

Accepted Manuscript

Validation of Recently Proposed Colorectal Cancer Susceptibility Gene Variants in an Analysis of Families and Patients—a Systematic Review

Peter Broderick, Sara E. Dobbins, Daniel Chubb, Ben Kinnersley, Malcolm G. Dunlop, Ian Tomlinson, Richard S. Houlston

PII: S0016-5085(16)35138-1
DOI: [10.1053/j.gastro.2016.09.041](https://doi.org/10.1053/j.gastro.2016.09.041)
Reference: YGAST 60735

To appear in: *Gastroenterology*
Accepted Date: 27 September 2016

Please cite this article as: Broderick P, Dobbins SE, Chubb D, Kinnersley B, Dunlop MG, Tomlinson I, Houlston RS, Validation of Recently Proposed Colorectal Cancer Susceptibility Gene Variants in an Analysis of Families and Patients—a Systematic Review, *Gastroenterology* (2016), doi: 10.1053/j.gastro.2016.09.041.

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.



Validation of Recently Proposed Colorectal Cancer Susceptibility Gene Variants in an Analysis of Families and Patients—a Systematic Review

Peter Broderick^{1*}, Sara E Dobbins^{1*}, Daniel Chubb^{1*}, Ben Kinnersley¹, Malcolm G Dunlop², Ian Tomlinson³, Richard S Houlston^{1,4}

1. Division of Genetics and Epidemiology, The Institute of Cancer Research, London SM2 5NG, UK.
2. Colon Cancer Genetics Group, Institute of Genetics and Molecular Medicine, Western General Hospital, Edinburgh EH4 2XU
3. Molecular and Population Genetics Laboratory and NIHR Biomedical Research Centre, Oxford Centre for Cancer Gene Research, Wellcome Trust Centre for Human Genetics, University of Oxford, Oxford OX3 7BN, UK.
4. Division of Pathology, The Institute of Cancer Research, London SM2 5NG, UK

* These authors contributed equally to this work.

Correspondence to: Richard S Houlston, E-mail: richard.houlston@icr.ac.uk; Tel: +44 (0) 208 722 4175; Fax: +44 (0) 722 4365.

CONFLICT OF INTEREST

Peter Broderick, Sara E Dobbins, Daniel Chubb, Ben Kinnersley, Malcolm G Dunlop, Ian Tomlinson, Richard S Houlston: None to declare

ACKNOWLEDGEMENTS

This work was supported by Cancer Research UK Research (C1298/A8362, Bobby Moore Fund for Cancer Research UK) and the European Union (FP7/2007-2013) under Grant No. 258236, FP7 collaborative project SYSCOL. D.C. was funded by a grant from Bloodwise. Additional support was

provided by the National Cancer Research Network and the National Health Service (NHS). In Oxford, the work was funded by the Oxford Comprehensive Biomedical Research Centre core infrastructure support to the Wellcome Trust Centre for Human Genetics, Oxford (Wellcome Trust 090532/Z/09/Z). In Scotland, the work was funded by a Cancer Research UK (C348/A12076) and Medical Research Council Grant (MR/KO18647/1). This study makes use of the ICR1000 UK exome series data generated by Professor Nazneen Rahman's team at The Institute of Cancer Research, London. This work made use of samples generated by the 1958 Birth Cohort. Access to these resources was enabled via the 58READIE Project funded by Wellcome Trust and Medical Research Council (grant numbers WT095219MA and G1001799). This publication is supported by COST Action BM1206.

AUTHOR CONTRIBUTIONS

Conception and design: Peter Broderick, Daniel Chubb, Sara E Dobbins, Richard S. Houlston

Collection and assembly of data: Peter Broderick, Daniel Chubb, Sara E. Dobbins, Malcolm G Dunlop, Ian Tomlinson, Richard S. Houlston

Data analysis and interpretation: Peter Broderick, Daniel Chubb, Ben Kinnersley, Sara E. Dobbins, Richard S. Houlston

Manuscript writing: All authors

Final approval of manuscript: All authors

ABSTRACT

High-throughput sequencing analysis has accelerated searches for genes associated with risk for colorectal cancer (CRC); germline mutations in *NTHL1*, *RPS20*, *FANCM*, *FAN1*, *TP53*, *BUB1*, *BUB3*, *LRP6*, and *PTPN12* have been recently proposed to increase CRC risk. We attempted to validate the association between variants in these genes and development of CRC in a systematic review of 11 publications, using sequence data from 863 familial CRC cases and 1604 individuals without CRC (controls). All cases were diagnosed at an age of 55 years or younger and did not carry mutations in an established CRC predisposition gene. We found sufficient evidence for *NTHL1* to be considered a CRC predisposition gene—members of 3 unrelated Dutch families were homozygous for inactivating p.Gln90Ter mutations; a Canadian woman with polyposis, CRC, and multiple tumors was reported to be heterozygous for the inactivating *NTHL1* p.Gln90Ter/c.709+1G>A mutations; and a man with polyposis was reported to carry p.Gln90Ter/p.Gln287Ter; whereas no inactivating homozygous or compound heterozygous mutations were detected in controls. Variants that disrupted *RPS20* were detected in a Finnish family with early-onset CRC (p.Val50SerfsTer23), a 39-year old individual with metachronous CRC (p.Leu61GlufsTer11 mutation), and a 41-year-old individual with CRC (missense p.Val54Leu), but not in controls. We therefore found published evidence to support the association between variants in *NTHL1* and *RPS20* with CRC, but not of other recently reported CRC susceptibility variants. We urge the research community to adopt rigorous statistical and biological approaches coupled with independent replication before making claims of pathogenicity.

KEYWORDS

Colon cancer, inherited, Germline, Exome Sequencing

ARTICLE

Understanding the genetics of familial CRC is clinically important to discriminate between high- and low-risk groups. Mutations in eleven genes are well-established to confer significant increases in CRC risk and testing for these is common in clinical practice. Despite this in many CRC families no genetic diagnosis can be made. While the availability of high-throughput-sequencing has

accelerated searches for new CRC genes there are challenges in assigning pathogenicity to identified variants.

Here we reviewed the data supporting recent assertions that *NTHL1*, *RPS20*, *FANCM*, *FAN1*, *TP53*, *BUB1*, *BUB3*, *LRP6*, and *PTPN12* are CRC susceptibility genes using an evidence-based framework (Supplementary-Material)¹⁻⁷. To search for independent evidence of a role in CRC risk we analyzed sequencing data on 863 familial CRC cases and 1,604 controls⁸. All cases were diagnosed aged ≤ 55 and were mutation-negative for known CRC genes.

Evidence for variation in *NTHL1*, which like *MUTYH* performs base-excision-repair (BER), as a cause of recessive-CRC has been provided by three unrelated Dutch families homozygous for the rare inactivating p.Gln90Ter mutation (Supplementary-Material, Supplementary-Table 1)⁶. The tumor mutation spectrum was enriched for C>T transitions, consistent with defective BER. Subsequently compound heterozygosity for inactivating *NTHL1* p.Gln90Ter/c.709+1G>A mutations was identified in a Canadian woman diagnosed with polyposis, CRC and multiple tumors⁹. Tumors were again enriched for somatic C>T transitions. While we found no p.Gln90Ter homozygotes amongst our WES cases, a 41-year old male case with co-incident polyposis harbored p.Gln90Ter/p.Gln287Ter. No inactivating homozygotes or compound heterozygotes were seen among our 1,604 controls.

Whole-exome sequencing (WES) of a Finnish Amsterdam-positive family demonstrated significant segregation of *RPS20* p.Val50SerfsTer23 with early-onset CRC (LOD score=3.0; Supplementary-Material, Supplementary-Table 1)³. No disruptive *RPS20* variants have been catalogued by the Exome-Aggregation-Consortium (ExAC), which contains WES data for 60,706 individuals of diverse ancestries¹⁰ suggesting the gene is intolerant to mutation. Hence, it is notable that in our WES series we identified the disruptive p.Leu61Glu fsTer11 mutation in a 39-year old with metachronous CRC. Furthermore we identified the deleterious missense p.Val54Leu in an Amsterdam-positive 41-year old case. No rare missense/disruptive mutations identified in the 1,604 controls.

Smith *et al.* identified *FANCM* p.Arg1931Ter in two sporadic CRC cases with cancers showing loss of the wild-type allele (LOH)⁵. p.Arg1931Ter has been shown to induce exon skipping resulting in

decreased DNA-repair (Supplementary-Material, Supplementary-Table 1). In our WES series we detected p.Arg1931Ter in four cases and one control ($P=0.02$; Supplementary-Table 3). To seek further evidence for an association between p.Arg1931Ter and CRC, we investigated the frequency of this specific variant in two additional UK series totaling 5,552 cases and 6,792 population controls (published Illumina-Exome-BeadChip data¹¹; Supplementary-Material). Combining these data provided no evidence for an association (Meta-analysis $P=0.22$; Supplementary Figure 1).

FAN1 mutations have been reported as a cause of CRC in Amsterdam-positive families⁴, but evidence for segregation was weak ($P=0.125$) and the evidence for any functional effect of mutation was only shown in non-colonic tissue (Supplementary-Material, Supplementary-Table 1). In our WES series we found no significant increase in the burden of *FAN1* mutations in cases (Table 1; Supplementary-Tables 2&3).

Germline mutation of *TP53*, archetypically associated with Li-Fraumeni syndrome, has recently been suggested to cause familial CRC at a frequency comparable to *APC*⁷. The assertion was, however, based on the flawed assumption that all rare missense changes seen were disease-causing with no consideration of mutation burden in controls (Supplementary-Material, Supplementary-Table 1). In our data no over-representation of *TP53* mutation was seen in cases (Table 1, Supplementary-Tables 2&3).

By WES small numbers of early-onset CRC, *BUB1*, *BUB3*, *LRP6* and *PTPN12* have been proposed as CRC predisposition genes^{1,2}. The published evidence to support assertions is minimal (Supplementary-Material, Supplementary-Table 1) with no evidence of segregation or LOH. Moreover, of the two *BUB1* mutation carriers, one also carried a *MLH1* mutation which, unlike *BUB1*, segregated with colorectal tumors. Only for *PTPN12* did the authors demonstrate an increase in the burden of mutation in cases versus controls ($P=0.039$; Supplementary-Material). While we also observed an enrichment of missense *PTPN12* mutation in our WES cases ($P=0.039$; Table 1, Supplementary-Table 3), in light of the number of genes investigated, the evidence for a role in CRC predisposition remains weak.

In conclusion a role for *NTHL1* as a *bona fide* CRC gene is supported by multiple lines of evidence. While compelling, the assertion that mutation of *RPS20* causes CRC remains to be established as this observation is based on a single family and the mechanism by which ribosomal proteins might

predispose to CRC is unclear. In contrast, evidence to support other genes as risk factors is currently lacking.

Investigators must remember that private variants are common; of the 7,404,909 variants listed in ExAC, 54% are observed only once¹⁰, therefore novel variants should be considered benign until proved otherwise. A studies power to detect a statistically significant association with any rare variant is typically weak, therefore additional evidence must be considered including segregation of the genotype with disease in families, somatic mutation and functional studies with relevance to CRC biology. Critically, where multiple variants are considered within a gene, the burden of variation within controls must also be considered. Since the frequency of variants can be highly population-specific it is essential that controls used for comparison are well matched.

While there is a strong rationale for seeking to identify new CRC genes, well powered studies are required to mitigate against erroneous findings being asserted as causative and subsequently included in databases from which they are seldom deleted. The WES data we have generated represents the largest cohort of CRC exomes sequenced to date. The use of this dataset, which is publically available, to validate observations from small sequencing studies should act to limit the reporting of false positive results. Finally, the evidence framework we have implemented to assess the validity of proposed CRC genes, provides a robust strategy for establishing clinically actionable genes.

TABLES AND FIGURES

Table 1: Gene Burden analysis. Number of cases (n=863) and controls (n=1,604) with rare (MAF<1%) mutations in postulated CRC genes. *P*-values calculated using Fishers exact test, *P*-values <0.05 are emboldened.

REFERENCES

1. **de Voer RM, Geurts van Kessel A, Weren RD, Ligtenberg MJ**, et al. *Gastroenterology* 2013;145:544-7.
2. **de Voer RM, Hahn MM**, et al. *PLoS Genet* 2016;12:e1005880.
3. Nieminen TT, et al. *Gastroenterology* 2014;147:595-598 e5.
4. **Segui N, Mina LB**, et al. *Gastroenterology* 2015;149:563-6.
5. Smith CG, et al. *Hum Mutat* 2013;34:1026-34.
6. Weren RD, et al. *Nat Genet* 2015;47:668-71.
7. Yurgelun MB, et al. *JAMA Oncol* 2015;1:214-21.
8. **Chubb D, Broderick P, Dobbins SE**, et al. *Nat Commun* 2016;7:11883.
9. Rivera B, et al. *N Engl J Med* 2015;373:1985-6.
10. Lek M, et al. *Nature* 2016;536:285-91.
11. **Timofeeva MN, Kinnersley B**, et al. *Sci Rep* 2015;5:16286.

Author names in bold designate shared co-first authorship

Table 1: Gene Burden analysis. Number of cases (n=863) and controls (n=1,604) with rare (MAF<1%) mutations in postulated CRC genes. *P*-values calculated using Fishers exact test, *P*-values <0.05 are emboldened.

Gene	Previously Reported	Disruptive mutations (stop-gain, frameshift)			Damaging mutations (disruptive, predicted-damaging, splice acceptor/donors)			All coding non-synonymous variants		
		Cases	Control	<i>P</i> _{Fisher}	Cases	Control	<i>P</i> _{Fisher}	Cases	Control	<i>P</i> _{Fisher}
<i>BUB1</i>	Disruptive	0	4	0.31	1	8	0.17	18	30	0.76
<i>BUB3</i>	Missense	0	2	0.55	0	4	0.31	1	5	0.67
<i>FAN1</i>	Disruptive /Missense	0	2	0.55	15	17	0.19	32	45 [#]	0.23
<i>FANCM</i>	Disruptive /Missense	5	1	0.02	23	33	0.33	51 [§]	67 [§]	0.06
<i>LRP6</i> (BPD*)	Missense	0	0	-	6 (4)	17 (13)	0.51 (0.45)	17 (8)	37 (21)	0.67
<i>PTPN12</i>	Missense	0	1	1.00	6	5	0.21	12	9	0.04
<i>RPS20</i>	Disruptive	1	0	0.35	2	0	0.12	2	0	0.12
<i>TP53</i>	Missense	1	0	0.35	1	1	1.00	1	4	0.66

* Number of variants within β -Propellor domain. All 3 variants identified by de Voer *et al* were within BPD.

Total number of variants in controls = 46; 1 sample has 2 *FAN1* missense

§ Totals number of variants in cases = 52, in controls =69; 3 samples have 2 *FANCM* missense

SUPPLEMENTARY METHODS AND MATERIALS

METHODS

INDEPENDENT EVALUATION

Whole-exome sequencing data: To search for independent evidence and to contextualize the impact of each purported CRC gene we made use of recently published whole-exome sequencing (WES) data on 1,006 early-onset familial CRC cases and 1,609 healthy controls¹. Cases were of European Ancestry recruited to the UK National Study of Colorectal Cancer Genetics (NSCCG)². All cases were diagnosed with CRC aged ≤ 55 and had at least one first-degree relative diagnosed with CRC. Controls were individuals with no history of malignancy selected from the 1958 Birth Cohort (1958BC), a longitudinal study following the lives of people born in England, Scotland and Wales during the week of 3-9 March 1958³. Full details of WES have been published previously. Briefly, paired end fastq files were aligned to build 37 (hg19) of the human reference genome and alignments were processed using the Genome Analysis Tool Kit (GATKv3) pipeline according to best practices⁴. The Variant Effect Predictor⁵ was used to provide annotations on the predicted impact of each variant together with functional classifications and assessment of deleteriousness from the CONDEL⁶ algorithm. Samples (cases and controls) with a variant in an established high-penetrance CRC gene which was predicted to be disruptive (stop-gain, frameshift) or previously catalogued as being pathogenic or likely-pathogenic by InSiGHT (The International Society for Gastrointestinal Hereditary Tumours) were removed. Specifically: *APC*: 19 cases, 1 control; *MLH1*: 46 cases; *MSH2*: 46 cases; *MSH6*: 13 cases, 1 control; *MUTYH*: 9 cases; *PMS2*: 6 cases, 2 controls; *POLD1*: 1 case, 1 control; *POLE*: 3 cases; *BMPR1A*, *SMAD4*, *STK11*: 0 cases. Thus for the analysis presented in this manuscript we made use of whole-exome sequencing data on 863 cases and 1,604 controls.

Gene Burden Analysis: With currently attainable sample sizes, a studies power to detect a statistically significant association with any rare variant is typically weak. Here we use WES data described above, to look for an enrichment of variation in cases versus controls, for each postulated CRC gene as a whole. As the power of such comparison depends critically on the ability to distinguish between pathogenic and non-pathogenic variation, we defined and compared a number of variant classes: (1) Disruptive mutations (stop-gain, frameshift); (2) Disruptive and predicted damaging mutations (stop-gain, frameshift, missense predicted to be damaging by CONDEL, splice site acceptor/donors); (3) All coding non-synonymous variants. We assessed rare (minor allele frequency [MAF] <1%) and very rare (MAF<0.1%) mutations in each variant class. Comparisons were made using (a) all 863 cases (b) 159 cases with Amsterdam-II positive family histories (Amsterdam-I n=146). Thus in total we performed 12 comparisons for each gene.

Further analysis of the recurrent *FANCM* p.Arg1931Ter (rs144567652) variant: We studied the association of the recurrent variant *FANCM* p.Arg1931Ter (rs144567652) with CRC by analyzing published Illumina Infinium Human Exome BeadChip 12v1.0 or 12v1.1 exon array data⁷. Specifically, we made use of UK case/control data (excluding samples also included in our WES data) comprising: (i) 3,537 English CRC cases and 4,811 control patients; (ii) 2,015 Scottish CRC cases and 1,981 Scottish controls.

EVIDENCE FRAMEWORK

To assess the validity of purported CRC genes, accounting for varying study design, we collated the following evidence where appropriate (Supplementary Table 1):

- (1) Where gene/variants were identified through the analysis of multiple members of a single family we evaluated the strength of segregation data – co-inheritance of the mutation with affection status (CRC or polyps) in the family. If not formally quantified in the published report we calculated non-parametric linkage (NPL) statistic P -values⁸ using the family information provided.
- (2) Where a specific CRC risk variant was reported: we looked for reported evidence of a statistically significant enrichment in CRC cases versus controls. In conjunction with

our WES data we examined frequency data on the mutation in an ethnically appropriate subset of the Exome-Aggregation-Consortium (ExAC) database⁹; a catalog of exome sequencing data for 60,706 individuals of diverse ancestries (non-Finnish European (NFE) 33,370 exomes, East Asian (EAS) 4,327 exomes, Finnish (FIN) 3,307 exomes).

- (3) Where numerous variants are identified in a specific gene: we looked for evidence of gene burden testing in cases versus controls, and if performed, evidence of statistically significant enrichment of mutation in cases.
- (4) Where recessive inheritance was suspected or indicated: homozygosity or compound-heterozygosity for pathogenic mutations in the proposed CRC gene was assessed in cases and controls.
- (5) Computational data on the presumptive effect of the variant.
- (6) Functional data – demonstration that mutation has a functional effect and the relevance to CRC biology.
- (7) Other information – evidence of a highly-specific phenotype associated with CRC, evidence of somatic mutation of the wild-type allele in cancers from carriers consistent with tumor suppressor gene function.

REVIEW OF EXISTING LITERATURE FOR EACH GENE ASSESSED

NTHL1

Evidence for variation in *NTHL1* as a cause of recessive CRC has been provided by three unrelated Dutch families homozygous for the rare (ExAC NFE MAF=0.0023, homozygosity ~1/75,000) inactivating p.Gln90Ter mutation¹⁰. Multiple colorectal adenomas with or without CRC were diagnosed in all seven homozygotes. The tumor mutation spectrum was significantly enriched for C:G>T:A transitions, consistent with the mutation spectra observed in *NTHL1* double-knockout mice. Subsequent to this report, compound heterozygosity for inactivating *NTHL1* p.Gln90Ter/c.709+1G>A mutations was identified in a 41-year old Canadian woman diagnosed with polyposis, CRC and multiple tumors¹¹. Tumors were again enriched for somatic C:G>T:A transitions.

Evidence Summary: Multiple reports associating homozygosity and compound-heterozygosity with CRC. Evidence of functional effect in CRC.

RPS20

In seven affected members (average age 52, range 24-75) of a four-generation Finnish Amsterdam-positive FCCTX family, Nieminen *et al.* identified a heterozygous 1-bp duplication, resulting in a frameshift and premature termination (p.Val50SerfsTer23), in *RPS20*¹². The mutation, identified through genetic linkage analysis and WES, showed full cosegregation with microsatellite-stable early-onset CRC thus providing statistically significant evidence (reported LOD score=3.0; calculated NPL=5.35, $P=0.0078$) for germline mutation in *RPS20* as a cause of CRC. The mutation was absent in 292 population controls and is not reported in the ExAC database, which includes 3,307 Finnish individuals. Tumors from mutation carriers did not show loss of the wildtype allele (LOH^{WT}). Lymphoblastoid cells (LCLs) from cases carrying p.Val50SerfsTer23 mutation showed a marked increase in 21S pre-rRNAs compared to controls (P -value not calculated), consistent with a late pre-rRNA processing defect and suggestive of *RPS20* haploinsufficiency. Germline *RPS20* mutations were not found in 25 additional Finnish FCCTX families, 292 population controls or in tumor DNA from 50 primary CRC and 11 CRC cell lines.

Evidence Summary: Statistically significant evidence of segregation, absence of gene mutation in controls. Functional evidence in non-colon tissue. No evidence of somatic mutation or functional effect in CRC.

FANCM

By searching within tumorigenesis genes for rare/novel truncating mutations, which also showed LOH^{WT} within the tumor, Smith *et al.* identified the *FANCM* mutation p.Arg1931Ter (rs144567652) in one of 50 sporadic UK CRC cases¹³. As *FANCM* is functionally linked to *MSH2/MSH6* they sought further evidence for the role of this variant by genotyping an additional case-control series identifying the mutation in 1 of 2,207 CRC cases and 1 of 2,176 controls. The tumor of the additional case again showed LOH^{WT}. Combining discovery and replication samples showed no significant enrichment of the mutation in CRC (cases

2/2,257, controls 1/2,176, $P=0.57$). Smith *et al.* presented no segregation or functional data. However, a subsequent report by Peterlongo *et al.* proposing p.Arg1931Ter as a familial breast cancer risk factor showed that p.Arg1931Ter induces exon skipping resulting in decreased DNA repair activity in mouse embryonic fibroblast cells¹⁴.

Evidence Summary: Loss of wild-type allele in tumors. Functional evidence in non-colon tissue. No evidence of segregation or functional effect in CRC. No significant enrichment of the mutation in CRC.

FAN1

Through WES of three individuals from a Spanish MMR-proficient, Amsterdam-positive CRC family Sequi *et al.* identified a novel *FAN1* truncating mutation p.Cys47Ter¹⁵. Evidence of segregation with CRC was limited (calculated NPL=0.95, $P=0.25$). Tumors developed by p.Cys47Ter mutation carrier showed no reduction in the expression of wild-type RNA or *FAN1* protein. Screening an additional 247 Spanish Amsterdam/Bethesda positive cases for rare (MAF<0.01; dbSNP135) variants identified an additional truncating mutation (p.Arg952Ter) and three missense mutations (p.Asp140Thr and p.Arg591Trp - predicted to be damaging by SIFT and CONDEL algorithms; p.Pro340Ser - predicted to be benign). No *FAN1* mutations were identified in 250 population individuals without CRC. Whilst suggestive of an enrichment in cases for the overall burden of variation in *FAN1*, the size of the control population is insufficient to provide statistically robust support (combined 5/248 cases, 0/250 controls, $P=0.061$). Using all five families the calculated NPL segregation score was non-significant (NPL=1.05, $P=0.125$, Supplementary Table 1). There was no evidence of *FAN1* LOH or somatic mutation in tumors from any of the five *FAN1* mutation carriers. LCLs derived from p.Cys47Ter and p.Asp140Thr carriers showed greater sensitivity to high doses of mitomycin C (MMC) compared to cells from a wild-type individual (p.Cys47Ter: $P=0.01$ Wilcoxon rank sum test; p.Asp140Thr: P -value not calculated). Transfection of a *FAN1* knockout human embryonic kidney cell line with p.Asp140Thr failed to reverse its MMC sensitivity.

Evidence Summary: Limited evidence of segregation and an increase in mutational burden in cases. Functional evidence for 2/4 mutations in non-colon tissue. No evidence of somatic mutation, LOH or functional effect in CRC.

TP53

Yurgelun *et al.* examined the frequency of rare germline *TP53* missense mutations in 457 patients with early-onset CRC (median age 36, range 15-40) and without a known hereditary cancer syndrome¹⁶. In six of the patients (1.3%), they identified missense changes in *TP53*. No comparison was made to the burden of *TP53* mutations in controls.

Based on this data they concluded that the frequency of *TP53* mutations is comparable with the proportion of inherited CRC thought to be attributable to germline *APC* mutations. This is a false comparison:

- (1) *TP53* missense variants are assumed to be deleterious. However, of the six variants they identified only one was predicted to be damaging by both SIFT and PolyPhen-2, with three being predicted benign by both algorithms. There was no other evidence that mutations had deleterious functional effect.
- (2) The proportion of inherited CRC thought to be attributable to germline *APC* mutations (1%) is not the same as the frequency of samples with *APC* missense changes (in our original 1,006 cases, not screened for mutations in known genes, 94 [9.3%] had rare [MAF<1%] missense changes in *APC*).

Evidence Summary: No comparison was made to mutational burden in controls. No evidence (even *in silico*) mutations have functional effect. No evidence of segregation or somatic mutation in CRC.

BUB1 and BUB3

BUB1: Disruptive mutations (p.Gln16Ter, p.Gln949Argfs) were identified in two of 23 Chinese early-onset (age at diagnosis ≤ 45) CRC cases; both variants were absent from 700 population controls but no comparison was made with the burden of mutations in controls¹⁷. No evidence of segregation was demonstrated and notably one of the mutation

carriers harbored a *MLH1* mutation (InSiGHT Class4: likely pathogenic), which unlike the *BUB1* variant was also carried by a sibling with polyps. The functional effect of mutations, LOH or somatic mutation within tumors was not assessed. No *BUB1* mutations were identified among 184 Netherlands/German cases.

Evidence Summary: Absence of identified mutations in controls, but no comparison made to mutational burden in controls. No functional evidence. No evidence of segregation or somatic mutation in CRC.

***BUB3*:** WES identified a novel damaging (as predicted by PolyPhen-2 SIFT, HOPE) missense mutation (p.Phe264Leu) in one of 10 Dutch early-onset (age at diagnosis ≤ 45) CRC cases (0/23 Chinese cases)¹⁷. Sequencing *BUB3* in 174 Netherlands/German CRC cases identified two further missense variants (p.Lys21Asn, p.Arg149Gln) that were predicted (by at least one algorithm) to be damaging, although PolyPhen-2 and SIFT predicted p.Arg149Gln to be benign. All three variants were absent in 1,154 controls but no comparison was made with burden of mutations in controls. No evidence of segregation, LOH or somatic mutation. Lymphocytes and primary skin fibroblasts from p.Phe264Leu and p.Lys21Asn mutation carriers showed significant enrichment of aneuploidy and structural abnormalities versus controls (p.Arg149Gln was not assessed/presented).

Evidence Summary: Absence of identified mutations in controls, but no comparison made to mutational burden in controls. Functional evidence in non-colon tissue. No evidence of somatic mutation, LOH or functional effect in CRC.

LRP6* and *PTPN12

Using WES de Voer *et al.* looked for genes recurrently affected by damaging missense mutations, assessed using a single prediction algorithm PhyloP, in 55 Dutch non-polyposis MMR-proficient early-onset (age at diagnosis ≤ 45) CRC cases¹⁸:

***LRP6*:** Damaging *LRP6* missense mutations (p.Trp239Leu, p.Asn789Ser, p.Thr867Ala) were identified in three cases. All three variants were within β -propeller domains. In mutation

carriers LRP6 protein showed no difference in protein expression or subcellular localization compared to wild-type. In Chinese Hamster Ovary cells p.Asn789Ser and p.Thr867Ala induced significant increases in WNT signaling activity ($P < 0.001$). Analysis of 174 additional Netherlands/German CRC and 2,329 population controls identified no additional damaging missense mutations in cases, with 18 identified in controls including p.Thr867Ala in three controls. By using only the 55 original cases and the 2,329 controls, de Voer *et al.* reported significant increase in mutation burden in cases versus controls (cases 3/55, controls 18/2,329, $P = 0.01$). However using all cases there was no significant enrichment of LRP6 mutation in cases versus controls (cases 3/229, controls 18/2,329, $P = 0.43$).

Evidence Summary: No increase in mutational burden in cases. Functional evidence in non-colon tissue. No evidence of somatic mutation, LOH or functional effect in CRC.

PTPN12: WES identified two damaging missense mutations (p.Arg522Met, p.Ser684Leu) in three of the 55 cases. Analysis of 174 additional Netherlands/German CRC and 2,329 population controls identified a new variant (p.Ala105Val) in one case and 11 variants in controls including previously identified p.Arg522Met in two controls and p.Ser684Leu in three controls. The burden of mutation in cases was significantly enriched versus controls (cases 4/229, controls 11/2329, $P = 0.039$) albeit non-significant after adjustment for multiple testing (three candidate genes, $P = 0.12$).

Evidence Summary: Limited evidence of an increase in mutational burden in cases. Variants not absent from controls. No evidence of somatic mutation, LOH or functional effect in CRC.

SUPPLEMENTARY FIGURES AND TABLES

Supplementary Figure 1: Forest plot of allelic odds ratio associated with *FANCM* p.Arg1931Ter (rs144567652) genotype and CRC. Studies [SMITH: original publication (2,207 cases, 2,176 controls)¹³; WES: whole-exome sequencing analyzed in this manuscript (863 cases, 1,604 controls); ENG: English Illumina Exome-BeadChip replication series (3,537 cases, 4,811 controls); SCOT: Scottish Illumina Exome-BeadChip replication series (2,015 cases, 1,981 controls)] were weighted according to the inverse of the variance of the log of the odds ratio (OR) calculated by unconditional logistic regression. Meta-analysis under a fixed-effects model was conducted using standard methods. Cochran's *Q* statistic to test for heterogeneity and the *I*² statistic to quantify the proportion of the total variation due to heterogeneity were calculated. Horizontal lines indicate 95% confidence intervals (CIs). Boxes indicate OR point estimate; its area is proportional to the weight of the study. Diamond (and broken line) indicates overall summary estimate, with CI given by its width. Unbroken vertical line indicates null value (OR=1.0).

Supplementary Table 1: Evidence for genes and variants being associated with CRC risk.

Analysis of the evidence presented in publications linking *NTHL1*, *RPS20*, *FANCM*, *FAN1*, *TP53*, *BUB1*, *BUB3*, *LRP6* and *PTPN12* with the risk of developing CRC

Abbreviations: AOD: Age of CRC diagnosis; EAS: East Asian; FCCTX: Familial CRC Cancer Type X; FS: frameshift; Het: heterozygous; Hom: homozygous; LOH: loss of heterozygosity; MMC: Mitomycin C; MS: missense; MSS: lack of mismatch repair deficiency (MMR) tested through either microsatellite stability or no loss of MMR proteins; NFE: Non-Finnish European; NPL: non-parametric linkage; NT: Not tested; SG: stop-gain; TS: target sequencing; WES: whole exome sequencing; WT: wild-type

Supplementary Table 2: Gene Burden analysis. Number of WES cases (n=863) and controls (n=1,604) with rare mutations in genes suggested to increase CRC risk. We considered three sets of variants: Class-1, disruptive mutations -stop-gain, frameshift; Class-2, predicted damaging -disruptive plus missense predicted to be damaging by CONDEL and splice site acceptor/donors; Class-3, all coding non-synonymous variants. Tables show -A all cases - n=863 and controls -n=1,604 with very rare -MAF<0.1% mutations; -B Amsterdam-II positive cases -n=159 and controls with rare -MAF<1% mutations; -C Amsterdam-II positive cases and controls with very rare -MAF<0.1% mutations. For each gene and variant class, numbers of cases and controls were compared and *P*-values calculated using Fishers exact test.

Supplementary Table 3: *BUB1*, *BUB3*, *FAN1*, *FANCM*, *LRP6*, *PTPN12*, *RPS20* and *TP53* variants -MAF<1% identified in 863 CRC cases and 1,604 controls. -See excel file

ACCEPTED MANUSCRIPT

REFERENCES

1. **Chubb D, Broderick P, Dobbins SE**, et al. Nat Commun 2016;7:11883.
2. Penegar S, et al. Br J Cancer 2007;97:1305-9.
3. Power C, et al. Int J Epidemiol 2006;35:34-41.
4. McKenna A, et al. Genome Res 2010;20:1297-303.
5. McLaren W, et al. Bioinformatics 2010;26:2069-70.
6. Gonzalez-Perez A et al. Am J Hum Genet 2011;88:440-9.
7. **Timofeeva MN, Kinnersley B**, et al. Sci Rep 2015;5:16286.
8. Kruglyak L et al. Am J Hum Genet 1996;58:1347-63.
9. Lek M, et al. Nature 2016;536:285-91.
10. Weren RD, et al. Nat Genet 2015;47:668-71.
11. Rivera B, et al. N Engl J Med 2015;373:1985-6.
12. Nieminen TT, et al. Gastroenterology 2014;147:595-598 e5.
13. Smith CG, et al. Hum Mutat 2013;34:1026-34.
14. Peterlongo P et al. Hum Mol Genet 2015;24:5345-55.
15. **Segui N, Mina LB**, et al. Gastroenterology 2015;149:563-6.
16. Yurgelun MB et al. JAMA Oncol 2015;1:214-21.
17. **de Voer RM, Geurts van Kessel A, Weren RD, Ligtenberg MJ**, et al. Gastroenterology 2013;145:544-7.
18. **de Voer RM, Hahn MM**, et al. PLoS Genet 2016;12:e1005880

Author names in bold designate shared co-first authorship

Supplementary Table 1: Evidence for genes and variants being associated with CRC risk. Analysis of the evidence presented in publications linking *NTHL1*, *RPS20*, *FANCM*, *FAN1*, *TP53*, *BUB1*, *BUB3*, *LRP6* and *PTPN12* with the risk of developing CRC

Abbreviations: AOD: Age of CRC diagnosis; EAS: East Asian; FCCTX: Familial CRC Cancer Type X; FS: frameshift; Het: heterozygous; Hom: homozygous; LOH: loss of heterozygosity; MMC: Mitomycin C; MS: missense; MSS: lack of mismatch repair deficiency (MMR) tested through either microsatellite stability or no loss of MMR proteins; NFE: Non-Finnish European; NPL: non-parametric linkage; NT: Not tested; SG: stop-gain; TS: target sequencing; WES: whole exome sequencing; WT: wild-type

Gene	CRC Gene Burden	Control Gene Burden	Gene Functional Data	Chr	Position (GRCh37)	c.DNA Change	Protein Change	Class	dbSNP	Segregation	Variant Case Frequency	Variant Control Frequency	ExAC Allele Frequency	Other Information
<i>NTHL1</i>	WES Hom 3/51 APC-MUTYH mutation-negative polyposis patients	17 Het/2329		16	2096239	c.268C>T	p.Gln90Ter	SG	rs150766139	Family1:2/2 Family2:3/3 Family3:2/2		17 Het/2329	NFE: 0.0023	Tumors significantly enriched for C:G>T:A transitions
<i>RPS20</i>	WES 4 individuals Finnish FCCTX family; TS 0/25 Finnish FCCTX	NT	RPS20 depletion in HeLa cells and LCLs from patients carrying c.147dupA showed increase in 21S pre-rRNA vs controls	8	56986283	c.147dupA	Val50SerfsTer23	FS		NPL=5.35, P=0.0078	1/26 FCCTX Family	0/584	Absent	No LOH (0/2)
<i>FANCM</i>	WES 1/50	NT		14	45667921	c.5791C>T	p.Arg1931Ter	SG	rs144567652	NT	2/2,258	1/2176	NFE: 0.0009	LOH WT 2/2
<i>FAN1</i>	WES 3 individuals Spanish FCCTX family. TS 4/247 Spanish FCCTX cases	0/250	LCL from p.C47Ter and p.Asp140Tyr carriers showed greater sensitivity to high doses of MMC. Transfection of HEK293T (Human Embryonic Kidney) FAN1 with p.Asp140Tyr failed to reverse its MMC sensitivity	15	31197007	c.141C>A	p.Cys47Ter	SG	rs144469584	NPL=0.94, P=0.25	1/248 FCXX Family	0/538	Absent	No LOH/somatic mutation. No reduction in FAN1 RNA or protein.
				15	31197284	c.418G>T	p.Asp140Tyr	MS	rs761776412	P/C, NPL=0, P=1	1/248 FCXX Family	0/250	NFE: 1.50E-5	Predicted to be damaging by SIFT/CONDEL. No LOH/somatic mutation.
				15	31197884	c.1018C>T	p.Pro340Ser	MS	rs771206220	NPL=1.05, P=0.125	1/248 FCXX Family	0/250	NFE: 1.50E-5	Predicted to be benign by SIFT/CONDEL

Gene	CRC Gene Burden	Control Gene Burden	Gene Functional Data	Chr	Position (GRCh37)	c.DNA Change	Protein Change	Class	dbSNP	Segregation	Variant Case Frequency	Variant Control Frequency	ExAC Allele Frequency	Other Information
<i>FAN1</i> (cont'd)				15	31206254	c.1771C>T	p.Arg591Trp	MS	rs3774 18523	NT	1/248 FCXX Family	0/250	NFE: 1.50E-5	Predicted to be damaging by PolyPhen- 2/SIFT/CONDEL
				15	31222812	c.2854C>T	p.Arg952Ter	SG	rs1847 45027	NPL=0.70, P=0.25	1/248 FCXX Family	0/250	NFE: 9.02E-5	
<i>TP53</i>	TS 6/457 North- American, Australian, New Zealand non- polyposis, AOD≤40	NT		17	7572973	c.1136G>A	p.Arg379His	MS		NT	1/457	NT	Absent	Predicted to be benign by PolyPhen- 2/SIFT
				17	7577069	c.869G>A	p.Arg290His	MS	rs5581 9519	NT	1/457	NT	NFE: 0.0002	Predicted to be benign by PolyPhen- 2/SIFT
				17	7577088	c.850A>T	p.Thr284Ser	MS	s14434 0710	NT	1/457	NT	Absent	Predicted to be possibly damaging by PolyPhen-2, benign by SIFT
				17	7577091	c.847C>T	p.Arg283Cys	MS	rs1496 33775	NT	1/457	NT	NFE: 0.0002	Predicted to be benign by PolyPhen- 2, possibly damaging by SIFT
				17	7577577	c.704A>G	p.Asn235Ser	MS	rs1443 40710	NT	1/457	NT	NFE: 0.0003	Predicted to be benign by PolyPhen- 2/SIFT
				17	7578475	c.455C>T	p.Pro152Leu	MS	rs5877 82705	NT	1/457	NT	NFE: 4.50E-5	Predicted to be damaging by PolyPhen-2/SIFT
<i>BUB1</i>	WES 2/23 Han Chinese ≤40; WES 0/10 Dutch non-polyposis MSS ≤40; TS 0/174 non- polyposis MSS Dutch/German CRC	NT	Disruption of <i>BUB1</i> exon 1 in HCT- 116 (MSI human CRC) caused chromosomal segregation defects	2	111398721	c.2844delC	p.Gln949Argfs	FS		NT	1/23 Chinese; 0/184 European	0/700	Absent	
				2	111431923	c.46C>T	p.Gln16Ter	SG		No	1/23 Chinese; 0/184 European	0/700	Absent	Carries MLH1 splice donor mutation which segregates with colorectal tumors

Gene	CRC Gene Burden	Control Gene Burden	Gene Functional Data	Chr	Position (GRCh37)	c.DNA Change	Protein Change	Class	dbSNP	Segregation	Variant Case Frequency	Variant Control Frequency	ExAC Allele Frequency	Other Information
<i>BUB3</i>	WES 0/23 Han Chinese \leq 40; WES 1/10 Dutch non-polyposis MSS \leq 40; TS 2/174 non-polyposis MSS Dutch/German CRC	NT		10	124914496	c.63G>C	p.Lys21Asn	MS		NT	0/23 Chinese; 1/184 European	0/1154		Predicted to be damaging by PolyPhen-2/SIFT/HOPE
				10	124919951	c.446G>A	Arg149Gln	MS	rs371545161	NT	0/23 Chinese; 1/184 European	0/1154	NFE: 7.52E-5 EAS: Absent	Predicted to be damaging by HOPE, benign by PolyPhen-2/SIFT/HOPE
				10	124922163	c.790T>C	p.Phe264Leu	MS		NT	0/23 Chinese; 1/184 European	0/1154	Absent	Predicted to be damaging by PolyPhen-2/SIFT/HOPE
<i>LRP6</i>	WES 3/55 Dutch non-polyposis MSS AOD \leq 45; TS 0/174 non-polyposis MSS Dutch/German CRC	18/2,329	In Chinese Hamster Ovary cells: no effect on LRP6 protein expression or localization; overexpression of p.Asn789Ser and p.Thr867Ala induced increased WNT signaling vs WT	12	12311955	c.2599A>G	p.Thr867Ala	MS	rs141458215	NT	WES: 1/55; TS: 0/174	3/2329	NFE: 0.0002	Predicted to be damaging by PhyloP
				12	12312812	c.2366A>G	p.Asn789Ser	MS		NT	WES: 1/55; TS: 0/174	0/2329	Absent	Predicted to be damaging by PhyloP
				12	12339985	c.716G>T	p.Trp239Leu	MS		NT	WES: 1/55; TS: 0/174	0/2329	Absent	Predicted to be damaging by PhyloP
<i>PTPN12</i>	WES 3/55 Dutch non-polyposis MSS \leq 45; TS 1/174 non-polyposis MSS Dutch/German CRC	11/2,329		7	77212900	c.314C>T	p.Ala105Val	MS		NT	WES: 0/55; TS: 1/174	0/2329	Absent	Predicted to be damaging by PhyloP
				7	77256561	c.1565G>T	p.Arg522Met	MS	rs537562368	NT	WES: 1/55; TS: 0/174	2/2329	NFE: 1.50E-5	Predicted to be damaging by PhyloP
				7	77261719	c.2051C>T	p.Ser684Leu	MS	rs201001953	NT	WES: 2/55; TS: 0/174	3/2329	NFE: 0.0012	Predicted to be damaging by PhyloP

Supplementary Table 2: Gene Burden analysis. Number of WES cases (n=863) and controls (n=1,604) with rare mutations in genes suggested to increase CRC risk. We considered three sets of variants: Class-1, disruptive mutations -stop-gain, frameshift; Class-2, predicted damaging -disruptive plus missense predicted to be damaging by CONDEL and splice site acceptor/donors; Class-3, all coding non-synonymous variants. Tables show -A all cases - n=863 and controls -n=1,604 with very rare -MAF<0.1% mutations; -B Amsterdam-II positive cases -n=159 and controls with rare -MAF<1% mutations; -C Amsterdam-II positive cases and controls with very rare -MAF<0.1% mutations. For each gene and variant class, numbers of cases and controls were compared and *P*-values calculated using Fishers exact test.

A

Gene	Class-1			Class-2			Class-3		
	Cases	Controls	<i>P</i> _{Fisher}	Cases	Controls	<i>P</i> _{Fisher}	Cases	Controls	<i>P</i> _{Fisher}
<i>BUB1</i>	0	4	0.31	1	8	0.17	5	16	0.36
<i>BUB3</i>	0	2	0.55	0	4	0.31	1	5	0.67
<i>FAN1</i>	0	2	0.55	3	5	1.00	7	10 [#]	0.62
<i>FANCM</i>	1	0	0.35	7	8	0.42	14 [§]	21	0.59
<i>LRP6 -BPD*</i>	0	0	-	3 -1	9 -6	0.56	8 -2	20 -11	0.55
<i>PTPN12</i>	0	1	1.00	3	3	0.43	8	7	0.17
<i>RPS20</i>	1	0	0.35	2	0	0.12	2	0	0.12
<i>TP53</i>	1	0	0.35	1	1	1.00	1	3	1.00

*Number of variants within β -Propellor domain

Total number of variants in controls = 11; 1 sample has 2 *FAN1* missense

§ Total number of variants in cases = 15; 1 sample has 2 *FANCM* missense

B

Gene	Class-1			Class-2			Class-3		
	Cases	Controls	<i>P</i> _{Fisher}	Case	Control	<i>P</i> _{Fisher}	Case	Control	<i>P</i> _{Fisher}
<i>BUB1</i>	0	4	1.00	1	8	0.58	4	30	0.54
<i>BUB3</i>	0	2	1.00	0	4	1.00	0	5	1.00
<i>FAN1</i>	0	2	1.00	2	17	0.69	6	45	0.46
<i>FANCM</i>	0	1	1.00	1	33	0.36	7 [§]	67 [§]	0.84
<i>LRP6 -BPD*</i>	0	0	-	1 -1	17 -13	1.00	2 -2	37 -21	0.57
<i>PTPN12</i>	0	1	1.00	1	5	0.44	3	9	0.09
<i>RPS20</i>	0	0	-	1	0	0.09	1	0	0.09
<i>TP53</i>	0	0	-	0	1	1.00	0	4	1.00

*Number of variants within β -Propellor domain

§ Total number of variants in cases = 8 in controls = 69; 3 samples have 2 *FANCM* missense

C

Gene	Class-1			Class-2			Class-3		
	Cases	Controls	P_{Fisher}	Cases	Controls	P_{Fisher}	Cases	Controls	P_{Fisher}
<i>BUB1</i>	0	4	1.00	1	8	0.58	1	16	1.00
<i>BUB3</i>	0	2	1.00	0	4	1.00	0	5	1.00
<i>FAN1</i>	0	2	1.00	0	5	1.00	2	10	0.30
<i>FANCM</i>	0	0	-	0	8	1.00	2 [§]	21	1.00
<i>LRP6 -BPD*</i>	0	0	-	0	9-6	1.00	0	20-11	0.25
<i>PTPN12</i>	0	1	1.00	0	3	1.00	1	7	0.53
<i>RPS20</i>	0	0	-	1	0	0.09	1	0	0.09
<i>TP53</i>	0	0	-	0	1	1.00	0	3	1.00

*Number of variants within β -Propellor domain

§ Total number of variants in cases = 3; 1 sample has 2 *FANCM* missense

Gene
Chr
Position (GRCh37)
Ref
Alt
CONSEQUENCE
Class
c.DNA Change
Protein Change
Control_MAF
Case_MAF

ALT_Control

ALT_Cases

ALT_Ams
ExAC_Freq
ExAC_NFE
1000G_ALL
1000G_EUR
Condel
CADD_PHRED
SIFT_score
SIFT_pred
Polyphen2_HDIV_score
Polyphen2_HDIV_pred
dbSNP
COSMIC_ID
COSMIC_DIS
ClinVar_SIG
ClinVar_DIS

Supplementary Figure 1: Forest plot of allelic odds ratio associated with *FANCM* p.Arg1931Ter (rs144567652) genotype and CRC. Studies [SMITH: original publication (2,207 cases, 2,176 controls)¹; WES: whole-exome sequencing analyzed in this manuscript (863 cases, 1,604 controls); ENG: English Illumina Exome-BeadChip replication series (3,537 cases, 4,811 controls); SCOT: Scottish Illumina Exome-BeadChip replication series (2,015 cases, 1,981 controls)] were weighted according to the inverse of the variance of the log of the odds ratio (OR) calculated by unconditional logistic regression. Meta-analysis under a fixed-effects model was conducted using standard methods. Cochran's Q statistic to test for heterogeneity and the I^2 statistic to quantify the proportion of the total variation due to heterogeneity were calculated. Horizontal lines indicate 95% confidence intervals (CIs). Boxes indicate OR point estimate; its area is proportional to the weight of the study. Diamond (and broken line) indicates overall summary estimate, with CI given by its width. Unbroken vertical line indicates null value (OR=1.0).

