 successfully manufactured custom variants (described in Supplementary Note). An additional 8,725 individuals from five studies were genotyped on the Illumina HumanOmniExpressExome-8v1-2 array. Genotyping and calling for both projects were performed at the Center for Inherited Disease Research (CIDR) at Johns Hopkins University. Genotypic data that passed initial QC at CIDR subsequently underwent QC at the University of Washington Genetic Analysis Center (UW GAC) using standardized methods detailed in Laurie et al.<sup>88</sup>. The median call rate for the custom Infinium OncoArray-500K data was 99.97%, and error rate estimated from 301 sample duplicate pairs was  $9.99\times 10^{-7}$ . A relatively low number of samples ( $n=246$ ) had a missing call rate  $>2\%$ , with the highest being 3.48%, and were included in analysis. For the HumanOmniExpressExome-8v1-2 data, median call rate was 99.96%, and the error rate estimated from 179 sample duplicate pairs was  $2.65\times 10^{-6}$ . Some 30 samples had a missing call rate  $>2\%$ , with the highest being 3.79%, and were included in analysis. We excluded samples with discrepancies between reported and genotypic sex based on X chromosome heterozygosity and the mean values of sex chromosome probe intensities, unintentional duplicates, and close relatives, defined as individuals that are second-degree or more closely related. After further excluding individuals of non-European ancestry as determined by PCA (see below), the custom OncoArray data included in analysis comprised 11,852 CRC

cases and 11,895 controls, and the HumanOmniExpressExome-8v1-2 array data included in analysis comprised 4,439 CRC cases and 4,115 controls. Only variants passing QC were used for imputation. We excluded variants failing CIDR technical filters or UW GAC quality filters, which included missing call rate  $>2\%$ , discordant calls in sample duplicates, and departures from Hardy-Weinberg equilibrium (HWE) ( $P < 1\times 10^{-4}$ ) based on European-ancestry controls. The stage 2 analysis also included genotype data from the CORSA study (Supplementary Note). In total, 2,354 individuals were genotyped using the Affymetrix Axiom Genome-Wide Human CEU 1 Array. We called genotypes using the AxiomGT1 algorithm. All samples had missing call rate  $<3\%$ . We excluded samples with discrepancies between reported and genotypic sex ( $n=20$ ), close relatives, defined as individuals that are second-degree or more closely related ( $n=94$ ), as inferred using KING-robust<sup>87</sup>, and individuals of non-European ancestry ( $n=6$ ) as inferred from PCA. After QC, data included in analysis comprised 1,460 cases and 774 controls. Before phasing and imputation, we filtered out SNPs with missing call rate  $>2\%$ , or HWE  $P < 1\times 10^{-4}$ . Imputed genotype data were obtained from UK Biobank and QC and imputation are described elsewhere<sup>89</sup>. A nested case-control data set was constructed as described in Supplementary Note. We excluded individuals of non-European ancestry as inferred from PCA, and randomly dropped one individual from each pair that were more closely related than third-degree relatives as inferred using KING-robust. This led to excluding 137 samples. In total, 5,356 CRC ( $n=5,004$ ) or advanced adenoma ( $n=352$ ) cases and 21,407 matched controls were included in the replication analysis.

**Principal component analysis.** After excluding close relatives, we performed PCA using PLINK1.9 (ref. <sup>90</sup>) on LD-pruned sets of autosomal SNPs obtained by removing regions with extensive long-range LD<sup>91,92</sup>, SNPs with  $MAF < 5\%$ , HWE  $P < 1\times 10^{-4}$ , or any genotype missingness, and carrying out LD pruning using the PLINK option “-indep-pairwise 50 5 0.2.” To identify population outliers we merged in 1,092 individuals from 1000 Genomes Project Phase III and performed PCA using the intersection of variants<sup>93</sup>.

**Genotype imputation.** The 2,159 whole-genome sequences described above were used to create a phased imputation reference panel. After estimating haplotypes for all GWAS array data sets using SHAPEIT2 (ref. <sup>94</sup>), we used minimac3 (ref. <sup>95</sup>) to impute from this reference panel (19.6 million variants with  $MAC > 1$ ) into the GWAS data sets described above. We also imputed to the Haplotype Reference Consortium (HRC) panel<sup>25</sup> (39.2 million variants) using the University of Michigan Imputation Server<sup>95</sup>. To improve imputation accuracy for stage 1 data sets, phasing and imputation were performed after pooling studies or genotyping projects that used the same, or very similar, genotyping platforms (Supplementary Table 1). For stage 2, we performed phasing and imputation separately for each genotyping project data set and imputed to the HCR panel.

**Statistical analyses. Association testing of sequence data.** We tested variants with  $MAC \geq 5$  for CRC association using Firth's bias-reduced logistic regression as implemented in EPACTS (see URLs) and adjusted for sex, age, study and three principal components calculated from an LD-pruned set of genotypes. We performed rare variant aggregate tests at the gene and enhancer level using the Mixed effects Score Test (MiST)<sup>96</sup>. This unified test is a linear combination between unidirectional burden and bidirectional variance component tests that performs best in terms of statistical power across a range of architectures<sup>97</sup>.

**Association and meta-analysis.** Stage 1 comprised two large mega-analyses of pooled individual-level genotype data sets (Supplementary Table 12). The four stage-2 genotyping project data sets were analyzed separately. Within each data set, variants with an imputation accuracy  $r^2 \geq 0.3$  and  $MAC \geq 50$  were tested for CRC association using the imputed genotype dosage in a logistic regression model adjusted for age, sex and study or genotyping project-specific covariates, including principal components to adjust for population structure (Supplementary Table 12). To account for residual confounding within CORSA, we tested association with each variant using a linear mixed model and kinship matrix calculated from the data, as implemented in EMMAX<sup>98</sup>. To enable meta-analysis, we then calculated approximate allelic log[OR] estimates and corresponding standard errors as described in Cook et al.<sup>99</sup>

Next, we combined association summary statistics across analyses via fixed-effects inverse variance-weighted meta-analysis. Because Wald tests can be notably anti-conservative for rare variant associations, we also performed likelihood ratio-based tests, followed by sample-size weighted meta-analysis, as implemented in METAL<sup>100</sup>. In total, 16,900,397 variants were analyzed. To examine residual population stratification, we inspected quantile-quantile plots of test statistics (Supplementary Fig. 8) and calculated genomic control inflation statistics ( $\lambda_{GC}$ ).  $\lambda_{GC}$  for the combined meta-analysis was 1.105, and for stage 1 and stage 2 meta-analyses was 1.071 and 1.075, respectively. Because  $\lambda_{GC}$  increases with sample size for polygenic phenotypes, even in the absence of confounding biases<sup>101</sup>, we investigated the effect of confounding due to residual population stratification using LD score regression<sup>102</sup>. Because of limitations of LD score regression, this analysis is restricted to common variants ( $MAF \geq 1\%$ ) for which  $\lambda_{GC}$  was 1.188 in the combined meta-analysis. The LD score regression intercept was 1.067, which is

substantially less than  $\lambda_{GC}$ , indicating at most a small contribution of bias and that inflation in  $\chi^2$  statistics results mostly from polygenicity. We also calculated  $\lambda_{1,000}$ , which is the equivalent inflation statistic for a study with 1,000 cases and 1,000 controls<sup>103</sup>. For the combined meta-analysis,  $\lambda_{1,000}$  was 1.004 and for both stage 1 and stage 2 meta-analyses it was 1.003.

**Significance threshold for the replication genotyping experiment.** To protect against probe design failure, we built redundancy into the custom genotyping panel by including LD proxies of independently associated variants selected for follow-up. To determine the number of independent tests, we performed LD clumping of the 9,198 analyzed variants that were selected for replication genotyping based on the stage 1 meta-analysis, and that survived filters described above. Using an  $r^2$  threshold of 0.1 this translated to 6,438 independent tests and a Bonferroni significance threshold of  $0.05/6,438 = 7.8 \times 10^{-6}$ .

**Conditional and joint multiple-variant analysis.** To identify additional distinct association signals at CRC loci, we performed a series of conditional meta-analyses. At each locus attaining  $P < 5 \times 10^{-8}$ , we included the genotype dosage for the variant showing the strongest statistical evidence for association in the region in the combined meta-analysis, as an additional covariate in the respective logistic regression models. Association summary statistics for each variant in the region were then combined across studies by a fixed-effects meta-analysis. If at least one association signal attained a significance level of  $P < 1 \times 10^{-5}$  in this meta-analysis, we performed a second round of conditional meta-analysis, adding the variant showing the strongest statistical evidence for association in the region in the first round of conditional meta-analysis as a covariate to the logistic regression models used in the first round. We repeated this procedure and kept adding variants to the model until no additional variants at the locus attained  $P < 1 \times 10^{-5}$ . Finally, we performed a joint multiple-variant analysis in which we jointly estimated the effects of variants selected in each step and tested for each variant whether the  $P$  value from the joint multiple-variant analysis ( $P_j$ ) was  $< 1 \times 10^{-5}$ . Analyses were performed on 2-Mb windows centered on the most associated variant in the unconditional analysis. If windows overlapped, we performed the analysis on the collapsed genomic region. Because of extensive LD, we used a 4-Mb window for the MHC region.

**Definition of known loci.** We compiled a list of 62 previously reported genome-wide significant CRC association signals from the literature (Supplementary Table 3). Because of improved power and coverage of our study, we identified the most associated variant at each signal, and used these lead variants for further analyses, rather than the previously reported index variant.

**Refinement of association signals.** To refine new association signals, we constructed credible sets that were 99% probable, based on posterior probability, to contain the causal disease-associated SNP<sup>104</sup>. In brief, for each distinct signal, we retained a candidate set of variants by identifying all analyzed variants with  $r^2 \geq 0.1$  with the most associated variant within a 2-Mb window centered on the most associated variant. We calculated approximate Bayes' factors (ABFs)<sup>105</sup> for each variant as:

$$ABF = \sqrt{1-r} e^{z^2/2}$$

where  $r = 0.04/(s.e.^2 + 0.04)$ ,  $z = \beta/s.e.$ , and  $\beta$  and s.e. are the log[OR] estimate and its standard error from the combined meta-analysis, respectively. For loci with multiple distinct signals, results are based on conditional meta-analysis, adjusting for all other index variants in the region. We then calculated the posterior probability of being causal as  $ABF/T$  where  $T$  is the sum of  $ABF$  values over all candidate variants. Next, variants were ranked in decreasing order by posterior probabilities and the 99% credible set was obtained by including variants with the highest posterior probabilities until the cumulative posterior probability was  $\geq 99\%$ .

**Functional genomic annotation.** To nominate variants for future laboratory follow-up, we performed bioinformatic analysis at each new signal using our functional annotation database, and a custom University of California Santa Cruz (UCSC) analysis data hub. Using ANNOVAR<sup>106</sup>, we annotated lead variants and variants in LD ( $r^2 \geq 0.4$ ) with the lead variant, relative to features pertaining to: (1) gene-centric function (PolyPhen2<sup>107</sup>); (2) genome-wide functional prediction scores (CADD<sup>108</sup>, DANN<sup>109</sup> and EigenPC<sup>110</sup>); (3) disease relatedness (GWAS catalog<sup>66</sup>) and (4) CRC-relevant regulatory functions (enhancer, repressor, DNA accessible and transcription factor-binding site<sup>111,112</sup>; Supplementary Table 13). Supplementary Table 8 summarizes variant annotations relative to the Consensus Coding Sequence (CCDS) Project<sup>113</sup>, and reference genome GRCh37. Variants were maintained in Supplementary Table 8 if they met any of the following conditions: DANN score  $\geq 0.9$ , CADD phred score  $\geq 20$ , EigenPC phred score  $\geq 17$ , PolyPhen2 prediction of "probably damaging," stop-loss, stop-gain, or splicing variant, or positioned in a predicted regulatory element. We visually inspected loci overlapping with CRC-relevant functional genomic annotations. Variants positioned in enhancers with aberrant CRC activity were identified by comparing epigenomes of nondiseased colorectal tissues or colonic

crypt cells to epigenomes of primary CRC cell lines. We prioritized target genes for loci with predicted regulatory function. Evidence suggests that topological association domains can be used to map physical boundaries on gene promoter interactions with distal regulatory elements<sup>114–116</sup>. As such, we used GM12878 Hi-C Chromosome Conformation Capture data to identify gene promoters that were in the same topological association domains as risk loci using the WashU Epigenome Browser (see URLs). Genes in this list were further prioritized based on biological relevance and expression quantitative trait loci data from GTEx<sup>117</sup> using HaploReg v4.1 (ref. <sup>118</sup>).

**Assay for transposable-accessible chromatin using sequencing.** We generated high-resolution maps of DNA-accessible regions in normal colon mucosa samples using ATAC-seq. Using the updated omni-ATAC protocol for archival samples, we performed ATAC-seq in four colon mucosa biopsies from the Catalan Institute of Oncology (ICO)-biobank taken from participants undergoing screening at IDIBELL, Spain. Biopsies were cryopreserved by slow freezing using a solution of 10% dimethylsulfoxide, 90% medium and Mr. Frosty Cryo 1°C Freezing Containers (Thermo Scientific). ATAC-seq was implemented as prescribed with two exceptions. Instead of Dounce homogenizer, we used a tissue lyser and stainless bead system, pulverizing at 40 Hz for 2 min and pulsing at 50 Hz for 10–20 s. Secondly, Illumina library quantification was performed using picogreen quantitation and TapeStation instead of KAPA quantitative PCR. Libraries were sequenced to an average of 25 million paired end reads using Illumina HiSeq 2500. The ENCODE data processing pipeline was implemented (see URLs) aligning to hg19 (ref. <sup>119</sup>). QC results are summarized in Supplementary Table 14.

**Regulatory and functional information enrichment analysis.** We used GARFIELD<sup>74</sup> to identify cell types, tissues and functional genomic features relevant to CRC risk. This method tests for enrichment of association in features primarily extracted from ENCODE and Roadmap Epigenomics Project data, while accounting for confounding variables, including LD. We applied default settings and used the author-supplied data, which are suitable for analysis of GWAS results based on individuals of European ancestry.

**Pathway and gene set enrichment analysis.** We used MAGENTA to test predefined gene sets (for example, KEGG pathways) for enrichment for CRC risk associations<sup>75</sup>. We used combined meta-analysis results as input and applied default settings, which included removing genes that fall in the MHC region from analysis. Enrichment was tested at two gene  $P$ -value cutoffs: the 95th and 75th percentiles of all gene  $P$  values in the genome.

**Estimation of contribution of rare variants to heritability.** We used the LD- and MAF-stratified component GREML (GREML-LDMS) method as implemented in GCTA<sup>76</sup> to estimate the proportion of variation in liability for CRC explained by all imputed autosomal variants (that is, estimate of narrow-sense heritability  $h_g^2$ ), and the proportion contributed by rare variants (MAF  $\leq 1\%$ ). Because of computational limitations we analyzed a subset of 11,895 cases and 14,659 controls imputed to our WGS panel. We analyzed individual-level data for 17,649,167 imputed variants with MAC  $> 3$  and HWE test  $P \geq 10^{-6}$ . Following Yang et al.<sup>76</sup>, we did not filter based on imputation quality. In brief, we stratified variants into groups based on MAF (boundaries at 0.001, 0.01, 0.1, 0.2, 0.3 and 0.4) and mean LD score (boundaries at quartiles) calculated as described in Yang et al.<sup>76</sup>. We then calculated genetic relationship matrices for each of these 28 variant partitions and jointly estimated variance components for these partitions, adjusting for age, sex, study, genotyping batch and three genotype principal components. From the variance component estimates and their variance-covariance matrix we estimated the contribution of rare variants (MAF  $\leq 1\%$ ) and common variants (MAF  $> 1\%$ ), and calculated standard errors using the delta method. We tested significance of the contribution of rare variants using a likelihood ratio test. To calculate heritability on the underlying liability scale we interpreted  $K$  as lifetime risk<sup>120</sup> and used an estimate of 4.3% (Surveillance, Epidemiology, and End Results Program (SEER) Cancer Statistics, 2011–2013).

**Familial relative risk explained by genetic variants.** We assumed a multiplicative model within and between variants and calculated the proportion of familial RR explained by a given set of genetic variants as  $\frac{\sum_i \log \lambda_i}{\log \lambda_0}$ , where  $\lambda_0$  is the overall familial RR to first-degree relatives of cases,  $\lambda_i$  is the familial RR due to variant  $i$  calculated as  $\lambda_i = \frac{p_i^2 + q_i}{q_i(p_i + q_i)}$ , where  $p_i$  is the risk allele frequency for variant  $i$ ,  $q_i = 1 - p_i$ , and  $r_i$  is the estimated per-allele OR<sup>3,121</sup>. We adjusted the OR estimates of new association signals for the winner's curse following Zhong and Prentice<sup>38</sup>. We represented previously identified association signals by the variant showing the strongest statistical evidence of association in the combined meta-analysis, and assumed that the winner's curse was negligible. We assumed  $\lambda_0$  to be 2.2 (ref. <sup>122</sup>). Using the delta method, we computed the variance for the proportion of familial RR as follows:

$$\sum_i \text{Var}(r_i) \left[ \frac{1}{\log \lambda_0} \frac{1}{\lambda_i} \frac{2p_i q_i (r_i - 1)}{(p_i r_i + q_i)^3} \right]^2$$

**Absolute risk of CRC incidence and starting age of first screening.** We constructed a PRS as a weighted sum of the number of risk alleles carried by an individual, using the per-allele OR for each variant as weights. OR estimates for newly discovered variants were adjusted for the winner's curse to avoid potential inflation<sup>78</sup>. Assuming all genetic variants are independent, let  $X$  denote a PRS constructed based on  $K$  variants:  $X = \sum_{i=1}^K \hat{\beta}_i Z_i$ , where  $\hat{\beta}_i$  and  $Z_i$  are the estimated OR and the number of risk alleles for variant  $i$ . We assumed  $X$  follows a normal distribution  $N(\mu, \sigma^2)$ , where the estimates of mean and variance are computed as follows:

$$\hat{\mu} = \sum_{i=1}^K \hat{\beta}_i \times 2 \times \hat{p}_i \quad \text{and} \quad \hat{\sigma}^2 = \sum_{i=1}^K \hat{\beta}_i^2 \times 2 \times \hat{p}_i \times (1 - \hat{p}_i),$$

where  $\hat{p}_i$  is the risk allele frequency for variant  $i = 1, \dots, K$ . Then the baseline hazard at each age  $t$ ,  $\hat{\lambda}_0(t)$ , is computed as follows:

$$\hat{\lambda}_0(1) = \lambda^*(1) \frac{\int f(x) dx}{\int e^{x^2} f(x) dx}$$

$$\hat{\lambda}_0(t) = \lambda^*(t) \frac{\int \exp\left(-\sum_{i=1}^{t-1} \hat{\lambda}_0(i) e^{x^2}\right) f(x) dx}{\int \exp\left(-\sum_{i=1}^{t-1} \hat{\lambda}_0(i) e^{x^2}\right) e^{x^2} f(x) dx}$$

for  $t = 2, \dots, 100$ ,

and  $\lambda^*(t)$  are the incidence rates for non-Hispanic whites who have not taken an endoscopy before, derived from population incidence rates during 1992–2005 from the SEER Registry. Using these baseline hazard rates, we estimated the 10-year absolute risk of developing CRC given age and a PRS as described<sup>123</sup>. By setting a risk threshold as the average of the 10-year CRC risk for a 50-year old man (1.25%) and woman (0.68%) who have not previously received an endoscopy<sup>124</sup>, that is,  $(1.25\% + 0.68\%)/2 = 0.97\%$ , we estimated the recommended starting age of first screening given the PRS. Variants and OR estimates used in these analyses are in Supplementary Table 15.

**Reporting Summary.** Further information on research design is available in the Nature Research Reporting Summary linked to this article.

## Data availability

All whole-genome sequence data have been deposited in the database of Genotypes and Phenotypes (dbGaP), which is hosted by NCBI, under accession number [phs001554.v1.p1](#). All custom Infinium OncoArray-500K array data for the studies in the stage 2 meta-analysis have been deposited at dbGaP under accession number [phs001415.v1.p1](#). All Illumina HumanOmniExpressExome-8v1-2 array data for the studies in the stage 2 meta-analysis have been deposited at dbGaP under accession number [phs001315.v1.p1](#). Genotype data for the studies included in the stage 1 meta-analysis have been deposited at dbGaP under accession number [phs001078.v1.p1](#). The UK Biobank resource was accessed through application number 8614. CRC-relevant epigenome data were obtained from the NCBI Gene Expression Omnibus (GEO) database under accession number [GSE77737](#).

## References

82. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**, 1754–1760 (2009).
83. Jun, G., Wing, M. K., Abecasis, G. R. & Kang, H. M. An efficient and scalable analysis framework for variant extraction and refinement from population-scale DNA sequence data. *Genome Res.* **25**, 918–925 (2015).
84. Browning, B. L. & Yu, Z. Simultaneous genotype calling and haplotype phasing improves genotype accuracy and reduces false-positive associations for genome-wide association studies. *Am. J. Hum. Genet.* **85**, 847–861 (2009).
85. Li, H. A statistical framework for SNP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. *Bioinformatics* **27**, 2987–2993 (2011).
86. 1000 Genomes Project Consortium et al. A map of human genome variation from population-scale sequencing. *Nature* **467**, 1061–1073 (2010).
87. Manichaikul, A. et al. Robust relationship inference in genome-wide association studies. *Bioinformatics* **26**, 2867–2873 (2010).
88. Laurie, C. C. et al. Quality control and quality assurance in genotypic data for genome-wide association studies. *Genet. Epidemiol.* **34**, 591–602 (2010).
89. Bycroft, C. et al. Genome-wide genetic data on ~500,000 UK Biobank participants. Preprint at <https://www.biorxiv.org/content/early/2017/07/20/166298> (2017).
90. Chang, C. C. et al. Second-generation PLINK: rising to the challenge of larger and richer datasets. *Gigascience* **4**, 7 (2015).
91. Price, A. L. et al. Long-range LD can confound genome scans in admixed populations. *Am. J. Hum. Genet.* **83**, 132–135 (2008).
92. Weale, M. E. Quality control for genome-wide association studies. *Methods Mol. Biol.* **628**, 341–372 (2010).
93. 1000 Genomes Project Consortium. et al. An integrated map of genetic variation from 1,092 human genomes. *Nature* **491**, 56–65 (2012).
94. Delaneau, O., Howie, B., Cox, A. J., Zagury, J.-F. & Marchini, J. Haplotype estimation using sequencing reads. *Am. J. Hum. Genet.* **93**, 687–696 (2013).
95. Das, S. et al. Next-generation genotype imputation service and methods. *Nat. Genet.* **48**, 1284–1287 (2016).
96. Sun, J., Zheng, Y. & Hsu, L. A unified mixed-effects model for rare-variant association in sequencing studies. *Genet. Epidemiol.* **37**, 334–344 (2013).
97. Moutsianas, L. et al. The power of gene-based rare variant methods to detect disease-associated variation and test hypotheses about complex disease. *PLoS Genet.* **11**, e1005165 (2015).
98. Kang, H. M. et al. Variance component model to account for sample structure in genome-wide association studies. *Nat. Genet.* **42**, 348–354 (2010).
99. Cook, J. P., Mahajan, A. & Morris, A. P. Guidance for the utility of linear models in meta-analysis of genetic association studies of binary phenotypes. *Eur. J. Hum. Genet.* **25**, 240–245 (2017).
100. Willer, C. J., Li, Y. & Abecasis, G. R. METAL: fast and efficient meta-analysis of genomewide association scans. *Bioinformatics* **26**, 2190–2191 (2010).
101. Yang, J. et al. Genomic inflation factors under polygenic inheritance. *Eur. J. Hum. Genet.* **19**, 807–812 (2011).
102. Bulik-Sullivan, B. K. et al. LD Score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nat. Genet.* **47**, 291–295 (2015).
103. Michailidou, K. et al. Large-scale genotyping identifies 41 new loci associated with breast cancer risk. *Nat. Genet.* **45**, 353–361 (2013).
104. Wellcome Trust Case Control Consortium. et al. Bayesian refinement of association signals for 14 loci in 3 common diseases. *Nat. Genet.* **44**, 1294–1301 (2012).
105. Wakefield, J. A Bayesian measure of the probability of false discovery in genetic epidemiology studies. *Am. J. Hum. Genet.* **81**, 208–227 (2007).
106. Wang, K., Li, M. & Hakonarson, H. ANNOVAR: functional annotation of genetic variants from high-throughput sequencing data. *Nucleic Acids Res.* **38**, e164 (2010).
107. Adzhubei, I., Jordan, D. M. & Sunyaev, S. R. Predicting functional effect of human missense mutations using PolyPhen-2. *Curr. Protoc. Hum. Genet.* **Chapter 7**, Unit7.20 (2013).
108. Kircher, M. et al. A general framework for estimating the relative pathogenicity of human genetic variants. *Nat. Genet.* **46**, 310–315 (2014).
109. Quang, D., Chen, Y. & Xie, X. DANN: a deep learning approach for annotating the pathogenicity of genetic variants. *Bioinformatics* **31**, 761–763 (2015).
110. Ionita-Laza, I., McCallum, K., Xu, B. & Buxbaum, J. D. A spectral approach integrating functional genomic annotations for coding and noncoding variants. *Nat. Genet.* **48**, 214–220 (2016).
111. Roadmap Epigenomics Consortium et al. Integrative analysis of 111 reference human epigenomes. *Nature* **518**, 317–330 (2015).
112. Corradin, O. et al. Combinatorial effects of multiple enhancer variants in linkage disequilibrium dictate levels of gene expression to confer susceptibility to common traits. *Genome Res.* **24**, 1–13 (2014).
113. Pruitt, K. D. et al. The consensus coding sequence (CCDS) project: Identifying a common protein-coding gene set for the human and mouse genomes. *Genome Res.* **19**, 1316–1323 (2009).
114. Harmston, N. et al. Topologically associating domains are ancient features that coincide with Metazoan clusters of extreme noncoding conservation. *Nat. Commun.* **8**, 441 (2017).
115. Berlivet, S. et al. Clustering of tissue-specific sub-TADs accompanies the regulation of HoxA genes in developing limbs. *PLoS Genet.* **9**, e1004018 (2013).
116. Hu, Z. & Tee, W.-W. Enhancers and chromatin structures: regulatory hubs in gene expression and diseases. *Biosci. Rep.* **37**, BSR20160183 (2017).
117. GTEx Consortium. The Genotype-Tissue Expression (GTEx) project. *Nat. Genet.* **45**, 580–585 (2013).
118. Ward, L. D. & Kellis, M. HaploReg: a resource for exploring chromatin states, conservation, and regulatory motif alterations within sets of genetically linked variants. *Nucleic Acids Res.* **40**, D930–D934 (2012).
119. Landt, S. G. et al. ChIP-seq guidelines and practices of the ENCODE and modENCODE consortia. *Genome Res.* **22**, 1813–1831 (2012).
120. Witte, J. S., Visscher, P. M. & Wray, N. R. The contribution of genetic variants to disease depends on the ruler. *Nat. Rev. Genet.* **15**, 765–776 (2014).
121. Cox, A. et al. A common coding variant in *CASP8* is associated with breast cancer risk. *Nat. Genet.* **39**, 352–358 (2007).
122. Johns, L. E. & Houlston, R. S. A systematic review and meta-analysis of familial colorectal cancer risk. *Am. J. Gastroenterol.* **96**, 2992–3003 (2001).
123. Hsu, L. et al. A model to determine colorectal cancer risk using common genetic susceptibility loci. *Gastroenterology* **148**, 1330–1339.e14 (2015).
124. Jeon, J. et al. Determining risk of colorectal cancer and starting age of screening based on lifestyle, environmental, and genetic factors. *Gastroenterology* **154**, 2152–2164.e19 (2018).

## Reporting Summary

Nature Research wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Research policies, see [Authors & Referees](#) and the [Editorial Policy Checklist](#).

### Statistical parameters

When statistical analyses are reported, confirm that the following items are present in the relevant location (e.g. figure legend, table legend, main text, or Methods section).

n/a Confirmed

- The exact sample size ( $n$ ) for each experimental group/condition, given as a discrete number and unit of measurement
- An indication of whether measurements were taken from distinct samples or whether the same sample was measured repeatedly
- The statistical test(s) used AND whether they are one- or two-sided  
*Only common tests should be described solely by name; describe more complex techniques in the Methods section.*
- A description of all covariates tested
- A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons
- A full description of the statistics including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals)
- For null hypothesis testing, the test statistic (e.g.  $F$ ,  $t$ ,  $r$ ) with confidence intervals, effect sizes, degrees of freedom and  $P$  value noted  
*Give  $P$  values as exact values whenever suitable.*
- For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings
- For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes
- Estimates of effect sizes (e.g. Cohen's  $d$ , Pearson's  $r$ ), indicating how they were calculated
- Clearly defined error bars  
*State explicitly what error bars represent (e.g. SD, SE, CI)*

Our web collection on [statistics for biologists](#) may be useful.

### Software and code

Policy information about [availability of computer code](#)

Data collection

No software was used to collect data.

Data analysis

Genotype calling for new genotyping projects was performed using Illumina's GenomeStudio version 2011.1, Genotyping Module v1.9.4, GenTrain v 1.0 algorithm, and Affymetrix Power Tools version 1.18.2 and the AxiomGT1 algorithm. Alignment of whole-genome sequence reads was performed using BWA v0.6.2 and variant calling using the GotCloud pipeline, which uses Beagle and ThunderVCF. Software used for genotyping quality control included PLINK1.9, KING 2.0.9, and the R package GWASTools. Phasing and imputation were performed using SHAPEIT2 and minimac3, respectively. Imputation to the Haplotype Reference Consortium panel was performed using the University of Michigan Imputation Server. Association analyses were performed using R, EPACTS v3.3.0, and METAL. Variant annotation was performed using ANNOVAR. Gene prioritization was performed using HaploReg v4.1. Functional and pathway enrichment analyses were conducted using GARFIELD v2 and MAGENTA v2. ATAC-Seq peak calling was performed using the ENCODE pipeline. Heritability analyses were performed using the GCTA package version 1.25.2. LD score regression was performed using ldsc version 1.0.0. The familial relative risk explained and risk prediction calculations were performed using R.

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors/reviewers upon request. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Research [guidelines for submitting code & software](#) for further information.



## Data

Policy information about [availability of data](#)

All manuscripts must include a [data availability statement](#). This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A list of figures that have associated raw data
- A description of any restrictions on data availability

All whole-genome sequence data have been deposited at the database of Genotypes and Phenotypes (dbGaP), which is hosted by the U.S. National Center for Biotechnology Information (NCBI), under accession number phs001554.v1.p1. All custom Infinium OncoArray-500K array data for the studies in the Stage 2 meta-analysis have been deposited at dbGaP under accession number phs001415.v1.p1. All Illumina HumanOmniExpressExome-8v1-2 array data for the studies in the Stage 2 meta-analysis have been deposited at dbGaP under accession number phs001315.v1.p1. Genotype data for the studies included in the Stage 1 meta-analysis have been deposited at dbGaP under accession number phs001078.v1.p1. The UK Biobank resource was accessed through application number 8614.

## Field-specific reporting

Please select the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

Life sciences       Behavioural & social sciences       Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see [nature.com/authors/policies/ReportingSummary-flat.pdf](https://www.nature.com/authors/policies/ReportingSummary-flat.pdf)

## Life sciences study design

All studies must disclose on these points even when the disclosure is negative.

Sample size	This work used newly generated whole-genome sequencing data (WGS; n=2,159), and existing and newly generated GWAS data on colorectal cancer (CRC; total n=125,478). Both datasets comprise the largest and most comprehensive to date. Rationale for the WGS study design which balanced sample size and sequencing coverage, is detailed under Study Overview in the Results section of the article. No sample size calculations were performed for the GWAS meta-analysis. We attempted to include as many studies as possible in the meta-analysis to maximize statistical power for discovery. Given that we almost doubled the sample size compared to the previous largest CRC GWAS meta-analysis, and the improved and more comprehensive genotype imputation, we expected to have substantially higher statistical power to detect variants with smaller effect sizes or low minor allele frequencies.
Data exclusions	We excluded samples and variants, using standardized methods and predefined filtering criteria to remove unreliable genotype calls or problematic samples, as described in the Online Methods. Sample exclusions included samples with evidence of DNA contamination, samples with high missing genotype rates, unintentional duplicate pairs, sex discrepancies, and closely related individuals. To avoid confounding due to population stratification, we excluded ancestry outliers identified through principal components analysis. Variants were excluded according to predefined quality controls filtering criteria as detailed in the Online Methods. For the WGS dataset this involved quality filtering based on a support vector machine (SVM) classifier; for the genotype array data, this included missing call rate filters, Hardy-Weinberg equilibrium, and various technical filters. Imputed genotype data were filtered on imputation quality and minor allele count.
Replication	We present two sets of results: 1) new association signals for CRC risk that replicate (i.e., attain a Bonferroni significance threshold) in our independent custom replication genotyping experiment, and 2) new association signals for CRC risk based on all available data. Since 2) is based on the largest CRC GWAS dataset (>125,000 individuals), a well-powered replication dataset is currently not available.
Randomization	To limit potentially confounding experimental batch effects, we sequenced and genotyped cases and controls together and randomized samples across flow cells or plates. Likewise, variant calling, phasing and imputation were performed for cases and controls together. Since this is an observational genetic association study, randomization to experimental groups is not relevant.
Blinding	The laboratories conducting the sequencing or genotyping did not have access to phenotype data. Genotype calling was also blinded with respect to case-control status.

## Reporting for specific materials, systems and methods

## Materials &amp; experimental systems

n/a	Involvement	Involved in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/>	Unique biological materials
<input checked="" type="checkbox"/>	<input type="checkbox"/>	Antibodies
<input checked="" type="checkbox"/>	<input type="checkbox"/>	Eukaryotic cell lines
<input checked="" type="checkbox"/>	<input type="checkbox"/>	Palaeontology
<input checked="" type="checkbox"/>	<input type="checkbox"/>	Animals and other organisms
<input type="checkbox"/>	<input checked="" type="checkbox"/>	Human research participants

## Methods

n/a	Involvement	Involved in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/>	ChIP-seq
<input checked="" type="checkbox"/>	<input type="checkbox"/>	Flow cytometry
<input checked="" type="checkbox"/>	<input type="checkbox"/>	MRI-based neuroimaging

## Human research participants

Policy information about [studies involving human research participants](#)

## Population characteristics

Participants were adult colorectal cancer (CRC) or advanced adenoma patients, or healthy controls. The composition of the study in terms of sex and age reflects the underlying study populations and demographic characteristics for each participating study are summarized in Supplementary Table 1.

## Recruitment

This is an observational case control GWAS study for CRC. Details of participant enrollment, selection of matching healthy controls, and inclusion and exclusion criteria vary between studies and are described in the Supplementary Text for each of the 45 studies (for newly generated data), or references are given in the Online Methods (for published data). In brief, all CRC or advanced adenoma cases were histologically or clinically confirmed, and controls with a known history of cancer or reported family history of CRC were excluded.