

Genetic risk factors for colorectal cancer in multiethnic Indonesians

by Upik A Miskad

Submission date: 03-Feb-2023 08:31AM (UTC+0700)

Submission ID: 2005183207

File name: Scientific_Report_Prof_Irawan.pdf (1.04M)

Word count: 6536

Character count: 36062



OPEN

22

Genetic risk factors for colorectal cancer in multiethnic Indonesians

Irawan Yusuf^{1,3,7}, Bens Pardamean^{2,4,7}, James W. Baurley^{2,7}, Arif Budiarto^{2,5},
Upik A. Miskad¹, Ronald E. Lusikooy¹, Arham Arsyad¹, Akram Irwan¹, George Mathew³,
Ivet Suriapranata³, Rinaldy Kusuma³, Muhamad F. Kacamarga^{2,5}, Tjeng W. Cenggoro^{2,5},
Christopher McMahan⁶, Chase Joyner⁶ & Carissa I. Pardamean²

Colorectal cancer is a common cancer in Indonesia, yet it has been understudied in this resource-constrained setting. We conducted a genome-wide association study focused on evaluation and preliminary discovery of colorectal cancer risk factors in Indonesians. We administered detailed questionnaires and collecting blood samples from 162 colorectal cancer cases throughout Makassar, Indonesia. We also established a control set of 193 healthy individuals frequency matched by age, sex, and ethnicity. A genome-wide association analysis was performed on 84 cases and 89 controls passing quality control. We evaluated known colorectal cancer genetic variants using logistic regression and established a genome-wide polygenic risk model using a Bayesian variable selection technique. We replicate associations for rs9497673, rs6936461 and rs7758229 on chromosome 6; rs11255841 on chromosome 10; and rs4779584, rs11632715, and rs73376930 on chromosome 15. Polygenic modeling identified 10 SNP associated with colorectal cancer risk. This work helps characterize the relationship between variants in the *SCL22A3*, *SCG5*, *GREM1*, and *STXBP5-AS1* genes and colorectal cancer in a diverse Indonesian population. With further biobanking and international research collaborations, variants specific to colorectal cancer risk in Indonesians will be identified.

Colorectal cancer is one of the most common cancers in the world and a leading cause of cancer-related deaths^{1,2}. There is growing evidence that colorectal cancer rates are changing in Asian countries, but the causes are still under investigation^{3,4}. Colorectal cancer is now one of the top three cancers in many Asian countries⁴. Currently, Asia contributes to 48% of the total number of new colorectal cancer cases in the world, of which the majority are found in Eastern Asia⁵. Specifically in Indonesia, the age-standardized incidence for males and females has been reported as 15.9 and 10.1 per 100,000 respectively⁶.

The heritability of colorectal cancer is estimated to be between 12 and 35%. However, germline mutations that are highly penetrant contribute less than 5% to colorectal cancer⁷. Nonetheless, increasing evidence is finding that heritability plays a potential, crucial role in colorectal cancer pathogenesis. Currently, mutations in 14 genes are suspected to underlie different subtypes of colorectal cancer, including mutations in the APC that increases predisposition to familial adenomatous polyposis (FAP) and defects in mismatch repair genes associated with Lynch Syndrome⁷. Recent genome-wide association studies have identified common genetic variants linked to colorectal cancer predisposition, highlighting a greater association between heritable risk and the disease. Thus far, over 40 genetic variants have been identified, within several well-known biological pathways that have been shown to be highly relevant to oncogenesis, including the TGF-beta/BMP pathway and the mitogen-activated protein kinases (MAPK) pathway⁷.

However, many of these colorectal cancer genetic associations were discovered in European-ancestry populations but do not replicate well in other ancestry groups, demonstrating the need for studies in diverse populations worldwide⁸. The Asia Colorectal Cancer Consortium was initiated in 2009 among East Asian nations and has successfully identified novel relevant, genetic regions^{9,10}. However, colorectal cancer cases from South East Asian cohorts have been under represented.

¹Faculty Medicine, Hasanuddin University, Makassar, South Sulawesi, Indonesia. ²Bioinformatics & Data Science Research Center, Bina Nusantara University, Jakarta, DKI Jakarta, Indonesia. ³Mochtar Riady Institute for Nanotechnology, Pelita Harapan University, Tangerang, Banten, Indonesia. ⁴Computer Science Department, BINUS Graduate Program-Master of Computer Science Program, Bina Nusantara University, Jakarta, DKI Jakarta, Indonesia. ⁵Computer Science Department, School of Computer Science, Bina Nusantara University, Jakarta, DKI Jakarta, Indonesia. ⁶School of Mathematical and Statistical Sciences, Clemson University, Clemson, SC, USA. ⁷These authors contributed equally: Irawan Yusuf, Bens Pardamean, and James W. Baurley. ✉ email: bpardamean@binus.edu; baurley@binus.edu

	Cases	Controls	P
	N = 89	N = 84	
Age	53.8 (13.2)	50.5 (14.5)	0.12
Gender			> 0.99
Female	38 (42.7%)	36 (42.9%)	
Male	51 (57.3%)	48 (57.1%)	
Ethnicity			0.68
Bugis	39 (43.8%)	45 (53.6%)	
Makassar	24 (27.0%)	23 (27.4%)	
Mandar	2 (2.3%)	1 (1.2%)	
Toraja	10 (11.2%)	8 (9.5%)	
Non South Sulawesi	9 (10.1%)	4 (4.8%)	
Non Sulawesi	5 (5.6%)	3 (3.6%)	
BMI	21.2 (3.1)	24.5 (3.6)	< 0.01
Smoking status			< 0.01
Smoker	39 (43.8%)	15 (17.9%)	
Non smoker	50 (56.2%)	69 (82.1%)	
Ancestry (estimated)			
East Asian (EAS)	0.92	0.94	0.02
South Asian (SAS)	0.07	0.05	0.15
African (AFR)	< 0.01	< 0.01	0.02
European (EUR)	0.01	0.01	0.36
Cancer site			
Right colon	15 (16.9%)	–	
Transversum	9 (10.1%)	–	
Left colon	1 (1.12%)	–	
Sigmoid Rectum	26 (29.2%)	–	
	38 (42.7%)	–	
Staging			
I	3 (3.4%)	–	
II	9 (10.1%)	–	
III	62 (69.7%)	–	
IV	11 (12.4%)	–	

Table 1. Characteristics of South Sulawesi colorectal cancer cases and controls.

Given the changes in rectal cancer rates in Asia and the differences in risk factors present in ethnically diverse South East Asia, we present results of the first genomic association study of colorectal cancer in Indonesia. We present results from the initial phase of this study, focused on cases from South Sulawesi, Indonesia.

Results

Characteristics of study sample. The characteristics of the colorectal cancer cases and controls are summarized in Table 1. The mean age of the colorectal cancer cases was 54 years. The majority of cases were male (57%). Among ethnicities, most cases were self-reported Bugis (44%) or Makassar ethnicity (27%). Controls appeared to be adequately frequency matched to cases by age, sex, and ethnicity ($p > 0.05$). Colorectal cancer cases had lower average body mass index (BMI) and were more likely to be smokers than controls ($p < 0.01$). Estimated genetically, the majority of both cases and controls were of East Asian ancestry. 82% of the cases had late stage cancer (III or IV) which unfortunately is consistent with recent reports in Indonesia¹¹. As seen in other studies, the most common colorectal cancer site was rectum (43%)^{12,13}.

Genome-wide association analysis. As expected given the sample size, no SNPs met the historical cut-off set for genome-wide significance (Supplementary Figs. 6 and 7). The summaries for all variants with a marginal p -value $< 5E-5$ are included in the “Supplementary materials” (Table 4). These include two intergenic SNPs and two SNPs in the *MRO* gene on chromosome 18.

Results for previously reported colorectal cancer SNPs are presented in Fig. 1 and Supplementary Table 3. There is evidence of replication for the following genetic variants: rs9497673, rs6936461 and rs7758229 on chromosome 6; rs11255841 on chromosome 10; and rs4779584, rs11632715, and rs73376930 on chromosome 15. The regions are characterized in Figs. 2, 3, 4, and 5. The pattern of associations is rather diffuse in the *STXBP5-AS1* (*STXBP5* Antisense RNA 1) and *SLC22A3* genes of chromosome 6, representing the correlation among

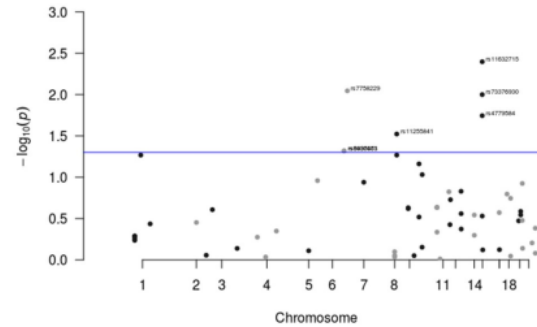


Figure 1. Results for known colorectal cancer susceptibility SNPs. Variants with p-values < 0.05 were flagged for further investigation.

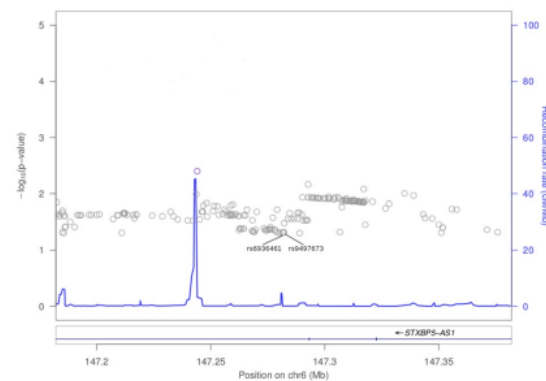


Figure 2. Association plot for 100 kb region flanking rs6936461 on chromosome 6.

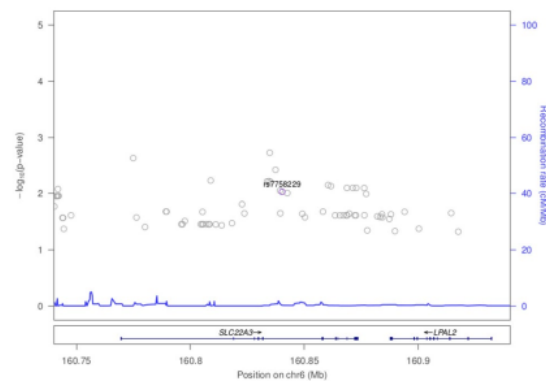


Figure 3. Association plot for 100 kb region flanking rs7758229 on chromosome 6.

the variants in these regions (Figs. 2 and 3). Similarly, the association pattern tapers along chromosome 10. The strongest association pattern can be found on chromosome 15. This region has 16 more defined peak than the other regions with associations spanning two genes: *SCG5* (secretogranin V) and *GREM1* (gremlin 1, DAN family BMP antagonist).

The polygenic 10 analysis identified 10 SNPs which appear to have a relatively strong association (i.e., large effect size) with the risk of developing colorectal cancer as can be seen in Table 2. These variants have marginal

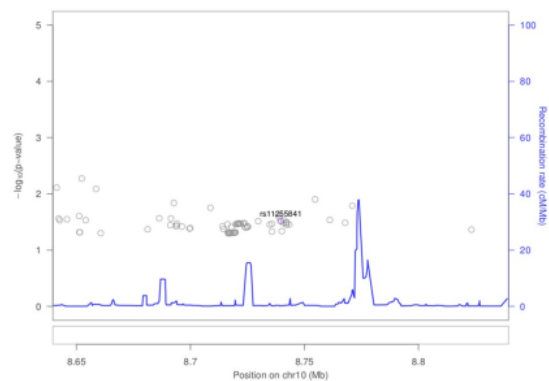


Figure 4. Association plot for 100 kb region flanking rs11255841 on chromosome 10.

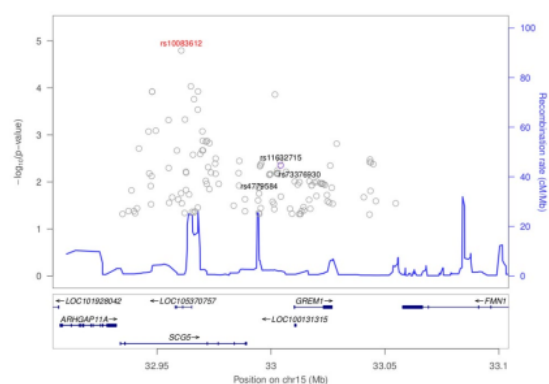


Figure 5. Association plot for 100 kb flanking rs11632715 on chromosome 15. The top associated SNP in the region was rs10083612.

Description	Chr	Position	Gene	Ref	MaF	Estimate
Intercept						0.90
Gender						0.00
Age						-3.75
BMI						0.00
Smoking						1.32
rs11919079	3	57086348	Intron:ARHGEF3	G	0.07	2.40
rs4888186	16	81947156	Intron:PLCG2	C	0.08	0.85
rs11016111	10	129963848	Intergenic	C	0.34	-1.32
rs77657157	5	98125016	Intron:RGMB	G	0.05	1.95
-	18	59822981	Deletion:PIGN	TC	0.19	-1.39
rs17066763	5	164113078	Intergenic	T	0.12	1.65
rs2446103	6	77328692	Intergenic	A	0.04	1.22
rs7219420	17	45800299	Intergenic	T	0.36	1.32
-	16	13018917	Insertion:SHISA9	C	0.11	1.67
rs78165118	3	12816282	Intergenic	A	0.03	2.13

Table 2. Polygenic risk model learned from colorectal cancer data. Presented results include the chromosome (Chr) and position of the significant genetic variants, the gene they lie on (Gene), reference allele (Ref), minor allele frequency (MaF), and estimated effect (Estimate).

p-values between 0.19 and $1.5E-5$ indicating some would have been overlooked in a standard analysis. Five of these SNPs lie in intergenic regions; three lie in introns of *ARHGEF3*, *PLCG2*, and *RGMB*; one is a deletion in *PIGN*; and one is an insertion in *SHISA9*.

Discussion

This preliminary study represents the first genome-wide analysis of a South Sulawesi population in Indonesia. We hope this work will motivate additional cancer research in this understudied and diverse population. Strengths of the study include the building of a colorectal cancer research program in Indonesia, the extensive questionnaire for assessing non-genetic risk factors, and genome-wide genotyping across diverse ethnicities.

Limitations of the study include the sample size due to the resource-constrained settings in Indonesia, which restricts the analysis to previously identified colorectal cancer markers and challenges shared by case-control study designs. For instance, the controls may represent different groups than cases. We attempted to account for this by frequency matching on age, sex, and ethnicity. Additionally, the timing of assessments need to be considered in interpreting the results. Given screening programs are still being developed in Indonesia, the majority of the cases had late stage colorectal cancer, stage III and IV. When BMI was assessed in these patients they already had significant weight loss, thus the direction of the effect is different than what one might expect.

Interestingly, the mean age of cases in this study was 54 which could imply a family history of cancer. Unfortunately we had limited data on family history because patients from the rural areas did not know the health history of their relatives. Indonesia also lacks a cancer registry which could also provide information on family histories of cancers. Also worth noting, the majority of the cases had rectal cancer. Recent work from Deng¹⁴ found that Asian countries appear to have higher rates of rectal cancer than western countries. Environmental factors are suspected to play a strong role, e.g., in this study we found that rectal cancer cases were more likely to be smokers.

For genome-wide imputation, an Indonesian population is not currently represented in common reference population such as the 1000 Genomes Project, thus some genetic markers relevant to colorectal cancer and specific to Indonesians may not impute well. However, the 1000 Genomes Project does have samples from Vietnam. There are genomic diversity studies underway in South East Asia which may offer a suitable reference panel for Indonesians in the future¹⁵.

Several previously identified colorectal cancer associated SNPs replicated in this population. And we can begin characterizing these regions by examining neighboring variants. The rs7758229 variant within *SLC22A3* on chromosome 6 was originally identified and subsequently replicated in large case-control study of a Japanese population (OR of 1.3)¹⁶. Interestingly, in a subsequent study in a Chinese population, this SNP was not associated with colorectal cancer (OR of 0.95)¹⁷. However, in S. Sulawesi, we detect a statistically significant association with colorectal cancer ($p = 0.009$, OR of 2.2). Given these difference among East Asians, further work to understand variation in *SLC22A3* and colorectal cancer is needed. *SLC22A3* encodes for the protein OCT3, which is an organic cationic transporter. While OCT3/*SLC22A3* is well characterized⁴⁴ in neurochemistry, it has been found to play a role within oncology as well. The upregulation of *SLC22A3* in head and neck squamous cell carcinoma is associated with improved prognosis while the downregulation of *SLC22A3* leads to enhanced metastasis and invasion of the tumor¹⁸. *SLC22A3* has also been implicated in the pathogenesis of prostate cancer and its expression is elevated in these neoplastic tissues¹⁹. The level of OCT3/*SLC22A3* expression has also been linked to the level of patient responsiveness towards cancer treatments²⁰; in particular, platinum-based cytotoxic cancer treatments in colorectal cancer²¹ patients, as well as head and neck squamous cell carcinoma patients¹⁸.

Intergenic variant rs11255841 on chromosome 10 was identified in an colorectal cancer GWAS of European ancestry individuals²² and has replicated in a Japanese study and a large meta-analysis with nearly 37,000 cases^{23,24}. With the risk allele of T, this variant had an odds ratio of 2.2 in our study, while previous reports had an odds ratio of 1.1–1.2.

The region on chromosome 15 nearby *SCG5* and *GRL42* have been flagged in multiple GWAS, e.g.,²⁵. We replicated colorectal cancer associations for rs4779584 ($p = 0.018$), rs11632715 ($p = 0.004$), and rs73376930 ($p = 0.010$). Interestingly, the smallest p-value in the region was rs10083612 within an intron of *SCG5* ($p = 1.61E-5$, see Fig. 5). The role of *SCG5* in colorectal cancer has not been well characterized, while much is known about its neighbor *GREM1*'s role in colorectal cancer. *GREM1*, which is one of the antagonists of the bone morphogenetic proteins (BMPs) found within the TGF-beta signaling pathway, has been found to be important for the survival and proliferation of several types of cancers²⁶. In particular, modulated expression of *GREM1* is found in cancer-associated stromal cells. *GREM1* is also found to be a proangiogenic factor, suggesting a role in cancer development when it is upregulated²⁷. *SCG5* and *GREM1* genes have been found¹⁹ to be associated with polyposis syndromes that are associated with colorectal cancer²⁸. A duplication⁴⁶ on that spans the 3' end of *SCG5* and the immediate, adjacent upstream region of *GREM1* is associated with hereditary mixed polyposis syndrome (HMPS) as well as tumorigenesis in juvenile polyposis. This duplication results in a 40-kb extra segment that leads to the upregulation of *GREM1* expression. The duplication is the basis for an autosomal dominant HMPS condition that is prevalent among the Ashkenazi Jewish population and is a recommended biomarker/genetic test to detect CRC in this population. Aberrant expression of *GREM1* has also been shown to underlie oncogenesis within the large intestines and colon²⁹.

Two of the previously identified colorectal cancer markers replicate in this study (rs6936461 and rs9497673; see Supplementary Table 3). These SNPs are located in the intronic regions of *STXBP5-AS1* on chromosome 6. Using bioinformatics tools, it is predicted that changes from T to A in rs6936461 and A to G in rs9497673, has the potential to alter the splicing of the gene³⁰. *STXBP5-AS1* is a long non-coding (lncRNA) gene. lncRNAs drive many important cancer phenotypes through their interactions with other cellular macromolecules including DNA, protein, microRNA and mRNA. The different expression of lncRNAs in colorectal cancer indicate that lncRNAs are involved in all stages of colorectal cancer. In colorectal cancer pathogenesis, lncRNAs are implicated

in a variety of signaling pathways including the Wnt/-catenin signaling pathway, epidermal growth factor receptor (EGFR)/insulin-like growth factor type I receptor (IGF-IR) signaling pathway, KRAS and phosphatidylinositol-3-kinase (PI3K) pathways, transforming growth factor-beta (TGF-) signaling pathway, p53 signaling pathway, and the epithelial-mesenchymal transition (EMT) pathway³¹. While it is still unclear how *STXBP5-AS1* contributes to colon carcinogenesis, in a study involving 1061³⁰ breast cancer samples, Guo et al. identified *STXBP5-AS1* among 27²⁷ RNA genes which play a role in predicting the prognostic survival with good sensitivity and specificity. The lncRNAs⁴⁷ may act as competing endogenous RNAs (ceRNAs) and interfere in the binding of miR-190b to certain targets such as ERG, STK38L, and FNDC3A and thus contribute to breast cancer pathogenesis³². *STXBP5-AS1* may act similarly in colorectal cancer; it may hinder the binding of microRNAs to their target genes and subsequently modulate colorectal cancer tumorigenesis.

Interestingly, *STXBP5-AS1* was identified among genes that are methylated in buccal samples in a genome-wide¹⁰ screen for cigarette smoke exposure, indicating its possible role in smoking-related diseases³³. Since there is a significant difference in smoking status between³⁹ cases and controls in our cohort, it is plausible that genetic variants associated with tobacco smoke are also associated with the presence of colorectal cancer in our study population.

The polygenic model represents a strategy for jointly modeling³² SNP effects in a GWAS and development of risk prediction models in a specific population. These models can be used to estimate an individual's risk of colorectal cancer based on easily obtainable genotypes. While most of the variants flagged in the polygenic model are novel, the gene *ARHGEF3* has been implicated in promoting nasopharyngeal carcinoma in Asians³⁴. *RGMB* has been shown to promote colorectal cancer growth³⁵. Additional samples will enable us to refine and validate a polygenic colorectal cancer risk model in Indonesians.

Methods

Study participants. Indonesia is an archipelago consisting of more than 14,000 islands. There are five major islands, and one of them is Sulawesi. Makassar is located in the southern part of Sulawesi. It is considered the largest city in eastern Indonesia. 162 colorectal cancer cases were recruited from seven hospitals throughout Makassar between 2014 and 2016. The hospitals were Wahidin Sudirohusodo Hospital, Hasanuddin University Hospital, 14¹⁴ Sina Hospital, Akademis Hospital, Grestelina Hospital, Stella Maris Hospital, and Hikmah Hos-
11¹¹pital. 193 controls were frequency matched to cases on age category, sex, and ethnicity. Informed consents were obtained from all subjects, and all methods were carried out in accordance with the relevant guidelines and regulations as determined by ethical review approved by the Hasanuddin University Ethical Committee (registration number: UH 15040389).

Data and DNA sample collection. Questionnaires and medical records were recorded into study data collection forms and entered into a study database. The case forms contained 382 questions and the control forms contained 319 questions. The forms included information on demographics, cancer history in the family, smoking behavior, alcohol use, and detailed dietary history. For colorectal cancer cases, the forms collected information on cancer symptoms, staging (post operation), tumor, location, histopathology, 29²⁹ type of surgery. The questionnaire is included as a "Supplementary file". The database was managed by the Bioinformatics and Data Science Research Center (BDSRC) at Bina Nusantara University (Jakarta, Indonesia). A blood sample was collected from the basilic/cephalic vein on all participants for genotyping. These blood samples were stored in Hasanuddin University Laboratory at -20° C.

Genotyping and imputation. DNA samples were collected at the hospital where surgery was performed (Wahidin Hospital). DNA was extracted from samples at Mochtar Riady Institute for Nanotechnology (MRIN) Laboratory <https://www.overleaf.com/project/5efa1240b367400001bf3549> (Tangerang, Indonesia). Genomic DNA was extracted from 200 μ L of whole blood sample using the QIAamp DNA Mini Kit (Qiagen, Hilden, Germany) according to the manufacturer's protocol. DNA concentration was determined using NanoDrop ND-1000 spectrophotometer, version 3.3 (Thermo Fisher Scientific, Wilmington, DE, USA) and adjusted to a concentration of 20 ng/ μ L. The quality of DNA extracted was verified by purity index of OD260/OD280 (1.8–2.0) and OD260/OD230 (> 1.5). The DNA was inspected through Gel Electrophoresis using 1% molecular biology grade Agarose (Biorad, Hercules, CA, USA). Two plates of samples (92 cases and 92 controls) were allocated for this preliminary study and filled based on the DNA quality. Extracted DNA were sent to RUCDR Infinite Biologics for genotyping (Piscataway, NJ, USA) under Material Transfer Agreement (MTA) approved by the 36³⁶ Indonesian Health Ministry (registration number: LB.02.01/I/12749/2016).

DNA samples from⁴ study cases and controls were genome-wide genotyped on the Smokescreen Genotyping Array³⁶. Using 200 ng of genomic DNA, array plates were prepared using the Axiom 2.0 Reagent⁴ kits and then processed on the GeneTitan MC instrument (Thermo Fisher Scientific, Wilmington, DE, USA). Analysis of the raw data was performed using Affymetrix Power tools (APT) v⁴ 1.6 according to the Affymetrix best practices workflow. 183 samples remained after completing these steps. Additional steps were performed using SNPfilter to identify and select best performing probe sets and high quality SNPs for downstream analysis. 524,765 SNPs remained after QC filtering. Additional sample quality control included verifying concordance of study replicates, checking for unintentional duplicates and unexpected relatives, and verifying genetic versus reported gender. After filtering samples with missing covariates, 173 samples (84 cases and 89 controls) remained for statistical anal-
33³³

Genome-wide imputation was performed on the Michigan Imputation Server v1.0.2³⁷. Briefly, quality controlled study genotypes were reported on the forward strand and uploaded in vcf format. 1000 Genomes Phase 3³⁸ was selected as a reference panel, phasing was performed using Eagle v2.3³⁹, and allele frequencies were

compared against the 1000 Genomes East Asian (EAS) populations. The server automatically excludes variants with alleles other than (A, C, T, G), variants with duplicate positions, indels, monomorphic sites, and allele mismatches with the reference panel.

Statistical analysis. *Ancestry analysis.* Ancestry categories were estimated from 5515 ancestry informative markers contained in the Smokescreen Genotyping Array using fastStructure 1.0⁴⁰. Combining study and reference data from the 1000 Genomes Project Phase 3, we estimated the ancestry proportions of East Asian (EAS), South Asian (SAS), European (EUR), and African (AFR).

Genome-wide association analysis. We filtered out variants with poor imputation quality (< 0.3) and rare variants (minor allele < 1%). We then performed a marginal analysis of the remaining SNP genotype dosages fitting logistic regression models, with sex, age, body mass index, smoking status and estimated ancestries proportions (i.e., SAS, EUR, AFR) as covariates. The threshold for statistical significance in the discovery scan was set at the historical traditional genome-wide value of 5E-8. This a multi-ethnic model was implemented using glm in R⁴¹.

We queried the scan results for markers previously reported to be associated with colorectal cancer. These variants were identified through previous genotyping in an independent sample of South Sulawesi colorectal cancer cases (R. Kusuma, I. Suriapranata, personal communication) and a recent catalog of colorectal cancer SNPs for a genome-wide association scan in Hispanics⁴². The source and annotation for these variants are provided in Supplementary Table 3. Variants with evidence of replication (p-value < 0.05) were flagged for further investigation. Regional association plots were generated in LocusZoom⁴³.

We also developed a polygenic model considering the joint effect of multiple genetic variants on colorectal cancer⁴⁴. We included a screening step as a practical way to keep the number of variants under consideration in the polygenic model close to the total sample size. In this screening step the top 200 genetic associations were selected, based on Bayes factors⁴⁵, as candidate predictors in this joint model. Bayes factors were computed for the marginal versus the null models for each SNP while controlling for gender, age, BMI, and smoking status. To jointly model these variants, we use a Bayesian variable selection technique. In particular, we fit a logistic regression model utilizing shrinkage priors for each of the explanatory variables; i.e., the covariates listed above as well as the remaining candidate SNPs. In this analysis, the generalized double Pareto shrinkage prior⁴⁶ was specified and the parameters of the joint model were estimated via a maximum a posteriori (MAP) estimator⁴⁶ which was obtained via an expectation-maximization (EM) algorithm⁴⁷. The MAP estimator under these specifications simultaneously completes parameter estimation and variable selection by obtaining a sparse estimator⁴⁸; i.e., some of the regression coefficients are estimated to be identically equal to zero thus removing the effect of the corresponding explanatory variable. The EM algorithm was developed following the techniques illustrated by Armagan et al.⁴⁶ and Polson et al.⁴⁹ and the regularization parameters were selected via the Bayesian information criterion⁵⁰. These algorithms were implemented in R and completed within 90 s on an Intel based laptop, see Joyner et al.⁴⁴ for details including the source code.

Conclusions

We demonstrate replication of several colorectal cancer genetic risk factors in an Indonesian population. This study overcame the many challenges of genomic research in resource-constrained settings and provides rational for additional data collection in this population to characterize these regions more precisely and identify genetic risk factors unique to this diverse population. The primary focus of this study was replicating associations of known colorectal cancer risk variants in an Indonesian population. A secondary focus was computing genome-wide summary statistics for contributions to international colorectal cancer consortia. With additional data collections in Indonesia, we may examine and report on environmental factors (e.g., dietary factors) as well as gene-environment interactions.

Received: 29 October 2020; Accepted: 14 April 2021
Published online: 11 May 2021

References

1. Torre, L. A. et al. Global cancer statistics, 2012. *CA Cancer J. Clin.* **65**, 87–108 (2015).
2. Siegel, R. L., Miller, K. D. & Jemal, A. Cancer statistics, 2016. *CA Cancer J. Clin.* **66**, 7–30 (2016).
3. Pardamean, B., Baurley, J. W., Pardamean, C. I. & Figueiredo, J. C. Changing colorectal cancer trends in Asians. *Int. J. Colorectal Disease* **31**, 1537 (2016).
4. Pourhoseingholi, M. A. Increased burden of colorectal cancer in Asia. *World J. Gastrointest. Oncol.* **4**, 68 (2012).
5. Ng, C. J., Teo, C. H., Abdullah, N., Tan, W. P. & Tan, H. M. Relationships between cancer pattern, country income and geographical region in Asia. *BMC Cancer* **15**, 613. <https://doi.org/10.1186/s12885-015-1615-0> (2015).
6. Ferlay, J. et al. Cancer incidence and mortality worldwide: Sources, methods and major patterns in GLOBOCAN 2012. *Int. J. Cancer* **136**, E359–E386, <https://doi.org/10.1002/ijc.29210> (2015).
7. Peters, U., Bien, S. & Zubair, N. Genetic architecture of colorectal cancer. *Gut* **64**, 1623–1636 (2015).
8. Haiman, C. A. & Stram, D. O. Exploring genetic susceptibility to cancer in diverse populations. *Curr. Opin. Genet. Dev.* **20**, 330–335 (2010).
9. Jia, W.-H. et al. Genome-wide association analyses in east Asians identify new susceptibility loci for colorectal cancer. *Nat. Genet.* **45**, 191 (2013).
10. Zhang, B. et al. Large-scale genetic study in east Asians identifies six new loci associated with colorectal cancer risk. *Nat. Genet.* **46**, 533 (2014).
11. Widjaja, S. & Yo, H. RM-049Colorectal cancer in Indonesia—A centre report. *Ann. Oncol.* **27**, ii97. <https://doi.org/10.1093/annonc/mdw201.46> (2016).

12. Phipps, A. I. *et al.* Colon and rectal cancer survival by tumor location and microsatellite instability: The Colon Cancer Family Registry. *Dis. Colon Rectum* **56**, 937–944. <https://doi.org/10.1097/DCR.0b013e31828f9a57> (2013).
13. Hemminki, K. *et al.* Tumor location and patient characteristics of colon and rectal adenocarcinomas in relation to survival and TNM classes. *BMC Cancer* **10**, 688. <https://doi.org/10.1186/1471-2407-10-688> (2010).
14. Deng, Y. Rectal cancer in asian vs. western countries: Why the variation in incidence?. *Curr. Treatment Options Oncol.* **18**, 1–8 (2017).
15. Consortium, G. *et al.* The genomeasia 100k project enables genetic discoveries across asia. *Nature* **576**, 106 (2019).
16. Cui, R. *et al.* Common variant in 6q26-q27 is associated with distal colon cancer in an Asian population. *Gut* **60**, 799–805 (2011).
17. Zhu, L. *et al.* Genetic variant rs7758229 in 6q26-q27 is not associated with colorectal cancer risk in a Chinese population. *PLoS ONE* **8**, e59256 (2013).
18. Hsu, C.-M. *et al.* Upregulated SLC22A3 has a potential for improving survival of patients with head and neck squamous cell carcinoma receiving cisplatin treatment. *Oncotarget* **8**, 74348–74358 (2017).
19. Grisanzio, C. *et al.* Genetic and functional analyses implicate the NUDT11, HNF1B, and SLC22A3 genes in prostate cancer pathogenesis. *Proc. Natl. Acad. Sci. USA* **109**, 11252–11257 (2012).
20. Li, Q. & Shu, Y. Role of solute carriers in response to anticancer drugs. *Mol. Cell Ther.* **2**, 15 (2014).
21. Yokoo, S. *et al.* Significance of organic cation transporter 3 (SLC22A3) expression for the cytotoxic effect of oxaliplatin in colorectal cancer. *Drug Metab. Dispos.* **36**, 2299–2306 (2008).
22. Whiffin, N. *et al.* Identification of susceptibility loci for colorectal cancer in a genome-wide meta-analysis. *Hum. Mol. Genet.* **23**, 4729–4737 (2014).
23. Tanikawa, C. *et al.* GWAS identifies two novel colorectal cancer loci at 16q24.1 and 20q13.12. *Carcinogenesis* **39**, 652–660 (2018).
24. Schmit, S. L. *et al.* Novel common genetic susceptibility loci for colorectal cancer. *J. Natl. Cancer Inst.* **111**, 146–157. <https://doi.org/10.1093/jnci/djy099> (2019).
25. Schumacher, F. R. *et al.* Genome-wide association study of colorectal cancer identifies six new susceptibility loci. *Nat. Commun.* **6**, 7138 (2015).
26. Sneddon, J. B. *et al.* Bone morphogenetic protein antagonist gremlin 1 is widely expressed by cancer-associated stromal cells and can promote tumor cell proliferation. *Proc. Natl. Acad. Sci. USA* **103**, 14842–14847 (2006).
27. Stabile, H. *et al.* Bone morphogenic protein antagonist drm/gremlin is a novel proangiogenic factor. *Blood* **109**, 1834–1840 (2007).
28. Ziai, J. *et al.* Defining the polyposis/colorectal cancer phenotype associated with the ashkenazi GREM1 duplication: Counselling and management recommendations. *Genet. Res.* **98**, e5 (2016).
29. Davis, H. *et al.* Aberrant epithelial GREM1 expression initiates colonic tumorigenesis from cells outside the stem cell niche. *Nat. Med.* **21**, 62–70 (2015).
30. Desmet, F. O. *et al.* Human Splicing Finder: An online bioinformatics tool to predict splicing signals. *Nucleic Acids Res.* <https://doi.org/10.1093/nar/gkp215> (2009).
31. Yang, Y., Junjie, P., Sanjun, C. & Ma, Y. Long non-coding RNAs in colorectal cancer: Progression and future directions. *J. Cancer.* <https://doi.org/10.7150/jca.19794> (2017).
32. Guo, W. *et al.* Transcriptome sequencing uncovers a three-long noncoding RNA signature in predicting breast cancer survival. *Sci. Rep.* <https://doi.org/10.1038/srep27931> (2016).
33. Wan, E. S. *et al.* Smoking-associated site-specific differential methylation in buccal mucosa in the COPDGenE study. *Am. J. Respir. Cell Mol. Biol.* **53**, 246–254. <https://doi.org/10.1165/rcmb.2014-0103OC> (2015).
34. Liu, T.-H. *et al.* The putative tumor activator ARHGEF3 promotes nasopharyngeal carcinoma cell pathogenesis by inhibiting cellular apoptosis. *Oncotarget* **7**, 25836–25848 (2016).
35. Shi, Y. *et al.* Dragon (repulsive guidance molecule b, RGMb) is a novel gene that promotes colorectal cancer growth. *Oncotarget* **6**, 20540–20554 (2015).
36. Baurley, J. W., Edlund, C. K., Pardamean, C. I., Conti, D. V. & Bergen, A. W. Smokescreen: A targeted genotyping array for addiction research. *BMC Genom.* **17**, 145. <https://doi.org/10.1186/s12864-016-2495-7> (2016).
37. Das, S. *et al.* Next-generation genotype imputation service and methods. *Nat. Genet.* **48**, 1284 (2016).
38. Consortium, G. P. *et al.* A global reference for human genetic variation. *Nature* **526**, 68 (2015).
39. Loh, P. *Eagle v2.4 user manual.* (Accessed 07 May 2018).
40. Raj, A., Stephens, M. & Pritchard, J. K. faststructure: Variational inference of population structure in large snp data sets. *Genetics* **197**, 573–589 (2014).
41. R Core Team. *GLM: Fitting Generalized Linear Models* (R Foundation for Statistical Computing, 2016).
42. Schmit, S. L. *et al.* Genome-wide association study of colorectal cancer in Hispanics. *Carcinogenesis* **37**, 547–556. <https://doi.org/10.1093/carcin/bgw046> (2016).
43. Pruim, R. J. *et al.* Locuszoom: Regional visualization of genome-wide association scan results. *Bioinformatics* **26**, 2336–2337 (2010).
44. Joyner, C., McMahan, C., Baurley, J. & Pardamean, B. A two-phase Bayesian methodology for the analysis of binary phenotypes in genome-wide association studies. *Biom. J.* **62**, 191–201. <https://doi.org/10.1002/bimj.201900050> (2020).
45. Raftery, A. E. Approximate Bayes factors and accounting for model uncertainty in generalised linear models. *Biometrika* **83**, 251–266 (1996).
46. Armagan, A., Dunson, D. B. & Lee, J. Generalized double pareto shrinkage. *Stat. Sinica* **23**, 119 (2013).
47. Dempster, A. P., Laird, N. M. & Rubin, D. B. Maximum Likelihood from Incomplete Data Via the EM Algorithm. *J. Royal Stat. Soc. Ser. B (Methodological)* **39**, 1–22. <https://doi.org/10.1111/j.2517-6161.1977.tb01600.x> (1977).
48. Friedman, J., Hastie, T. & Tibshirani, R. Regularization paths for generalized linear models via coordinate descent. *J. Stat. Softw.* **33**, 1 (2010).
49. Polson, N. G. & Scott, J. G. Data augmentation for non-gaussian regression models using variance-mean mixtures. *Biometrika* **100**, 459–471 (2013).
50. Konishi, S. & Kitagawa, G. Bayesian Information Criteria. 211–237. https://doi.org/10.1007/978-0-387-71887-3_9 (Springer, New York, NY, 2008).
51. Suryapranata, I. & Kusuma, R. (N.D.). Unpublished.
52. Peters, U. *et al.* Identification of genetic susceptibility loci for colorectal tumors in a genome-wide meta-analysis. *Gastroenterology* **144**, 799–807.e24. <https://doi.org/10.1053/j.gastro.2012.12.020> (2013).
53. Whiffin, N. *et al.* Identification of susceptibility loci for colorectal cancer in a genome-wide meta-analysis. *Hum. Mol. Genet.* **23**, 4729–4737. <https://doi.org/10.1093/hmg/ddu177> (2014).
54. Houlston, R. S. *et al.* Meta-analysis of three genome-wide association studies identifies susceptibility loci for colorectal cancer at 1q41, 3q26.2, 12q13.13 and 20q13.33. *Nat. Genet.* **42**, 973–977. <https://doi.org/10.1038/ng.670> (2010).
55. Schumacher, F. R. *et al.* Genome-wide association study of colorectal cancer identifies six new susceptibility loci. *Nat. Commun.* **6**, 7138. <https://doi.org/10.1038/ncomms8138> (2015).
56. Real, L. M. *et al.* A colorectal cancer susceptibility new variant at 4q26 in the Spanish population identified by genome-wide association analysis. *PLoS ONE* **9**, e101178. <https://doi.org/10.1371/journal.pone.0101178> (2014).
57. Dunlop, M. G. *et al.* Common variation near CDKN1A, POLD3 and SHROOM2 influences colorectal cancer risk. *Nat. Genet.* **44**, 770–776. <https://doi.org/10.1038/ng.2293> (2012).

58. Cui, R. *et al.* Common variant in 6q26-q27 is associated with distal colon cancer in an Asian population. *Gut* **60**, 799–805. <https://doi.org/10.1136/gut.2010.215947> (2011).
59. Zanke, B. W. *et al.* Genome-wide association scan identifies a colorectal cancer susceptibility locus on chromosome 8q24. *Nat. Genet.* **39**, 989–994. <https://doi.org/10.1038/ng2089> (2007).
60. Gruber, S. B. *et al.* Genetic variation in 8q24 associated with risk of colorectal cancer. *Cancer Biol. Ther.* **6**, 1143–1147 (2007).
61. Haiman, C. A. *et al.* A common genetic risk factor for colorectal and prostate cancer. *Nat. Genet.* **39**, 954–956. <https://doi.org/10.1038/ng2098> (2007) (NIHMS150003).
62. Tomlinson, I. P. *et al.* A genome-wide association study identifies colorectal cancer susceptibility loci on chromosomes 10p14 and 8q23.3. *Nat. Genet.* **40**, 623–630. <https://doi.org/10.1038/ng.111> (2008).
63. Hutter, C. M. *et al.* Characterization of the association between 8q24 and colon cancer: Gene–environment exploration and meta-analysis. *BMC Cancer* **10**, 670. <https://doi.org/10.1186/1471-2407-10-670> (2010).
64. Tenesa, A. *et al.* Genome-wide association scan identifies a colorectal cancer susceptibility locus on 11q23 and replicates risk loci at 8q24 and 18q21. *Nat. Genet.* **40**, 631–637. <https://doi.org/10.1038/ng.133> (2008) (NIHMS150003).
65. Wang, H. *et al.* Fine-mapping of genome-wide association study-identified risk loci for colorectal cancer in African Americans. *Hum. Mol. Genet.* **22**, 5048–5055. <https://doi.org/10.1093/hmg/ddt337> (2013).
66. Jaeger, E. *et al.* Common genetic variants at the CRAC1 (HMPS) locus on chromosome 15q13.3 influence colorectal cancer risk. *Nat. Genet.* **40**, 26–28. <https://doi.org/10.1038/ng.2007.41> (2008).
67. Broderick, P. *et al.* A genome-wide association study shows that common alleles of SMAD7 influence colorectal cancer risk. *Nat. Genet.* **39**, 1315–1317. <https://doi.org/10.1038/ng.2007.18> (2007).

Acknowledgements

We would like to acknowledge Bina Nusantara and Hasanuddin University for funding this study, MRIN Laboratory for DNA Extraction, RUCDR Infinite Biologics for DNA processing and genotyping, BioRealm for support of the Smokescreen Genotyping Array, Research credits from Amazon Web Services (AWS) and generous contributions from NVIDIA and the AI R&D Center at Bina Nusantara University for computing and database support.

Author contributions

Conceptualization, I.Y., U.M., R.L., G.M., B.P., and J.B.; methodology, J.B., M.K., A.B., C.M., and C.J.; software, M.K., A.B., T.C., C.M., and C.J.; validation, B. 43 C.P., C.M., and J.B.; formal analysis, A.B., C.M. 17 and C.J.; investigation, I.Y., U.M., R.L., G.M., I.S., B.P., A.B., and J.B.; data curation, A.I., A.A., R.K., and A.B.; writing—original draft preparation, I.S., R.K., A.B., T.C., C.M., C.J., and J.B.; writing—review and editing, I.Y., I.S., B.P., C.P., and J.B.; visualization, A.B.; supervision, I.Y., U. 6 R.L., and B.P.; project administration, A.I.; funding acquisition, I.Y., U.M., and B.P. All authors reviewed the manuscript.

Competing interests

The authors declare no competing interests.


Additional information

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41598-021-88805-4>.

Correspondence and requests for materials should be addressed to B.P. 18 or J.W.B.

Reprints and permissions information is available at 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) 2021

Genetic risk factors for colorectal cancer in multiethnic Indonesians

ORIGINALITY REPORT

19%

SIMILARITY INDEX

16%

INTERNET SOURCES

17%

PUBLICATIONS

8%

STUDENT PAPERS

PRIMARY SOURCES

1	bmcoralhealth.biomedcentral.com Internet Source	3%
2	www.readkong.com Internet Source	2%
3	jcancer.org Internet Source	1%
4	www.medrxiv.org Internet Source	1%
5	academic.oup.com Internet Source	1%
6	publikationen.uni-tuebingen.de Internet Source	1%
7	bmccancer.biomedcentral.com Internet Source	1%
8	www.researchsquare.com Internet Source	1%
9	Submitted to University of Brighton Student Paper	1%

10	"Hereditary Colorectal Cancer", Springer Science and Business Media LLC, 2018 Publication	<1 %
11	bmcmedgenomics.biomedcentral.com Internet Source	<1 %
12	www.newswise.com Internet Source	<1 %
13	Adam M. Schmitt, Howard Y. Chang. "Long Noncoding RNAs in Cancer Pathways", Cancer Cell, 2016 Publication	<1 %
14	Ulrike Peters, Shuo Jiao, Fredrick R. Schumacher, Carolyn M. Hutter et al. "Identification of Genetic Susceptibility Loci for Colorectal Tumors in a Genome-Wide Meta-analysis", Gastroenterology, 2013 Publication	<1 %
15	Submitted to Universitas Negeri Jakarta Student Paper	<1 %
16	Burnett-Hartman, A. N., P. A. Newcomb, C. M. Hutter, U. Peters, M. N. Passarelli, M. R. Schwartz, M. P. Upton, L.-C. Zhu, J. D. Potter, and K. W. Makar. "Variation in the Association Between Colorectal Cancer Susceptibility Loci and Colorectal Polyps by Polyp Type", American Journal of Epidemiology, 2014. Publication	<1 %

17 Inkyu Rhee, Jun-Seok Lee, Jong Kim, Jong-Ho Kim. "Nitrogen Oxides Mitigation Efficiency of Cementitious Materials Incorporated with TiO₂", Materials, 2018 <1 %
Publication

18 Jiayao Zhou, Yumeng Wang, Gaoxingyu Huang, Min Yang, Yumin Zhu, Chen Jin, Dan Jing, Kai Ji, Yigong Shi. "LilrB3 is a putative cell surface receptor of APOE4", Cell Research, 2023 <1 %
Publication

19 Marion Dhooge, Stéphanie Baert-Desurmont, Carole Corsini, Olivier Caron et al. "National recommendations of the French Genetics and Cancer Group - Unicancer on the modalities of multi-genes panel analyses in hereditary predispositions to tumors of the digestive tract", European Journal of Medical Genetics, 2020 <1 %
Publication

20 Submitted to University of Melbourne <1 %
Student Paper

21 downloads.hindawi.com <1 %
Internet Source

22 www.nature.com <1 %
Internet Source

23 www.sciencepub.net

<1 %

24

Anna Díez-Villanueva, Mireia Jordà, Robert Carreras-Torres, Henar Alonso et al.
"Identifying causal models between genetically regulated methylation patterns and gene expression in healthy colon tissue",
Clinical Epigenetics, 2021

Publication

<1 %

25

Submitted to Binus University International

Student Paper

<1 %

26

scirp.org

Internet Source

<1 %

27

www.karger.com

Internet Source

<1 %

28

Su, Zhan, Laura J Gay, Amy Strange, Claire Palles, Gavin Band, David C Whiteman, Francesco Lescai, Cordelia Langford, Manoj Nanji, Sarah Edkins, Anouk van der Winkel, David Levine, Peter Sasieni, Céline Bellenguez, Kimberley Howarth, Colin Freeman, Nigel Trudgill, Art T Tucker, Matti Pirinen, Maikel P Peppelenbosch, Luc J W van der Laan, Ernst J Kuipers, Joost P H Drenth, Wilbert H Peters, John V Reynolds, Dermot P Kelleher, Ross McManus, Heike Grabsch, Hans Prenen, Raf Bisschops, Kausila Krishnadath, Peter D

<1 %

Siersema, Jantine W P M van Baal, Mark Middleton, Russell Petty, Richard Gillies, Nicola Burch, Pradeep Bhandari, Stuart Paterson, Cathryn Edwards, Ian Penman, Kishor Vaidya, Yeng Ang, Iain Murray, Praful Patel, Weimin Ye, Paul Mullins, Anna H Wu, Nigel C Bird, Helen Dallal, Nicholas J Shaheen, Liam J Murray, Konrad Koss, Leslie Bernstein, Yvonne Romero, Laura J Hardie, Rui Zhang, Helen Winter, Douglas A Corley, Simon Panter, Harvey A Risch, Brian J Reid, Ian Sargeant, Marilie D Gammon, Howard Smart, An. "Common variants at the MHC locus and at chromosome 16q24.1 predispose to Barrett's esophagus", Nature Genetics, 2012.

Publication

29

conference.binus.ac.id

Internet Source

<1 %

30

cyberleninka.org

Internet Source

<1 %

31

www.scirp.org

Internet Source

<1 %

32

Korbinian Weigl, Hauke Thomsen, Yesilda Balavarca, Jacklyn N. Hellwege, Martha J. Shrubsole, Hermann Brenner. "Genetic Risk Score Is Associated With Prevalence of Advanced Neoplasms in a Colorectal Cancer

<1 %

Screening Population", Gastroenterology, 2018

Publication

33

Rebecca Sims, Sven J van der Lee, Adam C Naj, Céline Bellenguez et al. "Rare coding variants in PLCG2, ABI3, and TREM2 implicate microglial-mediated innate immunity in Alzheimer's disease", Nature Genetics, 2017

Publication

<1 %

34

Yusuke Takahashi, Keishi Sugimachi, Ken Yamamoto, Atsushi Niida et al. "Japanese genome-wide association study identifies a significant colorectal cancer susceptibility locus at chromosome 10p14", Cancer Science, 2017

Publication

<1 %

35

bmcoophthalmol.biomedcentral.com

Internet Source

<1 %

36

dspace.cuni.cz

Internet Source

<1 %

37

genepi.qimr.edu.au

Internet Source

<1 %

38

journals.plos.org

Internet Source

<1 %

39

onlinelibrary.wiley.com

Internet Source

<1 %

- | | | |
|----|---|------|
| 40 | www.mdpi.com
Internet Source | <1 % |
| 41 | www.oncotarget.com
Internet Source | <1 % |
| 42 | www.spandidos-publications.com
Internet Source | <1 % |
| 43 | "Big Data Analytics in Genomics", Springer Science and Business Media LLC, 2016
Publication | <1 % |
| 44 | Cheng-Ming Hsu, Pai-Mei Lin, Jan-Gowth Chang, Hsin-Ching Lin, Shau-Hsuan Li, Sheng-Fung Lin, Ming-Yu Yang. "Upregulated SLC22A3 has a potential for improving survival of patients with head and neck squamous cell carcinoma receiving cisplatin treatment", Oncotarget, 2017
Publication | <1 % |
| 45 | Nan Song, Aesun Shin, Ji Won Park, Jeongseon Kim, Jae Hwan Oh. "Common risk variants for colorectal cancer: an evaluation of associations with age at cancer onset", Scientific Reports, 2017
Publication | <1 % |
| 46 | Danielle B. McKenna, Jeroen Van Den Akker, Alicia Y. Zhou, Lauren Ryan et al. "Identification of a novel GREM1 duplication in | <1 % |

a patient with multiple colon polyps", Familial Cancer, 2018

Publication

47

Wenna Guo, Qiang Wang, Yueping Zhan, Xijia Chen, Qi Yu, Jiawei Zhang, Yi Wang, Xin-jian Xu, Liucun Zhu. "Transcriptome sequencing uncovers a three-long noncoding RNA signature in predicting breast cancer survival", Scientific Reports, 2016

Publication

<1 %

Exclude quotes On

Exclude matches < 5 words

Exclude bibliography On