

# Fine-mapping analysis including over 254,000 East Asian and European descendants identifies 136 putative colorectal cancer susceptibility genes

Received: 16 October 2023

Accepted: 26 March 2024

Published online: 26 April 2024

 Check for updates

A list of authors and their affiliations appears at the end of the paper

Genome-wide association studies (GWAS) have identified more than 200 common genetic variants independently associated with colorectal cancer (CRC) risk, but the causal variants and target genes are mostly unknown. We sought to fine-map all known CRC risk loci using GWAS data from 100,204 cases and 154,587 controls of East Asian and European ancestry. Our stepwise conditional analyses revealed 238 independent association signals of CRC risk, each with a set of credible causal variants (CCVs), of which 28 signals had a single CCV. Our cis-eQTL/mQTL and colocalization analyses using colorectal tissue-specific transcriptome and methylome data separately from 1299 and 321 individuals, along with functional genomic investigation, uncovered 136 putative CRC susceptibility genes, including 56 genes not previously reported. Analyses of single-cell RNA-seq data from colorectal tissues revealed 17 putative CRC susceptibility genes with distinct expression patterns in specific cell types. Analyses of whole exome sequencing data provided additional support for several target genes identified in this study as CRC susceptibility genes. Enrichment analyses of the 136 genes uncover pathways not previously linked to CRC risk. Our study substantially expanded association signals for CRC and provided additional insight into the biological mechanisms underlying CRC development.

Colorectal cancer (CRC) is one of the most common malignancies worldwide<sup>1</sup>. Inherited genetic factors play an important role in the development of CRC<sup>2</sup>. Since 2007, genome-wide association studies (GWAS) have identified over 200 common genetic variants independently associated with CRC risk<sup>3–7</sup>. These GWAS, however, typically only reported the most significantly associated variant (the lead variant) at each risk locus. Statistical fine-mapping analyses of known risk loci can identify additional association signals independent of the lead variant.

Approximately 90% of GWAS-identified risk variants for CRC are located in noncoding or intergenic regions, and target genes for most

of these risk variants remain unknown. Well-powered fine-mapping analyses, particularly those using data from multi-ancestry populations, can facilitate the identification of credible causal variants (CCVs) in each region. Previous genetic studies have provided strong evidence that regulatory variants in linkage disequilibrium (LD) with GWAS-identified risk variants drive the associations of genetic variants with cancer risk by modulating the expression of susceptibility genes<sup>8–11</sup>. Therefore, integrating functional genomic data to interrogate CCVs in each independent risk-associated signal could help to identify putative causal variants and target genes for CRC risk. Herein, we conducted large trans-ancestry fine-mapping analyses of all currently known CRC

✉ e-mail: [wei.zheng@vanderbilt.edu](mailto:wei.zheng@vanderbilt.edu)

risk regions, using GWAS data from 100,204 CRC cases and 154,587 controls of East Asian and European ancestry, to identify independent association signals and their target genes for CRC risk.

## Results

### Identification of independent association signals with CRC risk

We conducted fine-mapping analyses using GWAS summary statistics from 100,204 CRC cases and 154,587 controls (73% European and 27% East Asian ancestry) (Fig. 1, Supplementary Data 1). In our recent trans-ancestry meta-analysis of GWAS, we identified 205 genetic variants independently associated with CRC risk<sup>7</sup>. We aggregated regions flagged by these variants into 143 risk regions, each containing at least a 1 Mb interval centered on the most significant association (Supplementary Data 2). Among them, 40 regions harbor at least two reported independent risk associations. All risk regions were autosomal, except the one at Xp22.2. For subsequent analyses, we focused on the 142 regions located on the autosomes.

We used forward stepwise conditional analyses to identify independent association signals in each region in each population, conditioning on the most significant association from the trans-ancestral summary statistics (Supplementary Fig. 1, Methods). We then meta-analyzed the conditioned data using the fixed-effects inverse variance weighted model. We considered the threshold of conditional  $P < 1 \times 10^{-6}$  to determine independent significant associations to balance both Type 1 and 2 errors, as recommended by a previous fine-mapping study in breast cancer<sup>12</sup>. At this threshold, we identified 171 independent association signals in 122 regions (Fig. 2, Supplementary Data 3). To identify possible ancestry-specific association signals, we conducted similar analyses using only summary statistics from each

population, conditioning on the ancestry-specific most significant association. Using the same threshold, we identified 198 and 45 independent association signals in European and East Asian descendants, respectively (Supplementary Data 4 and 5). Of them, 60 signals in European and 7 in East Asian were not detected in the trans-ancestry analysis above, suggesting them as potential ancestry-specific risk signals (Fig. 2).

In total, we identified 238 independent association signals either from trans-ancestry or ancestry-specific analysis at these 142 regions (Fig. 2). A total of 94 regions (66.2%) contained only a single association signal, while the remaining 48 regions (33.8%) consisted of multiple independent association signals. Among the 238 independent association signals, 191 signals had lead variants that were correlated with previously GWAS-reported risk variants<sup>7</sup> ( $LD\ r^2 > 0.1$  in either of East Asian or European-ancestry population). The remaining 47 independent signals (19.7%) have not been previously reported, including 18 from trans-ancestry, 28 from European-specific, and one from East Asian-specific analyses (Fig. 2, Table 1). Among these 47 signals, 31 demonstrated significant associations with conditional  $P < 1 \times 10^{-7}$ , including 28 signals reached genome-wide significance.

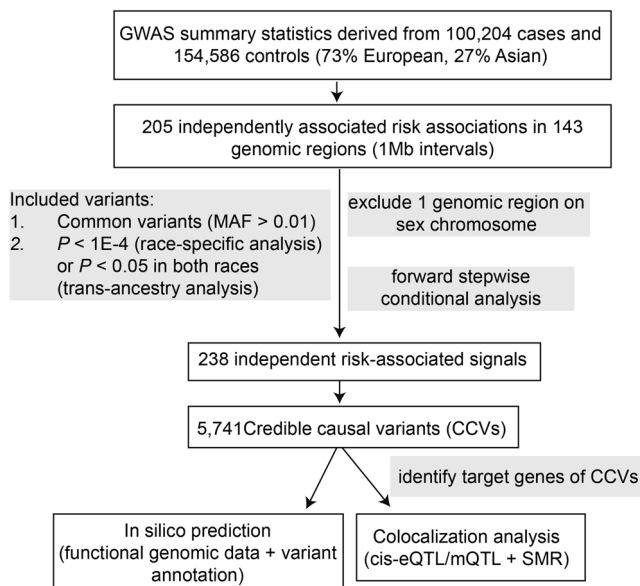
### Identification of credible causal variants (CCVs) for independent association signals

To identify CCVs for each independent association signal, we conducted conditional analysis with adjustment of the lead variants for other signals in the same risk region. We conducted this analysis for trans-ancestral independent signals separately for each population to account for differences in the LD structure and then meta-analyzed conditioned results. Using a similar approach conducted in breast cancer<sup>12</sup>, we defined variants as CCVs if they satisfied conditional  $P$  values within two orders of magnitude of the most significant association, conditioning on all other independent association signals. We identified a total of 5741 CCVs for the 238 signals, with the number of CCVs per signal ranging from 1 to 249 (median: 11 CCVs per signal) (Supplementary Data 6). For 28 risk signals, only a single CCV was identified, suggesting that these CCVs are likely to be the causal variants for these signals (Table 2).

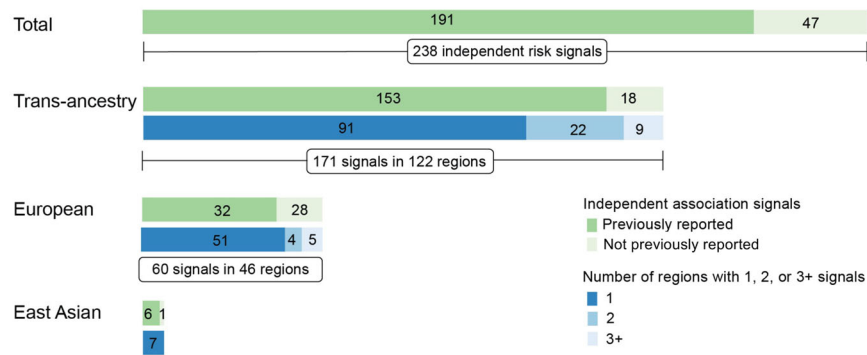
For the 138 independent association signals identified in both trans-ancestry and European-ancestry specific analyses (Supplementary Data 7), trans-ancestry analyses identified a smaller-sized set of CCVs (mean = 23.2, median = 8.5), compared with European-ancestry specific analysis (mean = 31.08, median = 15) (paired Wilcoxon test,  $P = 4.9 \times 10^{-7}$ ). Interestingly, a single CCV was identified for 10 signals in trans-ancestry analysis, while multiple CCV for them in European-ancestry specific analysis, highlighting the value of using multi-ancestry data to reduce the number of CCVs in fine-mapping analysis. For instance, signal 1 in region\_42 included 16 CCVs in the European set (lead variant: rs41302867), but only one variant in the trans-ancestry set (rs9379084). The variant rs9379084 is a predicted-deleterious missense variant (p.Asp1171Asn) of the *RREB1* gene which plays a regulatory role in Ras/Raf-mediated cell differentiation<sup>13</sup>, a pathway well known to be implicated in CRC development.

### Identification of target genes for CCVs

Of the 5741 CCVs identified in this study, 3716 (64.7%) are located in regions with at least one of six genomic features (open chromatin, transcribed regions of active genes, promoter, enhancer, repressed gene regulatory elements, and transcription factor (TF) binding sites) (Supplementary Data 6 and 8). To identify putative target genes of these CCVs, we used functional genomic data generated in CRC-related tissues/cells to conduct in-silico analyses with a modified INQUISIT pipeline<sup>12</sup> (Methods, Supplementary Data 9). We identified 72 putative target genes via CCVs located in distal enhancer elements (Supplementary Data 10), 48 genes via CCVs located in proximal promoter elements (Supplementary Data 11), and 19 genes that could be



**Fig. 1 | Schematic diagram of the study design.** We conducted fine-mapping analyses using GWAS summary statistics from 100,204 cases and 154,587 controls. All 205 genetic variants were aggregated to 143 risk regions containing at least a 1 megabase (Mb) interval centered on the most significant association. This study focused on 142 risk regions located on the autosomes. In forward stepwise conditional analysis, we included common variants (minor allele frequency (MAF) > 0.01) with associations at  $P < 0.05$  in both populations for the trans-ancestry analysis and with associations at  $P < 1 \times 10^{-4}$  in each population for race-specific analysis. The threshold of conditional  $P < 1 \times 10^{-6}$  was used to determine independent risk-associated signals. For credible causal variants (CCVs) for each independent signal, we conducted *in-silico* analyses with functional genomic data generated in CRC-related tissues/cells and colocalization of expression/methylation quantitative trait loci (e/mQTL) with GWAS signals to identify putative target genes for CCVs using the Summary-data-based Mendelian Randomization (SMR) approach.



**Fig. 2 | Independent association signals for colorectal cancer risk.** Numbers of fine-mapping regions and numbers of independent association signals identified through forward stepwise conditional analyses. The second bar for “Trans-ancestry”, “European” and “East Asian” also shows the number of regions with 1, 2, or 3+

signals per region. The green color indicates the number of independent association signals previously reported or not yet reported. The blue color indicates the number of independent association signals in each risk region.

targeted by CCVs in coding regions (i.e., deleterious missense, stop-gained, and start\_lost) (Supplementary Data 12). In total, we identified 128 genes associated with CCVs for 76 independent association signals, with a range from one to five putative target genes per signal. Of them, 52 independent association signals contain only a single putative target gene.

We also conducted cis-expression quantitative trait loci (cis-eQTL) analyses to identify target genes using four transcriptome datasets derived from either normal colon tissues or tumor-adjacent normal colon tissues from 1299 individuals from the Genotype-Tissue Expression (GTEx) project ( $n = 368$  individuals predominantly of European ancestry), the BarcUVA-Seq project ( $n = 423$  individuals of European ancestry), the Colonomics project ( $n = 144$  individuals of European ancestry), and the Asia Colorectal Cancer Consortium (ACCC) ( $n = 364$  individuals of East Asian ancestry) (Methods). At Bonferroni-corrected  $P < 0.05$ , we identified 153 genes associated with the lead variants, including 127 genes in 65 independent association signals and 30 in 15 signals identified from trans-ancestry and European-ancestry specific analyses, respectively. We also identified the *PPPIR21* gene in a potential Asian-specific risk signal (lead variant rs77272589) (Supplementary Data 13). Out of the 153 genes, 37 had been previously identified by eQTL analysis<sup>5,10,11</sup>. For independent association signals identified in European and trans-ancestry analyses, we further performed cis-methylation quantitative trait loci (cis-mQTL) analyses using two methylation datasets generated from 321 individuals from the GTEx project ( $n = 189$  individuals predominantly of European ancestry) and the Colonomics project ( $n = 132$  individuals of European ancestry). We found that DNA methylation levels at CpG sites for 84 genes were associated with 71 independent association signals, including 14 genes identified in previous mQTL analysis<sup>11</sup> (Supplementary Data 14).

We next conducted colocalization analyses for identified likely target genes in significant eQTL/mQTLs above using the Summary-data-based Mendelian Randomization (SMR) approach (Methods). Through the integration of eQTL/mQTL results and GWAS associations signals, we identified 205 genes at Bonferroni-corrected  $P_{SMR} < 0.05$  (Supplementary Data 15–19), including 150 genes from the eQTL analysis and 84 genes from the mQTL analysis. Of these, 45 (21.9%) genes were also identified as targets of CCVs by in-silico analyses based on functional genomic data as described above, and 29 genes were identified in both mQTL and eQTL analyses. That is in line with previous observations in the overlap fraction between mQTL and eQTL<sup>14</sup>. We considered genes with evidence of only mQTL colocalization, as the enrichment of mQTLs in gene regulatory elements, as well as their implications in other molecular phenotypes, such as chromatin accessibility<sup>14,15</sup>. Notably, of the 55 genes only identified in the mQTL

analysis, seven genes were supported by the above in silico analyses with functional genomic data, and 22 genes showed association with CRC risk in previous TWAS and eQTL colocalization analysis<sup>7,11,16,17</sup>.

In total, we identified 288 putative target genes for 140 independent association signals based on functional genomics data and/or colocalization analysis. For 35 of these signals, multiple target gene candidates were detected per signal, suggesting that some may be false positives (Supplementary Data 20). To minimize false positive findings, we further prioritized target gene candidates by analyzing associations of genes with CRC risk based on previous transcriptome-wide association studies (TWAS) and colocalizations between eQTL and CRC GWAS signals<sup>7,11,16,17</sup> (Methods). Finally, we obtained a credible set of 136 protein-coding genes for 124 independent association signals. Among them, 56 genes were not previously identified as potential targets for CRC risk associations, including nine genes in eight previously unreported association signals in this study (Table 3). The remaining 80 genes were previously reported as potential CRC susceptibility genes, and our study provided additional supporting evidence (Table 4)<sup>7,11,16,17</sup>.

### Using scRNA-seq data to evaluate gene expression pattern by cell types

To investigate potential underlying cell types of putative susceptibility genes that contribute to CRC development, we analyzed single-cell RNA-seq (scRNA-seq) datasets from normal colon tissues obtained from 31 participants included in the Colorectal Molecular Atlas Project<sup>18</sup> (Methods). Of the 136 identified genes, 17 genes exhibited significantly differential expression in specific cell types compared to the other cell types at  $|\log_2 \text{fold change (FC)}| > 1$  and a nominal  $P < 0.05$  (Supplementary Data 21). Nine of these genes (*DIP2B*, *CIB1*, *HPGD*, *CDKN2B*, *TMEM258*, *MYL12A*, *MYL12B*, *CDKN1A*, and *TMBIM1*) showed a distinct expression pattern in specific absorptive cells (ABS) cell, underscoring the relevance of this cell type underlying CRC development.

### Using whole exome sequencing data to evaluate pathogenic variants in target genes with CRC risk

We used whole exome sequencing data from 3362 CRC cases and 133,742 controls of European ancestry in the UK Biobank (UKBB) to evaluate the association of CRC risk with putative candidate genes identified our study using burden tests by aggregating either loss of function (pLOF) or pLOF and deleterious missense variants (Dmis) jointly in each gene (Methods). Of these 136 genes, *MLH1* was significantly associated with CRC risk with  $P = 1.35 \times 10^{-7}$  when considering only pLOF in tests (at Bonferroni-corrected threshold,  $0.05/136$  testing). Additional nine genes (*TNFSF18*, *LRP1*, *SMAD9*, *PDGFB*, *CIB1*,

**Table 1 | Independent association signals uncovered at known CRC risk loci in conditional analyses (conditional  $P < 1 \times 10^{-6}$ )**

Fine-mapping region	SNP	Chr	Position	Nearby gene	Alleles	AF	Single-SNP analysis		Joint analysis		Group
							OR (95% CI)	P value <sup>a</sup>	OR (95% CI)	P value <sup>b</sup>	
region_1	rs11579545	1	22249333	HSPG2	T/C	0.445	0.96 (0.95–0.98)	4.34E–07	0.96 (0.95–0.98)	5.63E–07	Trans-ancestry
region_1	rs112191583	1	22554378	MIR4418	T/C	0.974	0.88 (0.83–0.92)	1.19E–07	0.87 (0.83–0.92)	5.29E–08	Trans-ancestry
region_1	rs12137525	1	22584118	MIR4418	T/C	0.107	1.07 (1.04–1.09)	2.90E–08	1.08 (1.06–1.11)	1.14E–11	European
region_9	rs12122827	1	202172769	LGR6	T/G	0.715	1.04 (1.02–1.06)	9.44E–07	1.05 (1.03–1.06)	7.94E–08	European
region_22	rs2554878	3	41200064	RPT1-372H2.1	T/G	0.036	1.12 (1.08–1.16)	5.85E–09	1.12 (1.08–1.16)	3.75E–09	Trans-ancestry
region_27	rs9283588	3	133874566	RYK	A/G	0.715	1.06 (1.04–1.07)	3.43E–10	1.04 (1.03–1.06)	7.32E–07	Trans-ancestry
region_30	rs902443	4	105888417	RPT1-556H14.1	A/T	0.536	1.04 (1.03–1.06)	1.49E–11	1.04 (1.03–1.06)	1.26E–11	Trans-ancestry
region_36	rs582489	5	39908712	GCSHP1	T/C	0.570	0.97 (0.96–0.99)	8.23E–05	0.96 (0.94–0.97)	7.29E–09	European
region_36	rs77781678	5	40626064	SNORA63	T/C	0.020	0.84 (0.79–0.89)	2.09E–09	0.84 (0.79–0.89)	1.75E–09	European
region_43	rs4714081	6	11977905	RPT1-456H18.1	A/G	0.451	0.96 (0.95–0.97)	2.50E–09	0.96 (0.95–0.97)	1.21E–09	Trans-ancestry
region_43	rs4714350	6	12270290	EDN1	A/T	0.283	0.96 (0.94–0.97)	8.60E–09	0.96 (0.95–0.98)	4.43E–07	Trans-ancestry
region_43	rs17615624	6	12376025	RN7SKP293	C/G	0.975	0.87 (0.83–0.91)	7.29E–09	0.88 (0.84–0.92)	2.28E–07	European
region_44	rs3094576	6	29516242	OR21P	A/C	0.131	0.94 (0.92–0.96)	1.83E–07	0.94 (0.92–0.96)	2.26E–08	European
region_44	rs2517671	6	29937977	MICD	A/G	0.591	0.96 (0.95–0.98)	2.35E–08	0.96 (0.95–0.97)	2.87E–09	Trans-ancestry
region_45	rs6920820	6	30969938	MUC22	C/G	0.980	0.84 (0.79–0.9)	6.87E–08	0.8 (0.75–0.85)	1.89E–12	European
region_45	rs9264180	6	31219902	HLA-C	A/C	0.570	1.03 (1.02–1.05)	1.71E–06	1.04 (1.02–1.05)	5.62E–07	Trans-ancestry
region_45	rs9265501	6	31297568	XXbac-BPG248L24.10	A/G	0.678	0.88 (0.85–0.92)	3.05E–10	0.88 (0.84–0.91)	5.21E–11	European
region_45	rs116000952	6	32541270	HLA-DRB1	T/G	0.843	0.92 (0.89–0.96)	5.74E–06	0.9 (0.87–0.94)	1.50E–08	European
region_45	rs2858331	6	32681277	XXbac-BPG254F23.7	A/G	0.601	1.03 (1.02–1.05)	1.18E–05	1.05 (1.04–1.07)	2.67E–12	Trans-ancestry
region_50	rs13204733	6	55566108	RPT1-228O6.2	A/G	0.858	0.94 (0.92–0.96)	4.20E–08	0.93 (0.91–0.95)	1.17E–13	European
region_60	rs10089517	8	60178721	SNORA51	A/C	0.380	1.03 (1.02–1.05)	7.44E–07	1.03 (1.02–1.05)	2.81E–07	Trans-ancestry
region_61	rs117310502	8	117593052	EIF3H	A/G	0.048	0.92 (0.89–0.96)	9.36E–05	0.88 (0.85–0.92)	4.03E–10	European
region_61	rs72681666	8	117641754	EIF3H	T/C	0.043	1.09 (1.05–1.13)	1.57E–05	1.12 (1.08–1.17)	6.99E–10	European
region_61	rs1793717	8	118278575	SNORA31	A/C	0.629	1.03 (1.02–1.05)	6.90E–05	1.04 (1.03–1.06)	1.55E–07	European
region_62	rs79122086	8	128397907	CASC8	T/G	0.840	0.92 (0.9–0.93)	5.46E–20	0.94 (0.93–0.96)	9.34E–10	Trans-ancestry
region_62	rs77569096	8	128468955	CASC8	A/G	0.763	0.92 (0.9–0.94)	2.06E–15	0.93 (0.91–0.95)	2.67E–12	European
region_68	rs4994332	9	137117194	RPT1-145E17.2	T/C	0.423	0.97 (0.96–0.98)	4.05E–05	0.96 (0.95–0.97)	9.08E–08	European
region_74	rs117746067	10	101222300	RPT1-441O15.3	A/G	0.101	1.06 (1.03–1.08)	3.64E–06	1.08 (1.05–1.1)	1.74E–09	European
region_80	rs9795065	11	74376844	POLD3	T/C	0.981	1.19 (1.13–1.25)	5.37E–13	1.17 (1.12–1.23)	6.06E–11	Trans-ancestry
region_85	rs1003563	12	6424577	PLEKHG6	A/G	0.433	0.95 (0.94–0.97)	1.67E–12	0.95 (0.94–0.96)	1.23E–14	Trans-ancestry
region_106	rs68097734	14	92717447	RPT1-472N19.3	T/C	0.496	1.06 (1.03–1.09)	7.71E–06	NA	–	Asian
region_108	rs28630996	15	32993860	SCG5	A/T	0.713	0.9 (0.89–0.92)	1.25E–32	0.93 (0.91–0.94)	3.02E–17	Trans-ancestry
region_108	rs144674978	15	33149751	FMN1	T/C	0.013	1.34 (1.25–1.43)	1.11E–18	1.23 (1.15–1.31)	3.82E–10	European
region_109	rs3784710	15	68072458	MAP2K5	T/C	0.763	1.05 (1.03–1.07)	1.32E–07	1.05 (1.03–1.07)	1.34E–08	European
region_111	rs12913420	15	90797010	RPT1-697E2.6	C/G	0.376	1.04 (1.03–1.06)	2.29E–09	NA	–	Trans-ancestry
region_114	rs1117455	16	86179919	RPT1-805I24.4	T/C	0.181	1.04 (1.02–1.06)	7.52E–06	1.05 (1.03–1.07)	6.97E–07	European
region_115	rs73975588	17	816741	NXN	A/C	0.874	1.09 (1.07–1.12)	6.62E–16	1.07 (1.04–1.09)	6.08E–09	European
region_117	rs112592783	17	70633625	LINC00511	T/C	0.175	1.05 (1.03–1.07)	5.87E–09	1.05 (1.03–1.07)	5.95E–09	Trans-ancestry
region_120	rs4939821	18	46371993	CTIF	T/C	0.304	0.91 (0.89–0.92)	4.12E–32	0.96 (0.94–0.97)	1.89E–07	European

**Table 1 (continued) | Independent association signals uncovered at known CRC risk loci in conditional analyses (conditional  $P < 1 \times 10^{-6}$ )**

Fine-mapping region	SNP	Chr	Position	Nearby gene	Alleles	AF	Single-SNP analysis		Joint analysis		Group
							OR (95% CI)	P value <sup>a</sup>	OR (95% CI)	P value <sup>b</sup>	
region_121	rs72971616	19	987366	WDR18	T/G	0.063	0.92 (0.89–0.95)	4.89E–06	NA	–	European
region_126	rs12460535	19	49098750	SULT2B1	A/G	0.349	0.96 (0.95–0.98)	5.76E–07	NA	–	European
region_127	rs8099852	19	58895221	RPS5	T/C	0.546	1.03 (1.02–1.05)	5.15E–07	NA	–	Trans-ancestry
region_133	rs1971480	20	48897080	RPT1-290F20.3	T/G	0.672	0.96 (0.95–0.98)	8.61E–07	0.96 (0.94–0.97)	6.50E–09	European
region_133	rs149942633	20	48983073	RPT1-290F20.2	T/C	0.153	1.12 (1.08–1.16)	1.93E–08	1.1 (1.05–1.14)	3.94E–06	European
region_133	rs6126008	20	49075315	COX6CP2	A/T	0.660	0.96 (0.94–0.97)	1.07E–08	0.96 (0.94–0.97)	3.74E–09	European
region_134	rs34161672	20	56020599	RBW38	A/G	0.321	1.04 (1.02–1.05)	2.40E–06	1.04 (1.03–1.06)	1.64E–07	European
region_142	rs78106213	22	46121230	ATXN10	T/G	0.693	1.04 (1.03–1.06)	2.41E–07	1.05 (1.03–1.06)	8.20E–09	European

All independent association signals presented in this table are those not previously reported.  
 Chr and Position GRCh37, Alleles risk allele/Reference allele, AF Allele frequency, OR odds ratio, CI confidence interval.  
<sup>a</sup>P value derived from trans-ancestry or ancestry-specific meta-analysis under the fixed-effects inverse variance weighted model.  
<sup>b</sup>P value derived from conditional analysis conditioning on all other independent association signals in each fine-mapping region. “NA”—Only a single association signal was detected in the fine-mapping region in the analysis group.

*STK39*, *IGFBP3*, *FUT2*, and *FUT3*) showed nominal  $P < 0.05$  significance considering only pLoF or combination of pLoF and Dmis, whereas no significance was detected for the remaining genes.

**Biological significance of the target genes for CCVs**

We utilized Enrichr<sup>19–21</sup> to analyze multiple pathway databases and identify enriched biological pathways among the 136 credible target genes (Methods). At a false-discovery rate (FDR)  $< 0.05$ , 126 pathways showed significant enrichment (Supplementary Data 22). Our findings were in line with our prior study<sup>18</sup> and highlighted the enriched signaling pathways such as TGF- $\beta$ , BMP, Wnt, Hippo, and TNF- $\alpha$ /NF- $\kappa$ B, which are known to play a crucial role in the development and progression of colorectal cancer<sup>19,20</sup>. Of the 56 genes not previously reported, nine genes (*TGIF1*, *CDKN2B*, *MYC*, *BMP7*, *WNT7B*, *PRICKLE2*, *LGR6*, *CEBPB*, and *IRS2*) were mapped to these pathways (Table 5). Additionally, we identified several significant pathways, including those related to cancer, pluripotency of stem cells, epithelial–mesenchymal transition, extracellular matrix organization, adipogenesis, senescence, and autophagy in cancer. Interestingly, we also identified the glycolysis pathway, which provides energy support for cancer cells, as a significant pathway not previously reported. Four previously unreported genes, *GOT1*, *IGFBP3*, *IRS2*, and *LCT*, were mapped to glycolysis, supporting their association with CRC risk.

In addition, we performed functional annotation analysis on each credible target gene and assigned them to previously described cellular processes<sup>18</sup> (Supplementary Fig. 2). Of the 56 genes not previously reported, 26 were found to be involved in these cellular processes. Specifically, five genes were related to stemness/differentiation, one gene was linked to adhesion/migration, and six genes were associated with proliferation. Interestingly, we also identified an additional cellular process, post-translation modifications (PTMs) of protein, which included three genes (*DACF12*, *USP12*, and *SENP8*). These findings suggest potential critical roles of PTMs in the development of CRC.

**Discussion**

Our study, including approximately 254,000 individuals of East Asian and European ancestry, represents the largest study conducted to fine-map CRC risk-associated genomic regions using GWAS data. We identified 238 independent association signals at conditional  $P$  value  $< 1 \times 10^{-6}$ , including 47 signals not reported previously. Furthermore, integrating functional genomic data and results from cis-eQTL/mQTL and colocalization analyses, we identified 136 putative CRC susceptibility genes, including 56 genes that had not been previously reported. Notably, these identified genes are significantly enriched in several major CRC signaling pathways and other cancer-related pathways. Our findings not only significantly expanded the number of associated signals for CRC, but also provide substantial data to advance our understanding of CRC biology.

The integration of comprehensive functional genomic data from relevant colon tissues and cell lines, as well as genetic associations data, can facilitate the identification of potential target genes for CRC risk. Our study significantly extends previous efforts<sup>7,11,16,17</sup> by identifying 56 target gene candidates not previously reported for CRC risk, over half of which (29/56, 51.8%) are involved in the enriched biological pathways. For instance, eight target genes (*TGIF1*, *CDKN2B*, *LGR6*, *MYC*, *PRICKLE2*, *WNT7B*, *BMP7*, and *TBX3*) identified in this study may regulate normal intestinal homeostasis as they play roles in signaling pathways (i.e., Wnt and BMP) and pluripotency of stem cells. *LGR6*, for instance, is part of a G-protein-coupled receptor family and marks stem cells in the epidermis<sup>22</sup>. It activates a novel  $\beta$ -catenin/TCF7L2/*LGR6*-positive feedback loop in *LGR6*<sup>high</sup> cervical cancer stem cells (CSCs) to enhance the properties of cancer stem cells, including self-renewal, differentiation, and tumorigenicity<sup>23</sup>. Silencing of *LGR6* resulted in the inhibition of stemness by repressing Wnt/ $\beta$ -catenin signaling in ovarian cancer<sup>24</sup>. *TBX3*, a transcriptional repressor, regulates stem cell maintenance by

**Table 2 | Independent association signals with a single CCV**

Fine-mapping region	SNP	Chr	Position	Alleles	AF	OR (95% CI)	P value <sup>a</sup>	Putative target gene(s) <sup>b</sup>
<i>European-specific analysis</i>								
region_45	rs116000952	6	32541270	T/G	0.843	0.92 (0.89–0.96)	5.74E–06	–
region_45	rs6920820	6	30969938	C/G	0.980	0.84 (0.79–0.90)	6.87E–08	<i>LINCO0243</i>
region_61	rs72681666	8	117641754	T/C	0.043	1.09 (1.05–1.13)	1.57E–05	–
region_62	rs77569096	8	128468955	A/G	0.763	0.92 (0.90–0.94)	2.06E–15	–
region_84	rs3217810	12	4388271	T/C	0.127	1.13 (1.11–1.16)	1.96E–26	–
region_108	rs144674978	15	33149751	T/C	0.013	1.34 (1.25–1.43)	1.11E–18	–
region_133	rs149942633	20	48983073	T/C	0.153	1.12 (1.08–1.16)	1.93E–08	–
<i>Trans-ancestry analysis</i>								
region_1	rs112191583	1	22554378	T/C	0.974	0.88 (0.83–0.92)	1.19E–07	–
region_24	rs704417	3	64252424	T/C	0.546	1.05 (1.03–1.06)	4.35E–10	–
region_27	rs113569514	3	133748789	T/C	0.763	1.08 (1.07–1.10)	1.92E–21	<i>SLCO2A1</i>
region_29	rs2578155	4	94836291	C/G	0.503	1.04 (1.03–1.06)	1.09E–09	–
region_42	rs9379084	6	7231843	A/G	0.144	0.93 (0.91–0.95)	2.39E–12	<i>RREB1</i>
region_46	rs16878812	6	35569562	A/G	0.892	1.09 (1.07–1.12)	7.62E–15	<i>FKBP5</i>
region_48	rs6933790	6	41672769	T/C	0.788	1.08 (1.06–1.10)	2.66E–20	–
region_61	rs4129064	8	117735666	T/G	0.734	1.06 (1.04–1.07)	1.01E–09	–
region_62	rs6983267	8	128413305	T/G	0.508	0.86 (0.85–0.87)	1.65E–122	<i>MYC</i>
region_72	rs704017	10	80819132	A/G	0.473	0.92 (0.91–0.93)	1.97E–38	–
region_84	rs12818766	12	4376091	A/G	0.215	1.10 (1.08–1.12)	1.81E–29	–
region_89	rs7398375	12	57540848	C/G	0.651	1.07 (1.05–1.09)	3.70E–19	<i>LRP1</i>
region_94	rs11067228	12	115094260	A/G	0.560	0.95 (0.94–0.97)	2.50E–13	–
region_96	rs116964464	13	27543193	T/C	0.035	1.11 (1.07–1.15)	4.83E–09	<i>USP12</i>
region_99	rs7325844	13	73625133	A/G	0.639	1.05 (1.04–1.07)	1.28E–12	–
region_104	rs35107139	14	54419106	A/C	0.550	0.92 (0.91–0.93)	4.22E–36	–
region_105	rs8020436	14	59208437	A/G	0.370	1.06 (1.05–1.08)	1.27E–17	–
region_108	rs17816465	15	33156386	A/G	0.193	1.09 (1.07–1.10)	5.73E–20	–
region_116	rs1078643	17	10707241	A/G	0.765	1.09 (1.07–1.11)	2.31E–27	–
region_132	rs6066825	20	47340117	A/G	0.662	1.08 (1.07–1.10)	2.13E–32	–
region_136	rs1741640	20	60932414	T/C	0.208	0.88 (0.86–0.89)	8.15E–55	<i>LAMA5, CABLES2</i>

Chr and Position GRCh37, Alleles risk allele/Reference allele, AF Allele frequency, OR odds ratio, CI confidence interval. <sup>a</sup>P value derived from trans-ancestry or European-ancestry meta-analysis under the fixed-effects inverse variance weighted model; <sup>b</sup>– No target genes were prioritized for the variant in this study.

controlling stem cell self-renewal and differentiation, and reduced expression levels of *TBX3* are associated with reduced pluripotency of stem cells<sup>25,26</sup>. *MYC* and *WNT7B* are implicated in the signaling related to the self-renewal and differentiation of cancer stem cells<sup>27</sup>. Here, we linked *MYC* and *WNT7B* with credible causal variants of CRC risk associations through functional genomic interaction. Our findings also indicated the relevance of glycolysis to CRC risk associations, a metabolic pathway critical in early CRC tumorigenesis by supporting the energetic and biosynthetic demands of CRC cells<sup>28,29</sup>. It should be noted that future studies are needed to validate chromatin interactions between identified CCVs and their target genes in this study by employing chromatin conformation capture technology such as in situ Hi-C, Capture Hi-C (CHi-C), and HiChIP.

Additional evidence supports some of the candidate target genes identified in our study as possible CRC susceptibility genes. In our differential gene expression analysis among normal colon mucosa, adenoma, and adenocarcinoma using gene expression data from 135 normal colon mucosae, 218 colon adenomas, and 2760 colon adenocarcinomas, we observed that 26 genes showed significant differential expression between adenoma and normal colon tissues, while 31 genes showed significant differential expression between carcinoma and adenoma tissues (adjusted  $P < 0.05$ ) (Supplementary Data 20). Interestingly, three stemness/differentiation-related genes, including *LRRC34*, *CEBPB*, and *TBX3*, showed significant changes in their expression levels in adenoma compared to normal colon

mucosa. Additionally, 34 (60.7%) of not previously identified genes have been implicated in cancer-related functions in in vitro or in vivo functional experimental studies in CRC or other cancer types (Supplementary Data 20). These results provide further evidence supporting the potential involvement of these genes in CRC progression. Despite the above supportive evidence, it remains necessary to evaluate the functions of identified putative CRC susceptibility genes through both in vitro and in vivo assays in future investigations.

The trans-ancestry and ancestry-specific fine-mapping analyses conducted in this study not only enabled the discovery of independent association signals that are shared across populations of European and East Asian ancestry, but also revealed ancestry-specific signals. The larger sample size of the European-ancestry study enabled us to identify a larger number of independent association signals than the study conducted on Asians. However, there are some ancestry-specific signals identified in this study, which is most likely due to differences in LD structures and allele frequency between these two populations. Indeed, we observed distinct differences in the allele frequency for most ancestry-specific signals, as shown in Supplementary Data 4 and 5. For instance, the lead variant of 24 European ancestry-specific signals (40%, 24/60) is not detected among East Asian-ancestry populations. On the other hand, fine-mapping analyses capitalizing on ancestry differences in LD structure can substantially reduce the credible set size compared to European-ancestry specific analysis. This highlights the value of multi-ancestry fine-mapping over

**Table 3 | The 56 CRC susceptibility gene candidates not previously reported**

Fine-mapping region	Gene	Lead variant	Distal	Proximal	Coding	Colocalization (eQTL)	Colocalization (mQTL)
region_1	<i>CELA3B</i>	rs11579545				+	
region_1	<i>HSPG2</i>	rs11579545	+			+	+
region_5	<i>PTGER3</i>	rs2651244				+	
region_7	<i>TNFSF18</i>	rs10489274					+
region_9	<i>LGR6</i>	rs12122827					+
region_10	<i>CNTN2</i>	rs12078075		+			+
region_12	<i>FMN2</i>	rs2078095				+	
region_14	<i>PPP1R21</i>	rs77272589		+		+	
region_16	<i>LCT</i>	rs1446585					+
region_21	<i>GOLGA4</i>	rs1800734				+	
region_21	<i>MLH1</i>	rs1800734		+		+	
region_24	<i>ADAMTSS9</i>	rs6445418					+
region_24	<i>PRICKLE2</i>	rs704417				+	
region_27	<i>SLCO2A1</i>	rs113569514		+			
region_28	<i>LRRC34</i>	rs10936599			+		
region_28	<i>ACTRT3</i>	rs10936599	+	+			
region_28	<i>MYNN</i>	rs10936599	+		+		
region_34	<i>HPGD</i>	rs1426947				+	
region_42	<i>LY86</i>	rs1294438					+
region_44	<i>OR211P</i>	rs73402748				+	
region_46	<i>SRPK1</i>	rs16878812				+	
region_49	<i>RUNX2</i>	rs57939401					+
region_55	<i>IGFBP3</i>	rs80077929					+
region_62	<i>MYC</i>	rs4733655, rs6983267	+				
region_63	<i>CDKN2B</i>	rs7859362	+				
region_63	<i>MTAP</i>	rs7859362	+				
region_68	<i>VAV2</i>	rs7038489					+
region_73	<i>KIF20B</i>	rs140356782				+	
region_73	<i>PANK1</i>	rs140356782				+	+
region_74	<i>GOT1</i>	rs117746067		+			
region_75	<i>BORCS7</i>	rs12268849				+	
region_75	<i>AS3MT</i>	rs12268849		+		+	
region_79	<i>ANO1</i>	rs10751097					+
region_92	<i>NTN4</i>	rs11108175					+
region_93	<i>CUX2</i>	rs3858704					+
region_94	<i>TBX3</i>	rs7300312, rs11067228	+				+
region_96	<i>USP12</i>	rs116964464	+				
region_101	<i>IRS2</i>	rs1078563				+	
region_101	<i>COL4A2</i>	rs4773184					+
region_107	<i>BCL11B</i>	rs80158569				+	
region_108	<i>GOLGA8N</i>	rs56338436				+	
region_110	<i>SEN8</i>	rs8031386		+			+
region_111	<i>CIB1</i>	rs12913420		+		+	
region_111	<i>ZNF774</i>	rs7179095	+				
region_119	<i>MYL12A</i>	rs1612128	+				
region_119	<i>MYL12B</i>	rs1612128	+				
region_119	<i>TGIF1</i>	rs1612128	+				
region_125	<i>B3GNT8</i>	rs1963413				+	
region_133	<i>CEBPB</i>	rs1971480	+				
region_134	<i>RBM38</i>	rs34161672	+				
region_134	<i>BMP7</i>	rs6014965	+				+
region_138	<i>LSS</i>	rs9983528				+	+
region_138	<i>PCNT</i>	rs9983528			+	+	
region_138	<i>SPATC1L</i>	rs9983528				+	+
region_142	<i>WNT7B</i>	rs62228060					+
region_142	<i>ATXN10</i>	rs78106213				+	

The lead variant for each gene is presented by independent association signals. Supporting evidence for the likely target gene is presented as follows: “Distal”—the CCV(s) located in distal enhancer elements of the gene; “Proximal”—the CCV(s) located in proximal promoter element of the gene; “Coding”—the CCV is potential loss-of-function variants of the gene; “Colocalization (eQTL)”—target genes identified from eQTL colocalization analysis; “Colocalization (mQTL)”—target genes identified from mQTL colocalization analysis. “+” indicates the presence of supportive evidence.

**Table 4 | The 80 previously reported CRC susceptibility genes supported in this study**

Fine-mapping region	Gene	Lead variant	Distal	Proximal	Coding	Colocalization (eQTL)	Colocalization (mQTL)
region_1	<i>WNT4</i>	rs6426749				+	
region_2	<i>FHL3</i>	rs61776719				+	+
region_8	<i>LAMC1</i>	rs8179460	+	+			
region_9	<i>LMOD1</i>	rs12137232	+			+	
region_15	<i>ACTR1B</i>	rs11692435			+	+	
region_18	<i>STK39</i>	rs4668039	+			+	+
region_20	<i>TMBIM1</i>	rs3731861	+	+	+	+	
region_23	<i>SFMBT1</i>	rs2001732, rs2581817				+	+
region_26	<i>BOC</i>	rs73235124		+			
region_30	<i>TET2</i>	rs2047409, rs902443				+	+
region_31	<i>UGT8</i>	rs3924508		+			
region_35	<i>TERT</i>	rs2735940					+
region_40	<i>CDX1</i>	rs2302275		+			
region_41	<i>ERGIC1</i>	rs472959					+
region_42	<i>RREB1</i>	rs9379084			+		
region_43	<i>EDN1</i>	rs2070699					+
region_43	<i>HIVEP1</i>	rs4714081				+	+
region_47	<i>CDKN1A</i>	rs9470361					+
region_48	<i>TFEB</i>	rs6933790				+	
region_52	<i>DCBLD1</i>	rs6911915				+	
region_53	<i>TCF21</i>	rs151127921	+				
region_54	<i>GNA12</i>	rs1182197			+	+	+
region_55	<i>TBRG4</i>	rs67681615					+
region_55	<i>TNS3</i>	rs6948177	+				
region_56	<i>ABHD11</i>	rs7806956				+	+
region_57	<i>TRIM4</i>	rs2527927		+		+	
region_62	<i>POU5F1B</i>	rs6983267				+	
region_64	<i>DCAF12</i>	rs11557154			+		
region_68	<i>BRD3</i>	rs11789898				+	+
region_70	<i>BAMBI</i>	rs1773860	+				
region_71	<i>ASAH2B</i>	rs10740013				+	
region_72	<i>ZMIZ1</i>	rs704017					+
region_74	<i>ENTPD7</i>	rs35564340				+	
region_76	<i>TCF7L2</i>	rs4554812	+				
region_78	<i>TMEM258</i>	rs174570		+			
region_81	<i>TRPC6</i>	rs2186607				+	
region_82	<i>ARHGAP20</i>	rs3087967				+	
region_82	<i>FDX1</i>	rs3087967				+	
region_83	<i>BCL9L</i>	rs497916	+				
region_85	<i>PLEKHG6</i>	rs10849434, rs1003563	+	+			+
region_88	<i>CERS5</i>	rs11169572					+
region_88	<i>ATF1</i>	rs11169572		+			+
region_88	<i>DIP2B</i>	rs11169572				+	+
region_89	<i>LRP1</i>	rs7398375	+			+	+
region_91	<i>TSPAN8</i>	rs11178634			+	+	
region_98	<i>SMAD9</i>	rs12427846		+		+	+
region_99	<i>KLF5</i>	rs1304959, rs78341008	+				
region_102	<i>NIN</i>	rs1042266				+	
region_102	<i>ABHD12B</i>	rs1042266				+	+
region_102	<i>PYGL</i>	rs1042266				+	+
region_103	<i>NID2</i>	rs1151580				+	+
region_104	<i>BMP4</i>	rs1957628, rs35107139	+				+
region_105	<i>DACT1</i>	rs8020436				+	+
region_108	<i>GREM1</i>	rs16970016					+



**Table 4 (continued) | The 80 previously reported CRC susceptibility genes supported in this study**

Fine-mapping region	Gene	Lead variant	Distal	Proximal	Coding	Colocalization (eQTL)	Colocalization (mQTL)
region_109	SMAD6	rs3809570		+			+
region_109	SMAD3	rs56324967	+				
region_112	ZFP90	rs9924886				+	
region_112	CDH1	rs9924886	+	+			+
region_115	NXN	rs11247566					+
region_117	SOX9	rs112592783	+				
region_118	METRNL	rs35204860				+	+
region_120	SMAD7	rs4939821, rs2337113	+				+
region_122	FUT3	rs10409772			+	+	
region_124	RHPN2	rs28840750	+			+	
region_126	FUT2	rs12460535			+	+	
region_127	TRIM28	rs11670192				+	
region_127	ZNF584	rs8099852, rs11670192		+		+	
region_128	BMP2	rs990999				+	
region_130	MAP1LC3A	rs6059938				+	
region_130	MYH7B	rs6059938			+		
region_131	TOX2	rs6073241				+	+
region_132	PREX1	rs6066825					+
region_133	PARD6B	rs6091213				+	
region_133	PTPN1	rs6091213	+				
region_135	GNAS	rs8121252		+			
region_136	RBBP8NL	rs1741640				+	
region_137	STMN3	rs6089763					+
region_139	ZNRF3	rs4616575	+			+	
region_140	PDGFB	rs130651					+
region_142	RIBC2	rs6007600				+	+

The lead variant for each gene is presented by independent association signals. Supporting evidence for the likely target gene is presented as follows: “Distal”—the CCV(s) located in distal enhancer elements of the gene; “Proximal”—the CCV(s) located in proximal promoter element of the gene; “Coding”—the CCV is potential loss-of-function variants of the gene; “Colocalization (eQTL)”—target genes identified from eQTL colocalization analysis; “Colocalization (mQTL)”—target genes identified from mQTL colocalization analysis. “+” indicates the presence of supportive evidence.

single-ancestry analysis. Our analysis is limited to two ancestry groups. Further studies should increase the diversity of genetic data, including those from other racial groups.

In summary, our large trans-ancestry fine-mapping analysis has identified large numbers of not previously reported independent association signals for CRC risk and refined the majority of the previously reported association signals. By leveraging data from two ancestries, we further defined putative causal variants underlying CRC risk signals. Our study has also uncovered a credible set of target genes. These findings offer a significant advancement in our understanding of the genetic and biological processes underlying CRC and provide a roadmap for further investigation of variants and genes identified in our study.

## Methods

### GWAS data and meta-analysis

The GWAS data used in this study comprised 100,204 CRC cases and 154,587 controls (Supplementary Data 1), which were grouped into 31 GWAS analytical units based on the study or genotyping platform as consistent with the original reports. Of them, 17 datasets were derived from populations of European descent and 14 were from populations of Asian descent. These 31 GWAS datasets were meta-analyzed under the fixed-effects inverse variance weighted model implemented in METAL<sup>30</sup>. Further details regarding each analytical unit and meta-analysis were described in Supplementary Note.

### Identifying independent association signals

A total of 205 independent genetic associations have been reported for CRC risk by GWAS<sup>7</sup>. To define fine-mapping regions for CRC, we

aggregated these risk variants using *bedtools*. Specifically, we identified 1 megabase (Mb) intervals centered on the risk variants, and if there were regions of overlap, we combined them into a single interval over 1 Mb. In total, we determined 143 fine-mapping regions, including 142 on autosomes and one on chromosome X (Supplementary Data 2). Our fine-mapping analysis and downstream analyses focused on the 142 genomic risk regions on autosomes.

To identify distinct association signals within each risk region, we conducted a forward stepwise conditional analysis for summary statistics from the trans-ancestral meta-analysis, using GCTA-COJO<sup>31,32</sup>. We included common variants (MAF > 0.01) with associations at  $P < 0.05$  in both populations. To account for differences in the LD structure, we conducted conditional analysis in each population for each fine-mapping region, conditioning on the most significant association from the trans-ancestral summary statistics. We then meta-analyzed the conditioned results using the fixed-effects inverse variance weighted model with METAL. To identify potential ancestry-specific independent signals, we also performed conditional analysis in each population, conditioning on the ancestry-specific most significant association. Common variants (MAF > 0.01) with association at  $P < 1 \times 10^{-4}$  in each population were included. For LD estimation, we used genotyping data from 6684 unrelated samples of Asian descent<sup>33</sup>, and 503 European samples in the 1000 Genome project as the reference.

Following a previous study conducted for breast cancer<sup>12</sup>, we applied the conditional  $P$  value  $< 1 \times 10^{-6}$  to define the independent signal. For each region, we first adjusted for the most significant association and then added any additional variant that remained an independent signal at the conditional  $P$  value  $< 1 \times 10^{-6}$  to the

**Table 5 | Significant enrichment in biological pathways**

Pathways <sup>a</sup>	Genes <sup>b</sup>
TGF-beta signaling	<i>BAMBI</i> , <i>BMP2</i> , <i>BMP4</i> , <b><i>BMP7</i></b> , <i>CDH1</i> , <b><i>CDKN2B</i></b> , <i>GREM1</i> , <b><i>MYC</i></b> , <b><i>RUNX2</i></b> , <i>SMAD3</i> , <i>SMAD6</i> , <i>SMAD7</i> , <i>SMAD9</i> , <b><i>TGIF1</i></b>
Hippo signaling	<i>BMP2</i> , <i>BMP4</i> , <b><i>BMP7</i></b> , <i>CDH1</i> , <i>GNAS</i> , <b><i>MYC</i></b> , <i>PARD6B</i> , <i>SMAD3</i> , <i>SMAD7</i> , <i>TCF7L2</i> , <i>WNT4</i> , <b><i>WNT7B</i></b>
TNF-alpha Signaling via NF-kB	<b><i>TGIF1</i></b> , <i>BMP2</i> , <i>CDKN1A</i> , <i>EDN1</i> , <b><i>CEBPB</i></b> , <i>SMAD3</i> , <b><i>MYC</i></b> , <b><i>IRS2</i></b>
BMP signaling	<i>BMP2</i> , <i>SMAD6</i> , <b><i>RUNX2</i></b> , <i>SMAD9</i> , <i>SMAD7</i>
Pluripotency of stem cells	<i>POU5F1B</i> , <i>BMP4</i> , <i>SMAD3</i> , <b><i>MYC</i></b> , <b><i>WNT7B</i></b> , <i>SMAD9</i> , <b><i>TBX3</i></b> , <i>WNT4</i> , <i>PDGFB</i> , <i>SMAD6</i> , <i>SMAD7</i> , <i>TCF7L2</i>
Epithelial-mesenchymal transition	<i>SMAD3</i> , <i>CDH1</i> , <b><i>RUNX2</i></b> , <i>GREM1</i> , <b><i>COL4A2</i></b> , <i>LRP1</i> , <b><i>IGFBP3</i></b> , <i>LAMC1</i> , <i>NID2</i> , <b><i>WNT7B</i></b> , <i>WNT4</i>
Extracellular matrix organization	<i>BMP4</i> , <i>BMP2</i> , <b><i>COL4A2</i></b> , <i>PDGFB</i> , <b><i>NTN4</i></b> , <i>LAMC1</i> , <b><i>HSPG2</i></b> , <i>NID2</i> , <b><i>BMP7</i></b> , <b><i>ADAMTSS9</i></b>
Senescence and Autophagy	<i>BMP2</i> , <i>CDKN1A</i> , <b><i>CEBPB</i></b> , <i>SMAD3</i> , <i>MAP1LC3A</i> , <b><i>IGFBP3</i></b> , <b><i>CDKN2B</i></b> , <b><i>MYC</i></b> , <i>KLF5</i>
DNA damage response	<i>TCF7L2</i> , <i>CDKN1A</i> , <i>SMAD3</i> , <b><i>MYC</i></b> , <b><i>WNT7B</i></b> , <i>WNT4</i>
Cell cycle	<i>CDKN1A</i> , <b><i>CDKN2B</i></b> , <i>SMAD3</i> , <b><i>MYC</i></b>
Focal adhesion	<b><i>COL4A2</i></b> , <i>PDGFB</i> , <i>LAMC1</i> , <b><i>MYL12A</i></b> , <b><i>MYL12B</i></b> , <b><i>VAV2</i></b>
Adherens junction	<i>PTPN1</i> , <i>TCF7L2</i> , <i>SMAD3</i> , <i>CDH1</i>
Glycolysis	<b><i>GOT1</i></b> , <b><i>IGFBP3</i></b> , <b><i>IRS2</i></b> , <i>SOX9</i> , <i>PYGL</i> , <b><i>LCT</i></b>
Proteoglycans in cancer	<i>CDKN1A</i> , <b><i>MYC</i></b> , <b><i>WNT7B</i></b> , <b><i>HSPG2</i></b> , <i>WNT4</i> , <b><i>VAV2</i></b>
Androgen Response	<b><i>HPGD</i></b> , <i>ZMIZ1</i> , <i>STK39</i> , <b><i>MYL12A</i></b>
Sphingolipid Metabolism	<i>UGT8</i> , <i>CERS5</i>
Other cancer related pathways <sup>c</sup>	<i>TCF7L2</i> , <i>CDKN1A</i> , <i>EDN1</i> , <b><i>CDKN2B</i></b> , <i>SMAD3</i> , <b><i>WNT7B</i></b> , <i>PTGER3</i> , <i>PDGFB</i> , <i>LAMC1</i> , <b><i>MLH1</i></b> , <i>BMP4</i> , <i>BMP2</i> , <b><i>COL4A2</i></b> , <i>TERT</i> , <i>CDH1</i> , <b><i>MYC</i></b> , <i>GNAI2</i> , <i>GNAS</i> , <i>WNT4</i> , <i>ATF1</i> , <b><i>CEBPB</i></b> , <b><i>BCL11B</i></b> , <b><i>HPGD</i></b> , <b><i>IGFBP3</i></b> , <b><i>RUNX2</i></b> , <i>ZMIZ1</i> , <i>SMAD6</i> , <i>SMAD7</i> , <b><i>VAV2</i></b> , <i>TFEB</i>

<sup>a</sup>Genes from the same or similar pathway item in multiple databases were combined.

<sup>b</sup>Genes identified in this study for each pathway item are highlighted in bold.

<sup>c</sup>Genes from all general cancer-related pathways (i.e., pathway in cancer, colorectal cancer) identified in multiple databases were combined.

conditional set. We then repeated the conditional analysis until no more variants met the significance threshold. In regions with multiple independent signals, we determined the index variant for each signal through a process of conditional analysis, adjusting for the index variants of the other signals. This process was repeated until the set of index variants were stabilized. The variant with the strongest residual association was defined as the index for the signal.

For independent association signals identified in ancestry-specific analyses, we compared them with those from trans-ancestry analyses by assessing correlations between their lead variants within each risk region. If a signal was consistently found in both ancestry-specific and trans-ancestry analyses (i.e., the same lead variant or correlated lead variants with LD  $r^2 > 0.1$  in each corresponding population), we considered it as a sharing signal between Asian and European-ancestry populations. Otherwise, they were defined as ancestry-specific signals.

### Identifying a set of CCVs of each independent signal

To determine the CCVs of each independent signal, we used the approach described in a previous study for breast cancer<sup>12</sup>. Specifically, variants that have a conditional *P* value within two orders of magnitude of the most significant association, conditioning on all other independent association signals, were defined as CCVs.

### RNA-seq data analysis

We conducted mRNA sequencing on tumor-adjacent normal colon tissues obtained from 364 East Asians patients with colorectal cancer who participated in the ACCC. Furthermore, we included RNA-seq data from normal colon tissues from 423 individuals of European ancestry who participated in the BarcUva-Seq project. Included subjects, library preparation and sequencing of colon tissue samples in the ACCC and the BarcUva-Seq project have been presented in Supplementary Note.

The raw RNA-seq data were processed according to the pipeline of the GTEx Consortium. Sequencing reads were aligned to the reference genome GRCh37 (RNA-seq data from East Asians) or GRCh38 (RNA-seq data from the BarcUva-Seq project) with STAR (v2.5.4)<sup>34</sup>. Quality control of aligned samples was performed using RNA-SeQC (v2.3.5)<sup>35</sup>. Samples that met any of the following criteria were removed: (1) <10 million mapped reads; (2) read mapping rate < 0.2; (3) intergenic mapping rate

>0.4; (4) base mismatch rate >0.01 for read mate 1 or >0.02 for read mate 2; and (5) rRNA mapping rate >0.3. If the sample had replicated RNA-seq data, the one with the highest mapped reads was retained.

Gene-level expression quantification was performed using RNA-SeQC based on the GENCODE release 19 annotation (for RNA-seq data from East Asians) and the GENCODE release 26 annotation (for RNA-seq data from the BarcUva-Seq project)<sup>36</sup>. The read counts and TPM values of genes were calculated using aligned reads with the following criteria: (1) reads were uniquely mapped; (2) aligned reads were properly paired; (3) the read alignment distance was <6. The genes with expression thresholds of  $\geq 0.1$  TPM in  $\geq 20\%$  of samples and  $\geq 6$  reads (unnormalized) in  $\geq 20\%$  of samples were selected. Quantile normalization of the gene expression was performed. We further performed rank-based inverse normal transformation for the expressions of genes across samples.

### Cis-expression/methylation quantitative loci (cis-eQTL/mQTL) analysis

To identify target genes, we performed cis-eQTL analysis based on a linear regression framework<sup>10,11</sup>. Gene expression data from four expression datasets comprising a total of 1,299 individuals were used: 1) GTEx project of transverse colon tissues from 368 individuals predominantly of European ancestry, 2) Colonomics project of normal colon tissues or tumor-adjacent normal colon tissues from 144 individuals of European ancestry, 3) BarcUva-Seq project of normal colon tissues from 423 individuals of European ancestry, and 4) ACCC of tumor-adjacent normal colon tissue from 364 CRC patients of East Asian ancestry. We obtained available cis-eQTL results for CCVs and their nearby genes (within 1 Mb to CCV) from the GTEx database (version 8) and the Colonomics project. Details for gene expression data and eQTL analysis in the Colonomics project are explained elsewhere<sup>37</sup>. For the analyses using the remaining two datasets, we conducted a linear regression analysis to assess the associations between CCV and the normalized expression levels of nearby genes (within 1 Mb to CCV), adjusting for age, gender, and five top principal components.

We conducted cis-mQTL analysis for CCVs identified in European and trans-ancestry analyses. To do this, we included methylation data

obtained from a total of 321 individuals. These datasets consisted of 189 transverse colon tissues predominantly of European ancestry from GTEx, as well as normal colon tissues or tumor-adjacent normal colon tissues of 132 individuals of European ancestry from the Colonomics project. We extracted cis-mQTL results for CCVs and their nearby CpG sites (within 1 Mb to CCV) from the GTEx database (version 8)<sup>14</sup>. In the Colonomics project, a linear regression analysis was used to evaluate the associations between CCV and the normalized methylation levels of CpG sites (within 1 Mb to CCV), with adjustments of age, gender, and colon sites (right/left). Further details about the cis-mQTL in the Colonomics project can be found in previous studies<sup>37,38</sup>.

### Meta-analysis of cis-eQTL/mQTL results

We performed a meta-analysis to integrate the summary cis-eQTL/mQTL results based on beta and *p* values from different datasets<sup>10,11</sup>. In brief, we calculated the *z* score from function  $qnorm(p/2)*sign(beta)$  and further converted the standard *z* score derived from  $sum(z*sqrt(N))/sqrt(sum(N))$  with a normalized weighted sampled size. Here, beta and *p* value were derived from eQTL/mQTL results and *N* referred to the sample size for each dataset. The meta *p* value was derived from the standard *z* score. For independent signals detected in both European and Asian populations, the eQTL results from both populations were combined.

We adjusted the combined *p*-values of eQTL/mQTL results with the Bonferroni procedure. The procedure was conducted for index variants of independent association signals. The Bonferroni-adjusted  $P < 0.05$  was applied to identify potential target genes for each signal.

### Colocalization analyses between GWAS association signals and eQTL/mQTL signals

To identify putative target genes, we employed the SMR method to conduct a colocalization analysis<sup>39</sup>. We integrated GWAS summary statistics of CCVs and their associations with genes from eQTL/mQTL analysis described above. The results of meta-analyses on cis-eQTLs/mQTLs were used. Specifically, we have a statistic:

$$T_{SMR} = b_{xy}^2 / Var(b_{xy}) \approx \frac{Z_{zy}^2 Z_{zx}^2}{Z_{zy}^2 + Z_{zx}^2} \quad (1)$$

Here,  $Z_{zx}$  and  $Z_{zy}$  are the *Z* statistics for the GWAS summary statistics and the cis-eQTL/mQTL results, respectively.  $T_{SMR}$  is the  $\chi^2$  statistic, which tests the significance of  $b_{xy}$ . The significant colocalized signals were determined based on the threshold of the Bonferroni-corrected  $P_{SMR} < 0.05$  within each independent signal.

### Functional annotation of CCVs

We investigated whether each potential causal variant was mapped to gene regulatory regions (i.e., promoter or enhancer) (Supplementary Data 8). We obtained 351 chromatin immunoprecipitation sequencing (ChIP-seq) peak files for histone modification marks and transcription factors, and 25 DNase I hypersensitive sites sequencing peak files for chromatin accessibility, generated in normal colorectal epithelium and CRC cell lines from the Cistrome database<sup>40,41</sup>. Only peaks that met all six quality controls set recommended by Cistrome were analyzed. Additionally, we obtained available ChIP-seq data of histone modification marks from colon tissues, tumor tissues of CRC, and CRC cell lines from Gene Expression Omnibus (GEO), which included 16 from GSE133928<sup>42</sup>, 215 from GSE136889<sup>43</sup> and 233 from GSE156613<sup>44</sup>. To generate coverage tracks Bigwig (bw) files for ChIP-seq data, we converted them to bedGraph files and then identified peaks with the subcommand *bdgpeakcall* from *macs2*<sup>45</sup>. For each variant, we examined whether it was mapped to a peak region of histone modification marks, DNase I hypersensitive, or transcription factors binding sites using an *in-house* script.

### In silico prediction of regulatory element-to-gene

Since the majority of the CCVs are located outside protein-coding regions, genes can potentially be regulated by CCVs located in distal enhancer elements and proximal promoter elements. Hence, we identified an extensive set of functional genomic data from normal colon tissues or tumor tissues of colorectal cancer or colorectal cancer cell lines (Supplementary Data 9). Subsequently, we conducted an in-silico analysis for each CCV-gene pair.

We used a variety of experimental and computational functional genomic data to identify target genes of CCVs in regulatory elements. Specifically, for distal regulatory elements, we utilized chromatin-chromatin interaction data from experiments or computational predictions. To do this, we downloaded 13 experimental chromatin-chromatin interaction datasets under accessions GSE133928<sup>42</sup> and GSE136629<sup>43</sup> from GEO, as well as two promoter capture Hi-C datasets from the previous study<sup>46</sup>. We combined this data with ChIP-seq data of the histone modification H3K27ac (an active enhancer mark) to identify enhancer-promoter loops. We defined these loops as interactions where one fragment overlapped an H3K27ac peak (enhancer-like) and the other fragment overlapped the promoter of a gene (the region from downstream 1 kb to upstream 100 bp around the transcription start site).

In addition to this, we downloaded experimentally confirmed enhancer-gene pairs from the ENdb database. We also obtained computational enhancer-promoter interactions from IM-PET<sup>47</sup>, FANTOM5<sup>48,49</sup>, EnhancerAtlas<sup>50</sup>, and super-enhancer<sup>51,52</sup>. To further refine our analysis, we included topologically associating domain (TAD) boundaries in three colorectal cancer cell lines (HT29, LoVo, and DLD1)<sup>46,53</sup>. Finally, we examined the overlap between CCVs and enhancer elements. For proximal promoter elements, we analyzed CCVs located within gene promoter regions that intersected with ChIP-seq peaks of H3K4me3 (an activity promoter mark).

To identify potential loss-of-function variants and their corresponding targeted genes, we conducted variant annotation of CCVs using the Variant Effect Predictor (VEP) tool<sup>54</sup>. To predict the consequence of missense coding variants, we utilized PolyPhen-2 and SIFT. Furthermore, to evaluate splicing effects, MaxEntScan was used.

We scored CCV target genes using different criteria (Supplementary Data 9). For the potential target gene of CCV in distal enhancer elements, the gene was awarded two points or one point if there was evidence from experimental chromatin-chromatin interaction or computed interaction. The score was unweighted to three if both experimental and computational interaction were detected for the gene-CCV pair. If CCV interacted with genomic features (open chromatin, activity enhancer, and TF binding sites), the corresponding gene was further unweighted by one point. An additional point was awarded if there are at least two interactions for the CCV. If the gene were colorectal cancer or pan-cancer drivers<sup>55</sup>, they were up-weighted by an additional point. The score was down-weighted for the gene if the CCV-gene pair was separated by TAD or a lack of expression in colon tissues. Distal scores eventually ranged from 0 to 6. For the potential target gene of CCV in proximal promoter elements, the gene was awarded one point if CCV overlapped with binding sites of transcription factors. If genes were colorectal cancer or pan-cancer drivers, they were up-weighted with an additional point. A lack of its expression resulted in down-weighting to 10% as target genes. Proximal scores eventually ranged from 0 to 2. Genes predicted to be regulated targets of coding CCVs were awarded points based on the annotation as either of missense, nonsense, and predicted splicing alterations. The consequences of missense variants which probably are damaging or deleterious resulted in the addition of one point to the target gene. Further points were awarded to such a gene if it was colorectal cancer or pan-cancer drivers. A lack of expression reduced the score (the score was down to 10%). Coding scores ranged from 0 to 2. For the set of confident target genes, we defined such genes if it has a distal score >4 or a proximal score >1, or a coding score >1.

### Credible set of susceptibility genes

To determine a set of credible genes for CRC susceptibility, we combined information on gene-CRC risk associations through TWAS and colocalization of eQTL signal with GWAS risk signals for genes that were present in both our study and previous investigations. We used three sets of previously identified genes below: (A) 155 effector genes identified through GWAS, TWAS, TisWAS, and MWAS<sup>7</sup>; (B) 136, 26, and 48 genes identified through TWASs<sup>7,16,17</sup>; (C) 73 genes identified through colocalization analysis between eQTL and GWAS signals<sup>11</sup> or genes associated with CRC risk at nominal  $P < 0.05$  in the previous TWAS<sup>17</sup>. We considered the prioritization order as  $A > B > C$  for these three gene sets and focused on protein-coding genes outside the MHC region. For the independent association signals with multiple target gene candidates, we kept either genes with higher prioritization or all genes if there was no evidence from these three gene sets. For the independent association signals with a single gene, we kept it regardless of evidence from the gene sets.

### Single-cell RNA-sequencing data analysis

We included single-cell RNA-sequencing datasets from colon tissues of 31 individuals who participated in the Colorectal Molecular Atlas Project (COLON MAP)<sup>18</sup>. We analyzed gene expression dataset for each individual's cell and combined these datasets into a count matrix. We normalized the number of unique molecular identifiers (UMIs) per cell and converted it to transcripts per 10,000 transcripts (TP10K). Next, we applied a logarithmic transformation to the normalized values and got the  $\log_2(\text{TP10K} + 1)$  expression matrix for the downstream analyses. Further, we determined the 2000 most highly variable genes within the entire dataset and performed a principal component analysis (PCA). The top 30 and 40 principal components (PCs) were identified. Subsequently, we performed batch correction removal and utilized the top 40 batch-corrected components to construct a  $k$ -nearest neighbors graph of cell profiles with  $k = 9$ . We visualized the individual single-cell profiles using the Uniform Manifold Approximation and Projection (UMAP) and constructed the neighborhood graph using the Leiden graph-clustering method. Nine cell types were defined, including well-known major cell types such as absorptive cells (ABS), crypt top colonocytes (CT), enteroendocrine cells (EE), goblet cells (GOB), stem cells (STM), and others. We identified differentially expressed genes (DEGs) by comparing each cell type with all other cell types and calculated a  $P$ -value for each gene using Wilcoxon's rank-sum test. The criteria  $|\log_2 \text{fold change (FC)}| > 1$  and  $P < 0.05$  were applied to determine genes with significantly differential expression between cell types.

### Burden test for credible susceptibility genes

We annotated all variants in the UKBB WES 200 K cohort with functional annotations from ANNOVAR<sup>56</sup> based on the reference genome GRCh38. We only included rare loss-of-function (LoF) and deleterious missense (Dmis) variants with  $\text{MAF} < 0.01$  in our gene-based test. LoF variants were those predicted as frameshift insertion/deletion, splice-site alteration, stop gain, and stop loss by ANNOVAR, and deleterious missense (Dmis) variants were those predicted as deleterious by MetaSVM<sup>57</sup>. We considered both LoF sets and damaging sets (LoF+ Dmis) within a gene for testing. For a given set, we collapsed rare variants within a gene as a single combined 'mask' and tested the association between the 'mask' genotype and the CRC phenotype using logistic regression after adjusting for sex, age, the interaction of sex and age, and the top four principal components.

### Pathway analysis of credible susceptibility genes

To explore the potential biological roles of the identified CRC susceptibility genes, we analyzed their functional enrichment using the enrichR<sup>19-21</sup> and various pathway databases, including WikiPathway, KEGG, MSigDB, and Reactome. The biological pathways (adjusted  $P < 0.05$ ) were considered and presented.

### Reporting summary

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

### Data availability

The GWAS summary statistics are available at the GWAS catalog under accession number [GCST90129505](https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs000424.v8.p2). The RNA-seq data and genotype data of subjects of East Asian ancestry from the ACCC is being deposited to NCBI database of Genotypes and Phenotypes (dbGaP), accession number phs002813.v1.p1. All requests to access these data could also be made by contacting Drs. Wei Zheng ([wei.zheng@vanderbilt.edu](mailto:wei.zheng@vanderbilt.edu)) and Xingyi Guo ([xingyi.guo@vumc.org](mailto:xingyi.guo@vumc.org)). The data from the Genotype-Tissue Expression (GTEx, version 8) project used in this study are publicly available at the dbGaP under accession number phs000424.v8.p2 ([https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study\\_id=phs000424.v8.p2](https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs000424.v8.p2)). The transcriptome and genotype data as well as the sample covariates from the BarcUVA-Seq project can be accessed at the dbGaP under accession number phs003338.v1.p1 ([https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study\\_id=phs003338.v1.p1](https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs003338.v1.p1)). The access to data from the Colonomics project could be requested by submission of an inquiry to Dr. Victor Moreno ([v.moreno@iconcologia.net](mailto:v.moreno@iconcologia.net)). The CRC-relevant epigenome and functional genomic data were obtained from the NCBI's Gene Expression Omnibus database (GEO) under accession numbers: [GSE133928](https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE133928), [GSE136889](https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE136889), and [GSE156613](https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE156613). Enhancer-promoter interaction data were obtained from the ENdb database (<https://bio.liclab.net/ENdb/>), 4Dgenome (<https://4dgenome.research.chop.edu/>), FANTOM5 (<https://fantom.gsc.riken.jp/5/>), EnhancerAtlas 2.0 (<http://www.enhanceratlas.org/>) and Super-enhancers (<https://bio.liclab.net/sedb/> and [https://www.cell.com/fulltext/S0092-8674\(13\)01227-0#supplementaryMaterial](https://www.cell.com/fulltext/S0092-8674(13)01227-0#supplementaryMaterial)). Single-cell RNA-sequencing datasets from colon tissues of 31 individuals were obtained from the Colorectal Molecular Atlas Project (COLON MAP). Whole exome sequencing data from 137,104 individuals of European ancestry were obtained from the UK Biobank (<https://www.ukbiobank.ac.uk/>).

### Code availability

The code used in this study is available at the GitHub repository [https://github.com/zhishanchen/CRC\\_Finmapping](https://github.com/zhishanchen/CRC_Finmapping)<sup>58</sup>.

### References

- Sung, H. et al. Global cancer statistics 2020: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries. *CA Cancer J. Clin.* **71**, 209–249 (2021).
- Jiao, S. et al. Estimating the heritability of colorectal cancer. *Hum. Mol. Genet.* **23**, 3898–3905 (2014).
- Lu, Y. et al. Large-scale genome-wide association study of east asians identifies loci associated with risk for colorectal cancer. *Gastroenterology* **156**, 1455–1466 (2019).
- Huyghe, J. R. et al. Discovery of common and rare genetic risk variants for colorectal cancer. *Nat. Genet.* **51**, 76–87 (2019).
- Law, P. J. et al. Association analyses identify 31 new risk loci for colorectal cancer susceptibility. *Nat. Commun.* **10**, 2154 (2019).
- Lu, Y. et al. Identification of Novel Loci and New Risk Variant in Known Loci for Colorectal Cancer Risk in East Asians. *Cancer Epidemiol. Biomark. Prev.* **29**, 477–486 (2020).
- Fernandez-Rozadilla, C. et al. Deciphering colorectal cancer genetics through multi-omic analysis of 100,204 cases and 154,587 controls of European and east Asian ancestries. *Nat. Genet.* **55**, 89–99 (2023).
- Zeng, C. et al. Identification of independent association signals and putative functional variants for breast cancer risk through fine-scale mapping of the 12p11 locus. *Breast Cancer Res.* **18** (2016).
- Guo, X. et al. A comprehensive cis-eQTL analysis revealed target genes in breast cancer susceptibility loci identified in genome-wide association studies. *Am. J. Hum. Genet.* **102**, 890–903 (2018).

10. Chen, Z. et al. Identifying putative susceptibility genes and evaluating their associations with somatic mutations in human cancers. *Am. J. Hum. Genet.* **105**, 477–492 (2019).
11. Yuan, Y. et al. Multi-omics analysis to identify susceptibility genes for colorectal cancer. *Hum. Mol. Genet.* **30**, 321–330 (2021).
12. Fachal, L. et al. Fine-mapping of 150 breast cancer risk regions identifies 191 likely target genes. *Nat. Genet.* **52**, 56–73 (2020).
13. Thiagalingam, A. et al. RREB-1, a novel zinc finger protein, is involved in the differentiation response to Ras in human medullary thyroid carcinomas. *Mol. Cell. Biol.* **16**, 5335–5345 (1996).
14. Oliva, M. et al. DNA methylation QTL mapping across diverse human tissues provides molecular links between genetic variation and complex traits. *Nat. Genet.* **55**, 112–122 (2023).
15. Wu, Y. et al. Joint analysis of GWAS and multi-omics QTL summary statistics reveals a large fraction of GWAS signals shared with molecular phenotypes. *Cell Genom.* **3**, 100344 (2023).
16. Guo, X. et al. Identifying novel susceptibility genes for colorectal cancer risk from a transcriptome-wide association study of 125,478 subjects. *Gastroenterology* **160**, 1164–1178.e6 (2021).
17. Chen, Z. et al. Novel insights into genetic susceptibility for colorectal cancer from transcriptome-wide association and functional investigation. *J. Natl. Cancer Inst.* <https://doi.org/10.1093/jnci/djad178> (2023).
18. Chen, B. et al. Differential pre-malignant programs and micro-environment chart distinct paths to malignancy in human colorectal polyps. *Cell* **184**, 6262–6280.e26 (2021).
19. Chen, E. Y. et al. Enrichr: interactive and collaborative HTML5 gene list enrichment analysis tool. *BMC Bioinform.* **14**, 128 (2013).
20. Kuleshov, M. V. et al. Enrichr: a comprehensive gene set enrichment analysis web server 2016 update. *Nucleic Acids Res.* **44**, W90–W97 (2016).
21. Xie, Z. et al. Gene set knowledge discovery with Enrichr. *Curr. Protoc.* **1**, e90 (2021).
22. Jung, B., Staudacher, J. J. & Beauchamp, D. Transforming growth factor  $\beta$  superfamily signaling in development of colorectal cancer. *Gastroenterology* **152**, 36–52 (2017).
23. Feng, Q., Li, S., Ma, H.-M., Yang, W.-T. & Zheng, P.-S. LGR6 activates the Wnt/ $\beta$ -catenin signaling pathway and forms a  $\beta$ -catenin/TCF7L2/LGR6 feedback loop in LGR6high cervical cancer stem cells. *Oncogene* **40**, 6103–6114 (2021).
24. Ruan, X. et al. Silencing LGR6 attenuates stemness and chemoresistance via inhibiting Wnt/ $\beta$ -catenin signaling in ovarian cancer. *Mol. Ther. Oncolytics* **14**, 94–106 (2019).
25. Dong, L., Lyu, X., Faleti, O. D. & He, M.-L. The special stemness functions of Tbx3 in stem cells and cancer development. *Semin. Cancer Biol.* **57**, 105–110 (2019).
26. Russell, R. et al. A dynamic role of TBX3 in the pluripotency circuitry. *Stem Cell Rep.* **5**, 1155–1170 (2015).
27. Elbadawy, M., Usui, T., Yamawaki, H. & Sasaki, K. Emerging roles of C-Myc in cancer stem cell-related signaling and resistance to cancer chemotherapy: A potential therapeutic target against colorectal cancer. *Int. J. Mol. Sci.* **20**, 2340 (2019).
28. Satoh, K. et al. Global metabolic reprogramming of colorectal cancer occurs at adenoma stage and is induced by MYC. *Proc. Natl. Acad. Sci. USA* **114**, E7697–E7706 (2017).
29. Ong, E. S. et al. Metabolic profiling in colorectal cancer reveals signature metabolic shifts during tumorigenesis. *Mol. Cell. Proteom.* <https://doi.org/10.1074/mcp.M900551-MCP200> (2010).
30. Willer, C. J., Li, Y. & Abecasis, G. R. METAL: fast and efficient meta-analysis of genomewide association scans. *Bioinformatics* **26**, 2190–2191 (2010).
31. Yang, J. et al. Conditional and joint multiple-SNP analysis of GWAS summary statistics identifies additional variants influencing complex traits. *Nat. Genet.* **44**, 369–375 (2012).
32. Yang, J., Lee, S. H., Goddard, M. E. & Visscher, P. M. GCTA: a tool for genome-wide complex trait analysis. *Am. J. Hum. Genet.* **88**, 76–82 (2011).
33. 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).
34. Dobin, A. et al. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15–21 (2013).
35. Graubert, A., Aguet, F., Ravi, A., Ardlie, K. G. & Getz, G. RNA-SeqQC 2: Efficient RNA-seq quality control and quantification for large cohorts. *Bioinformatics* **37**, 3048–3050 (2021).
36. Frankish, A. et al. GENCODE: reference annotation for the human and mouse genomes in 2023. *Nucleic Acids Res.* **51**, D942–D949 (2023).
37. Díez-Villanueva, A. et al. Identifying causal models between genetically regulated methylation patterns and gene expression in healthy colon tissue. *Clin. Epigenetics* **13**, 162 (2021).
38. Díez-Villanueva, A. et al. DNA methylation events in transcription factors and gene expression changes in colon cancer. *Epigenomics* **12**, 1593–1610 (2020).
39. Zhu, Z. et al. Integration of summary data from GWAS and eQTL studies predicts complex trait gene targets. *Nat. Genet.* **48**, 481–487 (2016).
40. Zheng, R. et al. Cistrome Data Browser: expanded datasets and new tools for gene regulatory analysis. *Nucleic Acids Res.* **47**, D729–D735 (2019).
41. Mei, S. et al. Cistrome Data Browser: a data portal for ChIP-Seq and chromatin accessibility data in human and mouse. *Nucleic Acids Res.* **45**, D658–D662 (2017).
42. Johnstone, S. E. et al. Large-scale topological changes restrain malignant progression in colorectal cancer. *Cell* **182**, 1474–1489.e23 (2020).
43. Orouji, E. et al. Chromatin state dynamics confers specific therapeutic strategies in enhancer subtypes of colorectal cancer. *Gut* **71**, 938–949 (2022).
44. Li, Q.-L. et al. Genome-wide profiling in colorectal cancer identifies PHF19 and TBC1D16 as oncogenic super enhancers. *Nat. Commun.* **12**, 6407 (2021).
45. Zhang, Y. et al. Model-based analysis of ChIP-Seq (MACS). *Genome Biol.* **9**, R137 (2008).
46. Orlando, G. et al. Promoter capture Hi-C-based identification of recurrent noncoding mutations in colorectal cancer. *Nat. Genet.* **50**, 1375–1380 (2018).
47. He, B., Chen, C., Teng, L. & Tan, K. Global view of enhancer-promoter interactome in human cells. *Proc. Natl. Acad. Sci. USA* **111**, E2191–E2199 (2014).
48. Lizio, M. et al. Gateways to the FANTOM5 promoter level mammalian expression atlas. *Genome Biol.* **16**, 22 (2015).
49. Andersson, R. et al. An atlas of active enhancers across human cell types and tissues. *Nature* **507**, 455–461 (2014).
50. Gao, T. & Qian, J. EnhancerAtlas 2.0: an updated resource with enhancer annotation in 586 tissue/cell types across nine species. *Nucleic Acids Res.* **48**, D58–D64 (2020).
51. Hnisz, D. et al. Super-enhancers in the control of cell identity and disease. *Cell* **155**, 934–947 (2013).
52. Jiang, Y. et al. SEdb: a comprehensive human super-enhancer database. *Nucleic Acids Res.* **47**, D235–D243 (2019).
53. Lajoie, B. R., Dekker, J. & Kaplan, N. The Hitchhiker’s guide to Hi-C analysis: practical guidelines. *Methods* **72**, 65–75 (2015).
54. McLaren, W. et al. The ensembl variant effect predictor. *Genome Biol.* **17** (2016).
55. Bailey, M. H. et al. Comprehensive characterization of cancer driver genes and mutations. *Cell* **173**, 371–385.e18 (2018).
56. Yang, H. & Wang, K. Genomic variant annotation and prioritization with ANNOVAR and wANNOVAR. *Nat. Protoc.* **10**, 1556–1566 (2015).

57. Kim, S., Jhong, J.-H., Lee, J. & Koo, J.-Y. Meta-analytic support vector machine for integrating multiple omics data. *BioData Min.* **10**, 2 (2017).
58. zhishanchen. *zhishanchen/CRC\_Finemapping: crc\_finemapping*. (Zenodo, 2024). <https://doi.org/10.5281/ZENODO.10645372>

## Acknowledgements

This research was supported primarily by US National Institutes of Health (NIH) grant R01CA188214 (to W.Z.), Anne Potter Wilson Chair endowment from the Vanderbilt University School of Medicine (to W.Z.), and NIH grant R37CA227130 and R01CA269589 (to X.G.). Sample preparation and genotyping assays at Vanderbilt University were conducted at the Survey and Biospecimen Shared Resources and Vanderbilt Microarray Shared Resource, supported in part by the Vanderbilt-Ingram Cancer Center (grant P30CA068485). Data analyses were performed on servers maintained by the Advanced Computing Center for Research and Education (ACCRE) at Vanderbilt University.

## Author contributions

Wei Zheng and Xingyi Guo conceived and supervised the study, and acquired funding. Wei Zheng, Xingyi Guo, and Zhishan Chen designed the study with significant contribution from Ran Tao. Zhishan Chen carried out the main analysis. Chao Li, Quanhu Shen, and Ken S Lau contributed to single-cell RNA-seq analysis. Yuhan Xie and Hongyu Zhao contributed to whole exome sequencing analysis. Zhishan Chen, Xingyi Guo, and Wei Zheng interpreted results with help from other authors. Zhishan Chen, Xingyi Guo, Jeroen R Huyghe, Philip J Law, Ceres Fernandez-Rozadilla, Jie Ping, Guochong Jia, Maria N Timofeeva, Minta Thomas, Stephanie L Schmit, Virginia Díez-Obrero, Matthew Devall, Ferran Moratalla-Navarro, Juan Fernandez-Tajes, Sarah E W Briggs, Victoria Svinti, Kevin Donnelly, Yingchang Lu, Fredrick R Schumacher, Stephanie J Weinstein, Kala Visvanathan, Kostas K Tsilidis, Yu-Ru Su, Robert Steinfeld, Sonja I Berndt, Sushma S Thomas, Kimberly F Doheny, Tameka Shelford, Amit D Joshi, Anshul Kundaje, Christopher K Edlund, Andre Kim, Lori C Sakoda, Stephanie A Bien, Yi Lin, Conghui Qu, Chenxu Qu, Stuart Reid, and Li Hsu analyzed the data. Xingyi Guo, Ceres Fernandez-Rozadilla, Jirong Long, Matthew Devall, Claire Palles, Kitty Sherwood, Susan M Farrington, James Blackmur, Peter G. Vaughan-Shaw, Xiao-Ou Shu, Peter Broderick, James Studd, Tabitha A Harrison, David V Conti, Marilena Melas, Gad Rennert, Mireia Obón-Santacana, Vicente Martín-Sánchez, Jae Hwan Oh, Jeongseon Kim, Sun Ha Jee, Keum Ji Jung, Sun-Seog Kweon, Min-Ho Shin, Aesun Shin, Yoon-Ok Ahn, Dong-Hyun Kim, Isao Oze, Wanqing Wen, Keitaro Matsuo, Koichi Matsuda, Chizu Tanikawa, Zefang Ren, Yu-Tang Gao, Wei-Hua Jia, John L Hopper, Mark A Jenkins, Aung Ko Win, Rish K Pai, Jane C Figueiredo, Robert W Haile, Steven Gallinger, Michael O Woods, Polly A Newcomb, David Duggan, Jeremy P. Cheadle, Richard Kaplan, Rachel Kerr, David Kerr, Iva Kirac, Jan Böhm, Jukka-Pekka Mecklin, Pekka Jousilahti, Paul Knekt, Lauri A. Aaltonen, Harri Rissanen, Eero Pukkala, Johan G Eriksson, Tatiana Cajuso, Ulrika Hänninen, Johanna Kondelin, Kimmo Palin, Tomas Tanskanen, Laura Renkonen-Sinisalo, Satu Männistö, Demetrius Albanes, Edward Ruiz-Narvaez, Julie R Palmer, Daniel D Buchanan, Elizabeth A Platz, Cornelia M Ulrich, Erin Siegel, Stefanie Brezina, Andrea Gsur, Peter T Campbell, Jenny Chang-Claude, Michael Hoffmeister, Hermann Brenner, Martha L Slattery, John D Potter, Matthias B Schulze, Marc J Gunter, Neil Murphy, Antoni Castells, Sergi Castellví-Bel, Leticia Moreira, Volker Arndt, Anna Shcherbina, D. Timothy Bishop, Graham G Giles, Melissa C. Southey, Gregory E Idos, Kevin J McDonnell, Zomoroda Abu-Ful, Joel K Greenson, Katerina Shulman, Flavio Lejbkovicz, Kenneth Offit, Temitope O Keku, Bethany van Guelpen, Thomas J Hudson, Heather Hampel, Rachel Pearlman, Richard B Hayes, Marie Elena Martinez, Paul D. P. Pharoah, Susanna C Larsson, Yun Yen, Heinz-Josef Lenz, Emily White, Li Li, Elizabeth Pugh, Andrew T Chan, Marcia Cruz-Correa, Annika Lindblom, David J Hunter, Clemens Schafmayer, Peter C Scacheri, Robert E Schoen, Jochen Hampe, Zsafia K Stadler, Pavel Vodicka, Ludmila Vodickova, Veronika Vymetalkova, W. James Gauderman, David

Shibata, Amanda Toland, Sanford Markowitz, Stephen J Chanock, Franzel van Duijnhoven, Edith JM Feskens, Manuela Gago-Dominguez, Alicia Wolk, Barbara Pardini, Liesel M FitzGerald, Soo Chin Lee, Shuji Ogino, Charles Kooperberg, Christopher I Li, Ross Prentice, Stéphane Bézieau, Taiki Yamaji, Norie Sawada, Motoki Iwasaki, Loic Le Marchand, Anna H Wu, Caroline E McNeil, Gerhard Coetzee, Caroline Hayward, Ian J Deary, Sarah E Harris, Evropi Theodoratou, Marion Walker, Li Yin Ooi, Qiuyin Cai, Malcolm G Dunlop, Stephen B Gruber, Richard S Houlston, Victor Moreno, Graham Casey, Ulrike Peters, Ian Tomlinson, and Wei Zheng recruited patients and collected samples. Zhishan Chen, Xingyi Guo, and Wei Zheng wrote the manuscript with substantial contributions from Ceres Fernandez-Rozadilla, Jie Ping, Guochong Jia, Jirong Long, Xiao-Ou Shu, Richard S Houlston, and Ian Tomlinson. All authors have reviewed and approved the final manuscript.

## Competing interests

Antoni Castells is a consultant to Bayer Pharma AG, Boehringer Ingelheim and Pfizer Inc. for work unrelated to this manuscript. Anna Shcherbina is an employee at Insitro, including consulting fees from BMS. Heather Hampel is SAB for Invitae Genetics, Promega and Genome Medical, Stock/Stock options for Genome Medical and GI OnDemand. Rish K Pai collaborates with Eli Lilly, AbbVie, Allergan, Verily and Alimientiv, which includes consulting fees (outside the submitted work). Stephanie A Bien has a financial interest in Adaptive Biotechnologies. Stephen B Gruber is co-founder, Brogent International LLC. One of Zsafia K Stadler's immediate family members serves as a consultant in ophthalmology for Alcon, Adverum, Gyroscope Therapeutics Limited, Neurogene and RegenexBio (outside the submitted work). Victor Moreno has research projects and owns stocks of Aniling. The remaining authors declare no competing interests.

## Additional information

**Supplementary information** The online version contains supplementary material available at <https://doi.org/10.1038/s41467-024-47399-x>.

**Correspondence** and requests for materials should be addressed to Wei Zheng.

**Peer review information** *Nature Communications* thanks Jyotsna Batra and Juliet French for their contribution to the peer review of this work. A peer review file is available.

**Reprints and permissions information** is available at <http://www.nature.com/reprints>

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

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2024

Zhishan Chen<sup>1</sup>, Xingyi Guo<sup>1,2</sup>, Ran Tao<sup>3,4</sup>, Jeroen R. Huyghe<sup>5</sup>, Philip J. Law<sup>6</sup>, Ceres Fernandez-Rozadilla<sup>7,8</sup>, Jie Ping<sup>1</sup>, Guochong Jia<sup>1</sup>, Jirong Long<sup>1</sup>, Chao Li<sup>1</sup>, Quanhu Shen<sup>3</sup>, Yuhan Xie<sup>9</sup>, Maria N. Timofeeva<sup>10,11</sup>, Minta Thomas<sup>5</sup>, Stephanie L. Schmit<sup>12,13</sup>, Virginia Díez-Obrero<sup>14,15,16,17</sup>, Matthew Devall<sup>18</sup>, Ferran Moratalla-Navarro<sup>14,15,16,17</sup>, Juan Fernandez-Tajes<sup>7</sup>, Claire Palles<sup>19</sup>, Kitty Sherwood<sup>7</sup>, Sarah E. W. Briggs<sup>20</sup>, Victoria Svinti<sup>10</sup>, Kevin Donnelly<sup>10</sup>, Susan M. Farrington<sup>10</sup>, James Blackmur<sup>10</sup>, Peter G. Vaughan-Shaw<sup>10</sup>, Xiao-Ou Shu<sup>1</sup>, Yingchang Lu<sup>1</sup>, Peter Broderick<sup>6</sup>, James Studd<sup>6</sup>, Tabitha A. Harrison<sup>5</sup>, David V. Conti<sup>21</sup>, Fredrick R. Schumacher<sup>22,23</sup>, Marilena Melas<sup>24</sup>, Gad Rennert<sup>25,26,27</sup>, Mireia Obón-Santacana<sup>14,15,17</sup>, Vicente Martín-Sánchez<sup>15,28</sup>, Jae Hwan Oh<sup>29</sup>, Jeongseon Kim<sup>30</sup>, Sun Ha Jee<sup>31</sup>, Keum Ji Jung<sup>31</sup>, Sun-Seog Kweon<sup>32</sup>, Min-Ho Shin<sup>32</sup>, Aesun Shin<sup>33,34</sup>, Yoon-Ok Ahn<sup>34</sup>, Dong-Hyun Kim<sup>35</sup>, Isao Oze<sup>36</sup>, Wanqing Wen<sup>1</sup>, Keitaro Matsuo<sup>37,38</sup>, Koichi Matsuda<sup>39</sup>, Chizu Tanikawa<sup>40</sup>, Zefang Ren<sup>41</sup>, Yu-Tang Gao<sup>42</sup>, Wei-Hua Jia<sup>43</sup>, John L. Hopper<sup>44,45</sup>, Mark A. Jenkins<sup>44</sup>, Aung Ko Win<sup>44</sup>, Rish K. Pai<sup>46</sup>, Jane C. Figueiredo<sup>21,47</sup>, Robert W. Haile<sup>48</sup>, Steven Gallinger<sup>49</sup>, Michael O. Woods<sup>50</sup>, Polly A. Newcomb<sup>5,51</sup>, David Duggan<sup>52</sup>, Jeremy P. Cheadle<sup>53</sup>, Richard Kaplan<sup>54</sup>, Rachel Kerr<sup>55</sup>, David Kerr<sup>56</sup>, Iva Kirac<sup>57</sup>, Jan Böhm<sup>58</sup>, Jukka-Pekka Mecklin<sup>59</sup>, Pekka Jousilahti<sup>60</sup>, Paul Knekt<sup>60</sup>, Lauri A. Aaltonen<sup>61,62</sup>, Harri Rissanen<sup>63</sup>, Eero Pukkala<sup>64,65</sup>, Johan G. Eriksson<sup>66,67,68</sup>, Tatiana Cajuso<sup>61,62</sup>, Ulrika Hänninen<sup>61,62</sup>, Johanna Kondelin<sup>61,62</sup>, Kimmo Palin<sup>61,62</sup>, Tomas Tanskanen<sup>61,62</sup>, Laura Renkonen-Sinisalo<sup>69</sup>, Satu Männistö<sup>63</sup>, Demetrius Albanes<sup>70</sup>, Stephanie J. Weinstein<sup>70</sup>, Edward Ruiz-Narvaez<sup>71</sup>, Julie R. Palmer<sup>72,73</sup>, Daniel D. Buchanan<sup>74,75,76</sup>, Elizabeth A. Platz<sup>77</sup>, Kala Visvanathan<sup>77</sup>, Cornelia M. Ulrich<sup>78</sup>, Erin Siegel<sup>79</sup>, Stefanie Brezina<sup>80</sup>, Andrea Gsur<sup>80</sup>, Peter T. Campbell<sup>81</sup>, Jenny Chang-Claude<sup>82,83</sup>, Michael Hoffmeister<sup>84</sup>, Hermann Brenner<sup>84,85,86</sup>, Martha L. Slattery<sup>87</sup>, John D. Potter<sup>5,88</sup>, Kostas K. Tsilidis<sup>89,90</sup>, Matthias B. Schulze<sup>91,92</sup>, Marc J. Gunter<sup>93</sup>, Neil Murphy<sup>93</sup>, Antoni Castells<sup>94</sup>, Sergi Castellví-Bel<sup>94</sup>, Leticia Moreira<sup>94</sup>, Volker Arndt<sup>84</sup>, Anna Shcherbina<sup>95</sup>, D. Timothy Bishop<sup>96</sup>, Graham G. Giles<sup>44,97,98</sup>, Melissa C. Southey<sup>97,98,99</sup>, Gregory E. Idos<sup>100</sup>, Kevin J. McDonnell<sup>25,27,100</sup>, Zomoroda Abu-Ful<sup>26</sup>, Joel K. Greenson<sup>25,27,101</sup>, Katerina Shulman<sup>26</sup>, Flavio Lejbkowitz<sup>25,26,102</sup>, Kenneth Offit<sup>103,104</sup>, Yu-Ru Su<sup>105</sup>, Robert Steinfeldt<sup>5</sup>, Temitope O. Keku<sup>106</sup>, Bethany van Guelpen<sup>107,108</sup>, Thomas J. Hudson<sup>109</sup>, Heather Hampel<sup>110</sup>, Rachel Pearlman<sup>110</sup>, Sonja I. Berndt<sup>70</sup>, Richard B. Hayes<sup>111</sup>, Marie Elena Martinez<sup>112,113</sup>, Sushma S. Thomas<sup>114</sup>, Paul D. P. Pharoah<sup>115</sup>, Susanna C. Larsson<sup>116</sup>, Yun Yen<sup>117</sup>, Heinz-Josef Lenz<sup>118</sup>, Emily White<sup>5,119</sup>, Li Li<sup>22</sup>, Kimberly F. Doherty<sup>120</sup>, Elizabeth Pugh<sup>120</sup>, Tameka Shelford<sup>120</sup>, Andrew T. Chan<sup>121,122,123,124,125,126</sup>, Marcia Cruz-Correa<sup>127</sup>, Annika Lindblom<sup>128,129</sup>, David J. Hunter<sup>124,130</sup>, Amit D. Joshi<sup>123,124</sup>, Clemens Schafmayer<sup>131</sup>, Peter C. Scacheri<sup>132</sup>, Anshul Kundaje<sup>95,133</sup>, Robert E. Schoen<sup>134</sup>, Jochen Hampe<sup>135</sup>, Zsafia K. Stadler<sup>104,136</sup>, Pavel Vodicka<sup>137,138,139</sup>, Ludmila Vodickova<sup>137,138,139</sup>, Veronika Vymetalkova<sup>137,138,139</sup>, Christopher K. Edlund<sup>21</sup>, W. James Gauderman<sup>21</sup>, David Shibata<sup>140</sup>, Amanda Toland<sup>141</sup>, Sanford Markowitz<sup>142</sup>, Andre Kim<sup>21</sup>, Stephen J. Chanock<sup>70</sup>, Franzel van Duijnhoven<sup>143</sup>, Edith J. M. Feskens<sup>144</sup>, Lori C. Sakoda<sup>5,145</sup>, Manuela Gago-Dominguez<sup>146,147</sup>, Alicja Wolk<sup>116</sup>, Barbara Pardini<sup>148,149</sup>, Liesel M. FitzGerald<sup>97,150</sup>, Soo Chin Lee<sup>151</sup>, Shuji Ogino<sup>124,152,153,154</sup>, Stephanie A. Bien<sup>5</sup>, Charles Kooperberg<sup>5</sup>, Christopher I. Li<sup>5</sup>, Yi Lin<sup>5</sup>, Ross Prentice<sup>5,155</sup>, Conghui Qu<sup>5</sup>, Stéphane Bézieau<sup>156</sup>, Taiki Yamaji<sup>157</sup>, Norie Sawada<sup>158</sup>, Motoki Iwasaki<sup>157,158</sup>, Loic Le Marchand<sup>159</sup>, Anna H. Wu<sup>160</sup>, Chenxu Qu<sup>161</sup>, Caroline E. McNeil<sup>161</sup>, Gerhard Coetzee<sup>162</sup>, Caroline Hayward<sup>163</sup>, Ian J. Deary<sup>164</sup>, Sarah E. Harris<sup>164</sup>, Evropi Theodoratou<sup>165</sup>, Stuart Reid<sup>10</sup>, Marion Walker<sup>10</sup>, Li Yin Ooi<sup>10,166</sup>, Ken S. Lau<sup>167</sup>, Hongyu Zhao<sup>9,168,169</sup>, Li Hsu<sup>5,170</sup>, Qiuyin Cai<sup>1</sup>, Malcolm G. Dunlop<sup>10</sup>, Stephen B. Gruber<sup>100</sup>, Richard S. Houlston<sup>6</sup>, Victor Moreno<sup>14,15,16,17</sup>, Graham Casey<sup>18</sup>, Ulrike Peters<sup>5,171</sup>, Ian Tomlinson<sup>7,19</sup> & Wei Zheng<sup>1</sup>✉

<sup>1</sup>Division of Epidemiology, Department of Medicine, Vanderbilt-Ingram Cancer Center, Vanderbilt Epidemiology Center, Vanderbilt University Medical Center, Nashville, TN, USA. <sup>2</sup>Department of Biomedical Informatics, Vanderbilt University School of Medicine, Nashville, TN, USA. <sup>3</sup>Department of Biostatistics, Vanderbilt University Medical Center, Nashville, TN, USA. <sup>4</sup>Vanderbilt Genetics Institute, Vanderbilt University Medical Center, Nashville 37232 TN, USA. <sup>5</sup>Public Health Sciences Division, Fred Hutchinson Cancer Center, Seattle, WA, USA. <sup>6</sup>Division of Genetics and Epidemiology, Institute of Cancer Research, London, UK. <sup>7</sup>Edinburgh Cancer Research Centre, Institute of Genomics and Cancer, University of Edinburgh, Edinburgh, UK. <sup>8</sup>Genomic Medicine Group, Instituto de Investigacion Sanitaria de Santiago, Santiago de Compostela, Spain. <sup>9</sup>Department of Biostatistics, Yale School of Public Health, New Haven, CT, USA. <sup>10</sup>Colon Cancer Genetics Group, Medical Research Council Human Genetics Unit, Institute of Genetics and Cancer, University of Edinburgh, Edinburgh, UK. <sup>11</sup>Danish Institute for Advanced Study, Department of Public Health, University of Southern Denmark, Odense, Denmark. <sup>12</sup>Genomic Medicine Institute, Cleveland Clinic, Cleveland, OH, USA. <sup>13</sup>Population and Cancer Prevention Program, Case Comprehensive Cancer Center, Cleveland, OH, USA. <sup>14</sup>Colorectal Cancer Group, ONCOBELL Program, Bellvitge Biomedical Research Institute, Barcelona, Spain. <sup>15</sup>Consortium for Biomedical Research in Epidemiology and Public Health, Madrid, Spain. <sup>16</sup>Department of Clinical Sciences, Faculty of Medicine, University of Barcelona, Barcelona, Spain. <sup>17</sup>Oncology Data Analytics Program, Catalan Institute of Oncology, Barcelona, Spain. <sup>18</sup>Center for Public Health Genomics, Department of Public Health Sciences, University of Virginia, Charlottesville, VA, USA. <sup>19</sup>Institute of Cancer and Genomic Sciences, College of Medical and Dental Sciences, University of Birmingham, Birmingham, UK. <sup>20</sup>Department of Public Health, Richard Doll Building, University of Oxford, Oxford, UK. <sup>21</sup>Department of Preventive Medicine, USC Norris Comprehensive Cancer Center, Keck School of Medicine, University of Southern California, Los Angeles, CA, USA. <sup>22</sup>Case Comprehensive Cancer Center, Case Western Reserve University, Cleveland, OH, USA. <sup>23</sup>Department of Population and Quantitative Health Sciences, Case Western Reserve University, Cleveland, OH, USA. <sup>24</sup>The Steve and Cindy Rasmussen Institute for Genomic Medicine, Nationwide Children's Hospital, Columbus, OH, USA.

<sup>25</sup>Clalit National Cancer Control Center, Haifa, Israel. <sup>26</sup>Department of Community Medicine and Epidemiology, Lady Davis Carmel Medical Center, Haifa, Israel. <sup>27</sup>Ruth and Bruce Rappaport Faculty of Medicine, Technion-Israel Institute of Technology, Haifa, Israel. <sup>28</sup>Biomedicine Institute, University of León, León, Spain. <sup>29</sup>Center for Colorectal Cancer, National Cancer Center Hospital, National Cancer Center, Gyeonggi-do, South Korea. <sup>30</sup>Department of Cancer Biomedical Science, Graduate School of Cancer Science and Policy, National Cancer Center, Gyeonggi-do, South Korea. <sup>31</sup>Department of Epidemiology and Health Promotion, Graduate School of Public Health, Yonsei University, Seoul, South Korea. <sup>32</sup>Department of Preventive Medicine, Chonnam National University Medical School, Gwangju, South Korea. <sup>33</sup>Cancer Research Institute, Seoul National University, Seoul, South Korea. <sup>34</sup>Department of Preventive Medicine, Seoul National University College of Medicine, Seoul, South Korea. <sup>35</sup>Department of Social and Preventive Medicine, Hallym University College of Medicine, Okcheon-dong, South Korea. <sup>36</sup>Division of Cancer Epidemiology and Prevention, Aichi Cancer Center Research Institute, Nagoya, Japan. <sup>37</sup>Department of Epidemiology, Nagoya University Graduate School of Medicine, Nagoya, Japan. <sup>38</sup>Division of Molecular and Clinical Epidemiology, Aichi Cancer Center Research Institute, Nagoya, Japan. <sup>39</sup>Laboratory of Clinical Genome Sequencing, Department of Computational Biology and Medical Sciences, Graduate School of Frontier Sciences, University of Tokyo, Tokyo, Japan. <sup>40</sup>Laboratory of Genome Technology, Human Genome Center, Institute of Medical Science, University of Tokyo, Tokyo, Japan. <sup>41</sup>School of Public Health, Sun Yat-sen University, Guangzhou, China. <sup>42</sup>State Key Laboratory of Oncogenes and Related Genes and Department of Epidemiology, Shanghai Cancer Institute, Renji Hospital, Shanghai Jiaotong University School of Medicine, Shanghai, China. <sup>43</sup>State Key Laboratory of Oncology in South China, Cancer Center, Sun Yat-sen University, Guangzhou, China. <sup>44</sup>Centre for Epidemiology and Biostatistics, Melbourne School of Population and Global Health, University of Melbourne, Melbourne, VIC, Australia. <sup>45</sup>Department of Epidemiology, School of Public Health and Institute of Health and Environment, Seoul National University, Seoul, South Korea. <sup>46</sup>Department of Laboratory Medicine and Pathology, Mayo Clinic Arizona, Scottsdale, AZ, USA. <sup>47</sup>Department of Medicine, Samuel Oschin Comprehensive Cancer Institute, Cedars-Sinai Medical Center, Los Angeles, CA, USA. <sup>48</sup>Division of Oncology, Department of Medicine, Cedars-Sinai Cancer Research Center for Health Equity, Los Angeles, CA, USA. <sup>49</sup>Lunenfeld Tanenbaum Research Institute, Mount Sinai Hospital, University of Toronto, Toronto, ON, Canada. <sup>50</sup>Division of Biomedical Sciences, Memorial University of Newfoundland, St. John, ON, Canada. <sup>51</sup>School of Public Health, University of Washington, Seattle, WA, USA. <sup>52</sup>City of Hope National Medical Center, Translational Genomics Research Institute, Phoenix, AZ, USA. <sup>53</sup>Institute of Medical Genetics, Cardiff University, Cardiff, UK. <sup>54</sup>MRC Clinical Trials Unit, Medical Research Council, Cardiff, UK. <sup>55</sup>Department of Oncology, University of Oxford, Oxford, UK. <sup>56</sup>Radcliffe Department of Medicine, University of Oxford, Oxford, UK. <sup>57</sup>Department of Surgical Oncology, University Hospital for Tumors, Sestre milosrdnice University Hospital Center, Zagreb, Croatia. <sup>58</sup>Department of Pathology, Central Finland Health Care District, Jyväskylä, Finland. <sup>59</sup>Central Finland Health Care District, Jyväskylä, Finland. <sup>60</sup>Department of Health and Welfare, Finnish Institute for Health and Welfare, Helsinki, Finland. <sup>61</sup>Department of Medical and Clinical Genetics, University of Helsinki, Helsinki, Finland. <sup>62</sup>Genome-Scale Biology Research Program, University of Helsinki, Helsinki, Finland. <sup>63</sup>Department of Public Health and Welfare, Finnish Institute for Health and Welfare, Helsinki, Finland. <sup>64</sup>Faculty of Social Sciences, Tampere University, Tampere, Finland. <sup>65</sup>Finnish Cancer Registry, Institute for Statistical and Epidemiological Cancer Research, Helsinki, Finland. <sup>66</sup>Folkhälsan Research Centre, University of Helsinki, Helsinki, Finland. <sup>67</sup>Human Potential Translational Research Programme, National University of Singapore, Singapore, Singapore. <sup>68</sup>Unit of General Practice and Primary Health Care, University of Helsinki and Helsinki University Hospital, Helsinki, Finland. <sup>69</sup>Department of Surgery, Abdominal Centre, Helsinki University Hospital, Helsinki, Finland. <sup>70</sup>Division of Cancer Epidemiology and Genetics, National Cancer Institute, National Institutes of Health, Bethesda, MD, USA. <sup>71</sup>Department of Nutritional Sciences, School of Public Health, University of Michigan, Ann Arbor, MI, USA. <sup>72</sup>Department of Medicine, Boston University School of Medicine, Boston, MA, USA. <sup>73</sup>Slone Epidemiology Center at Boston University, Boston, MA, USA. <sup>74</sup>Colorectal Oncogenomics Group, Department of Clinical Pathology, University of Melbourne, Parkville, VIC, Australia. <sup>75</sup>Genomic Medicine and Family Cancer Clinic, Royal Melbourne Hospital, Parkville, VIC, Australia. <sup>76</sup>University of Melbourne Centre for Cancer Research, Victorian Comprehensive Cancer Centre, Parkville, VIC, Australia. <sup>77</sup>Department of Epidemiology, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD, USA. <sup>78</sup>Huntsman Cancer Institute and Department of Population Health Sciences, University of Utah, Salt Lake City, UT, USA. <sup>79</sup>Cancer Epidemiology Program, H. Lee Moffitt Cancer Center and Research Institute, Tampa, FL, USA. <sup>80</sup>Institute of Cancer Research, Department of Medicine I, Medical University Vienna, Vienna, Austria. <sup>81</sup>Department of Epidemiology and Population Health, Albert Einstein College of Medicine, New York, NY, USA. <sup>82</sup>Division of Cancer Epidemiology, German Cancer Research Center, Heidelberg, Germany. <sup>83</sup>University Medical Centre Hamburg-Eppendorf, University Cancer Centre Hamburg, Hamburg, Germany. <sup>84</sup>Division of Clinical Epidemiology and Aging Research, German Cancer Research Center, Heidelberg, Germany. <sup>85</sup>Division of Preventive Oncology, German Cancer Research Center and National Center for Tumor Diseases, Heidelberg, Germany. <sup>86</sup>German Cancer Consortium, German Cancer Research Center, Heidelberg, Germany. <sup>87</sup>Department of Internal Medicine, University of Utah, Salt Lake City, UT, USA. <sup>88</sup>Research Centre for Hauora and Health, Massey University, Wellington, New Zealand. <sup>89</sup>Department of Epidemiology and Biostatistics, School of Public Health, Imperial College London, London, UK. <sup>90</sup>Department of Hygiene and Epidemiology, University of Ioannina School of Medicine, Ioannina, Greece. <sup>91</sup>Department of Molecular Epidemiology, German Institute of Human Nutrition Potsdam-Rehbruecke, Nuthetal, Germany. <sup>92</sup>Institute of Nutritional Science, University of Potsdam, Potsdam, Germany. <sup>93</sup>Nutrition and Metabolism Branch, International Agency for Research on Cancer, World Health Organization, Lyon, France. <sup>94</sup>Gastroenterology Department, Hospital Clinic, Institut d'Investigacions Biomèdiques August Pi i Sunyer, Centro de Investigación Biomédica en Red de Enfermedades Hepáticas y Digestivas, University of Barcelona, Barcelona, Spain. <sup>95</sup>Department of Genetics, Stanford University, Stanford, CA, USA. <sup>96</sup>Leeds Institute of Medical Research at St James's, University of Leeds, Leeds, UK. <sup>97</sup>Cancer Epidemiology Division, Cancer Council Victoria, Melbourne, VIC, Australia. <sup>98</sup>Precision Medicine, School of Clinical Sciences at Monash Health, Monash University, Clayton, VIC, Australia. <sup>99</sup>Department of Clinical Pathology, University of Melbourne, Melbourne, VIC, Australia. <sup>100</sup>Department of Medical Oncology and Center For Precision Medicine, City of Hope National Medical Center, Duarte, CA, USA. <sup>101</sup>Department of Pathology, University of Michigan, Ann Arbor, MI, USA. <sup>102</sup>Clalit Health Services, Personalized Genomic Service, Lady Davis Carmel Medical Center, Haifa, Israel. <sup>103</sup>Clinical Genetics Service, Department of Medicine, Memorial Sloan-Kettering Cancer Center, New York, NY, USA. <sup>104</sup>Department of Medicine, Weill Cornell Medical College, New York, NY, USA. <sup>105</sup>Kaiser Permanente Washington Health Research Institute, Seattle, WA, USA. <sup>106</sup>Center for Gastrointestinal Biology and Disease, University of North Carolina, Chapel Hill, NC, USA. <sup>107</sup>Department of Radiation Sciences, Oncology Unit, Umeå University, Umeå, Sweden. <sup>108</sup>Wallenberg Centre for Molecular Medicine, Umeå University, Umeå, Sweden. <sup>109</sup>Ontario Institute for Cancer Research, Toronto, ON, Canada. <sup>110</sup>Division of Human Genetics, Department of Internal Medicine, Ohio State University Comprehensive Cancer Center, Columbus, OH, USA. <sup>111</sup>Division of Epidemiology, Department of Population Health, New York University School of Medicine, New York, NY, USA. <sup>112</sup>Department of Family Medicine and Public Health, University of California San Diego, La Jolla, CA, USA. <sup>113</sup>Population Sciences, Disparities and Community Engagement, University of California San Diego Moores Cancer Center, La Jolla, CA, USA. <sup>114</sup>Fred Hutchinson Cancer Center, Seattle, WA, USA. <sup>115</sup>Department of Public Health and Primary Care, University of Cambridge, Cambridge, UK. <sup>116</sup>Institute of Environmental Medicine, Karolinska Institutet, Stockholm, Sweden. <sup>117</sup>Taipei Medical University, Taipei, Taiwan. <sup>118</sup>Department of Medicine, Keck School of Medicine, University of Southern California, Los Angeles, CA, USA. <sup>119</sup>Department of Epidemiology, University of Washington School of Public Health, Seattle, WA, USA. <sup>120</sup>Center for Inherited Disease Research, Department of Genetic Medicine, Johns Hopkins University School of Medicine, Baltimore, MD, USA. <sup>121</sup>Broad Institute of Harvard and MIT, Cambridge, MA, USA. <sup>122</sup>Channing Division of Network Medicine, Brigham and Women's Hospital and Harvard Medical



School, Boston, MA, USA. <sup>123</sup>Clinical and Translational Epidemiology Unit, Massachusetts General Hospital and Harvard Medical School, Boston, MA, USA. <sup>124</sup>Department of Epidemiology, Harvard T.H. Chan School of Public Health, Harvard University, Boston, MA, USA. <sup>125</sup>Department of Immunology and Infectious Diseases, Harvard T.H. Chan School of Public Health, Harvard University, Boston, MA, USA. <sup>126</sup>Division of Gastroenterology, Massachusetts General Hospital and Harvard Medical School, Boston, MA, USA. <sup>127</sup>Comprehensive Cancer Center, University of Puerto Rico, San Juan, Puerto Rico. <sup>128</sup>Department of Clinical Genetics, Karolinska University Hospital, Stockholm, Sweden. <sup>129</sup>Department of Molecular Medicine and Surgery, Karolinska Institutet, Stockholm, Sweden. <sup>130</sup>Nuffield Department of Population Health, University of Oxford, Oxford, UK. <sup>131</sup>Department of General Surgery, University Hospital Rostock, Rostock, Germany. <sup>132</sup>Department of Genetics and Genome Sciences, Case Western Reserve University, Cleveland, OH, USA. <sup>133</sup>Department of Computer Science, Stanford University, Stanford, CA, USA. <sup>134</sup>Department of Medicine and Epidemiology, University of Pittsburgh Medical Center, Pittsburgh, PA, USA. <sup>135</sup>Department of Medicine I, University Hospital Dresden, Technische Universität Dresden, Dresden, Germany. <sup>136</sup>Department of Medicine, Memorial Sloan-Kettering Cancer Center, New York, NY, USA. <sup>137</sup>Department of Molecular Biology of Cancer, Institute of Experimental Medicine of the Czech Academy of Sciences, Prague, Czech Republic. <sup>138</sup>Faculty of Medicine and Biomedical Center in Pilsen, Charles University, Pilsen, Czech Republic. <sup>139</sup>Institute of Biology and Medical Genetics, First Faculty of Medicine, Charles University, Prague, Czech Republic. <sup>140</sup>Department of Surgery, University of Tennessee Health Science Center, Memphis, TN, USA. <sup>141</sup>Departments of Cancer Biology and Genetics and Internal Medicine, Comprehensive Cancer Center, Ohio State University, Columbus, OH, USA. <sup>142</sup>Departments of Medicine and Genetics, Case Comprehensive Cancer Center, Case Western Reserve University and University Hospitals of Cleveland, Cleveland, OH, USA. <sup>143</sup>Division of Human Nutrition and Health, Wageningen University and Research, Wageningen, The Netherlands. <sup>144</sup>Division of Human Nutrition, Wageningen University and Research, Wageningen, The Netherlands. <sup>145</sup>Division of Research, Kaiser Permanente Northern California, Oakland, CA, USA. <sup>146</sup>Genomic Medicine Group, Galician Public Foundation of Genomic Medicine, Servicio Galego de Saude, Santiago de Compostela, Spain. <sup>147</sup>Instituto de Investigación Sanitaria de Santiago de Compostela, Santiago de Compostela, Spain. <sup>148</sup>Candiolo Cancer Institute FPO-IRCCS, Candiolo, (TO), Italy. <sup>149</sup>Italian Institute for Genomic Medicine, Candiolo Cancer Institute FPO-IRCCS, Candiolo, (TO), Italy. <sup>150</sup>Menzies Institute for Medical Research, University of Tasmania, Hobart, TAS, Australia. <sup>151</sup>National University Cancer Institute, Singapore, Cancer Science Institute of Singapore, National University of Singapore, Singapore, Singapore. <sup>152</sup>Broad Institute of MIT and Harvard, Cambridge, MA, USA. <sup>153</sup>Cancer Immunology Program, Dana-Farber Harvard Cancer Center, Boston, MA, USA. <sup>154</sup>Department of Pathology, Brigham and Women's Hospital, Harvard Medical School, Boston, MA, USA. <sup>155</sup>Department of Biostatistics, University of Washington, Seattle, WA, USA. <sup>156</sup>Service de Génétique Médicale, Centre Hospitalier Universitaire Nantes, Nantes, France. <sup>157</sup>Division of Epidemiology, National Cancer Center Institute for Cancer Control, National Cancer Center, Tokyo, Japan. <sup>158</sup>Division of Cohort Research, National Cancer Center Institute for Cancer Control, National Cancer Center, Tokyo, Japan. <sup>159</sup>Cancer Center, University of Hawaii, Honolulu, HI, USA. <sup>160</sup>Preventative Medicine, University of Southern California, Los Angeles, CA, USA. <sup>161</sup>USC Norris Comprehensive Cancer Center, Keck School of Medicine, University of Southern California, Los Angeles, CA, USA. <sup>162</sup>Van Andel Research Institute, Grand Rapids, MI, USA. <sup>163</sup>MRC Human Genetics Unit, Institute of Genomics and Cancer, University of Edinburgh, Edinburgh, UK. <sup>164</sup>Lothian Birth Cohorts group, Department of Psychology, University of Edinburgh, Edinburgh, UK. <sup>165</sup>Centre for Global Health, Usher Institute, University of Edinburgh, Edinburgh, UK. <sup>166</sup>Department of Pathology, National University Hospital, National University Health System, Singapore, Singapore. <sup>167</sup>Epithelial Biology Center and Department of Cell and Developmental Biology, Vanderbilt University School of Medicine, Nashville, TN, USA. <sup>168</sup>Department of Genetics, Yale School of Medicine, New Haven, CT, USA. <sup>169</sup>Program in Computational Biology and Bioinformatics, Yale University, New Haven, CT, USA. <sup>170</sup>Department of Biostatistics, School of Public Health, University of Washington, Seattle, WA, USA. <sup>171</sup>Department of Epidemiology, University of Washington, Seattle, WA, USA. ✉ e-mail: [wei.zheng@vanderbilt.edu](mailto:wei.zheng@vanderbilt.edu)