

ROLE OF RARE VARIANTS IN UNDETERMINED MULTIPLE ADENOMATOUS POLYPOSIS AND EARLY ONSET COLORECTAL CANCER.

Jérémie H. Lefevre ^{1,2,3}, Carolina Bonilla ^{1,2#}, Chrystelle Colas ^{3, 4}, Bruce Winney ¹, Elaine Johnstone ⁵, Susan Tonks ¹, Tammy Day ¹, Katarzyna Hutnik ¹, Abdelhamid Boumertit ¹, Florent Soubrier ⁴, Yann Parc ³, Walter F. Bodmer ^{1, 2 *}.

¹Cancer and Immunogenetics Laboratory, Department of Clinical Pharmacology, University of Oxford, Old Road Campus Research Building, Old Road Campus, Headington, Oxford OX3 7DQ, United Kingdom

²Weatherall Institute of Molecular Medicine (WIMM), John Radcliffe Hospital, University of Oxford, Headington, Oxford OX3 9DS, United Kingdom

³Department of Digestive Surgery, Hôpital Saint-Antoine, AP-HP, University Paris VI (Pierre and Marie Curie), Paris, France

⁴Laboratory of Angiogenetics and Oncogenetics, Hôpital Pitié-Salpêtrière, AP-HP, University Paris VI (Pierre and Marie Curie), Paris, France

⁵Translational Oncology Group, Department of Clinical Pharmacology, Old Road Campus Research Building, Old Road Campus, Headington, Oxford OX3 7DQ, United Kingdom

Short title: Rare variants in polyposis and early colorectal cancer

#Current address:

School of Social and Community Medicine

University of Bristol

Oakfield House

Oakfield Grove

Bristol BS8 2BN

United Kingdom

***Corresponding author:**

Walter F. Bodmer

Tel: +44 1865 222422

Fax: +44 1865 222431

Email: walter.bodmer@hertford.ox.ac.uk

ABSTRACT

INTRODUCTION: Some 15-20% of multiple adenomatous polyposis have no genetic explanation and 20-30% of colorectal cancer cases (CRC) are thought to be due to inherited multifactorial causes. Accumulation of the deleterious effects of a series of low frequency dominant and independently acting variants of a variety of different genes may be a partial explanation for such patients. The aim of this study was to type a selection of rare and low frequency variants (<5%) in a set of patients with undetermined adenomatous polyposis or early onset CRC in order to elucidate the role of these variants in CRC susceptibility.

MATERIALS AND METHODS: 1181 subjects were included in the study (866 controls and 315 cases). Cases comprised UK (n=184) and French (n=131) patients with adenomatous polyposis of unknown origin (n=187) or early onset CRC (n=128). Seventy variants in 17 different genes were examined in cases and controls. The effect of each variant on protein function was investigated *in silico*.

RESULTS: In the UK, out of the 70 variants typed, 36 (51%) were tested for association with CRC. Twenty-one variants were considered rare (minor allele frequency in the control population <1%). Four rare variants were found to have a significantly higher minor allele frequency in cases than in controls (EXO1-12, MLH1-1, CTNNB1-1 and BRCA2-37, $p < 0.05$). Pooling all rare variants with a minor allele frequency lower than 0.5% showed an excess risk in cases (OR 3.2; 95% CI 1.1-9.5; $p = 0.04$).

CONCLUSION: Rare variants are important risk factors in CRC and as such, should be systematically assayed alongside common variation in the search for the genetic basis of complex diseases.

INTRODUCTION

Familial Adenomatous Polyposis (FAP) and MYH Associated Polyposis (MAP) are two inherited syndromes that show a high incidence of adenomatous polyps and an elevated risk of developing colorectal cancer (CRC). They account for a small fraction of CRC, less than 4%¹. However, despite the fact that these two syndromes are caused by deleterious highly penetrant mutations in *APC* (GeneID 324)² and *MUTYH* (GeneID 4595)³, around 15-20% of patients with polyposes have no known genetic risk factors. This is especially so for multiple polyposis patients who carry between 3 and 100 adenomatous polyps. Moreover, 20-30% of CRC is thought to be due to inherited multifactorial causes⁴. In the absence of identification of a new deleterious mutation, CRC may in part be due to the summation of the deleterious effects of a series of low frequency dominant and independently acting variants of a variety of different genes, each, conferring a moderate but readily detectable increase in relative risk⁴. This 'rare variant' hypothesis was based upon the observation by Frayling *et al.* of the APC I1307K and E1317Q variants in patients with multiple adenomas⁵. I1307K is found in the Ashkenazi Jewish population at a frequency of ~6-7% while it is absent from non-Jewish populations, and confers an increased risk of multiple adenomas and colorectal cancer⁵. This variant implies an amino I-K (Isoleucine to Lysine) substitution in a region involved in protein binding leading to a mild dominant-negative effect. The E1317Q variant substitution may also affect the function of the APC protein presumed to translate into a slight but definitive advantage for the growth of a tumor⁶.

Following these observations other rare variants have been tested. The candidate variants were selected because of their known involvement in sporadic or hereditary colorectal cancer or adenomas. Fearnhead *et al.* (2004) observed a cumulative effect of 13 rare variants on five different genes in a cohort of 124 patients with adenomatous polyposis with an overall odds ratio of 2.2 ($p=0.0001$) when compared with a control set⁷.

Since this publication, several variants in different CRC susceptibility genes such as *hMLH1*, *hMSH6* have been reported to increase the risk of CRC but not cause Lynch syndrome^{8,9} while *CHEK2* confers a higher colorectal cancer risk in HNPCC/HNPCC related families¹⁰.

Rare variants are defined by a minor allele frequency (MAF) lower than 1% in the general population and are unlikely to be identified by Genome Wide Association Studies (GWAS) due to their low frequency and their small contribution to the overall susceptibility of a disease⁴. Only variants with a frequency higher than 5% are detected in these large case-control association studies. Rare variants are best identified in studies with selected cases and selected candidate genes already known to be likely to be functionally relevant^{1,5}. Patients with early onset CRC (before the age of 50) and multiple polyposis (3-100 polyps) with no known mutations in *APC* or *MYH* are ideal candidates to demonstrate an elevated predisposition to disease due to the accumulation of rare variants, as they are likely to involve inherited susceptibility.

The aim of this study was therefore to type a selection of rare (MAF<1%) and low frequency variants (MAF 1-5%) in a relatively large set of patients with undetermined multiple polyposis (3-100 polyps) or early onset CRC (diagnosed before 50 years of age) in order to elucidate the wider role of such variants in CRC susceptibility.

MATERIALS AND METHODS

A total of 315 cases and 866 controls, 1181 subjects in all, were included in this study. Collection of blood samples from cases and controls and clinicopathological information from patients were undertaken with appropriate individual informed consent and local ethical committee approvals.

Controls

The controls comprised 866 individuals collected in 10 different regions across the UK as part of the People of the British Isles (PoBI) study¹¹ (see below), and were unselected with respect to disease status.

Cases

The UK patient group consisted of 112 individuals with 3 to 100 histologically proven synchronous or metachronous adenomatous polyps and 72 individuals with colorectal cancer diagnosed before 50 years of age. Sixty-three individuals with early onset disease were obtained through the VICTOR clinical trial, a Phase III double-blind placebo controlled study of rofecoxib in Dukes stage B or C CRC patients following potentially curative therapy, while the remaining nine cases were recruited through the John Radcliffe and Churchill hospitals' gastrointestinal clinics. With the exception of one Black Caribbean and one Indian individual, ethnic origin was White British for all UK patients for whom information was available. Non-white individuals were excluded from further analysis. No patient fulfilled the criteria for Familial Adenomatous Polyposis (FAP), autosomal recessive *MUTYH*-associated polyposis (MAP), or Hereditary NonPolyposis Colorectal Cancer (HNPCC) on clinical grounds. Some of these patients had already been screened for germline mutations in the *APC* and *MYH* genes in previous studies^{5,12}.

We also collected samples from 131 French patients including 75 with multiple adenomas and 56 with early onset CRC, who were recruited in the Department of Digestive Surgery at the Hospital Saint-Antoine in Paris using the criteria described above. Cases were selected from those who underwent a colectomy or total colectomy for CRC or polyposis. Patients diagnosed with CRC before age 50 or with more than three polyps detected after 2005, were referred for a consultation with the geneticist. Immunohistochemical staining to determine loss of expression of the genes *MLH1* and *MSH2* and microsatellite status was performed **for** all patients with a CRC diagnosed

before the age of 50. Microsatellite instability was confirmed by PCR. Sequencing of the entire *MUTYH* and *APC* genes was carried out in patients with more than three adenomatous polyps¹³. Only patients with no indication of HNPCC, MAP or FAP were included in this study. No ethnic identification was available for the French patients.

All UK and French cases had histological confirmation of adenomatous polyps but not all of them had the precise number of polyps determined. For 24 UK and 14 French adenoma patients only “multiple” was recorded. Within both the UK and French patient groups, individuals with attenuated familial adenomatous polyposis (AFAP) may be included, as they were not purposely eliminated from the study.

Variant selection

Rare and low frequency variants were chosen based on prior literature reports that suggested a putative association with CRC, and with features related to colorectal disease, gastric cancer or other cancers (mostly of the breast and prostate). Variants in the *BRCA* genes were specifically selected from those that were classified either as non-pathogenic or as of unknown significance (see Breast Cancer Information Core database). Common variants, such as *CDH1* rs16260, *MTHFR* rs1801133 and *TP53* rs1042522, were genotyped because it has been suggested that they are associated with several types of cancer, including CRC¹⁴⁻¹⁸.

DNA extraction and processing

Genomic DNA was extracted from patients’ peripheral venous blood using standard techniques. The PoBI control blood samples were transported at room temperature to the laboratory, where the peripheral blood lymphocytes were separated under sterile conditions within two days of collection. DNA was prepared from the 10ml blood residue remaining after sterile separation using either magnetic beads (GeneCatcherTM, Invitrogen, Carlsbad, CA, USA) or spin columns (Qiagen, Valencia, CA, USA). DNA concentration was determined using Pico Green¹⁹ and normalised for

genotyping to 25ng/uL. Samples from UK cases underwent whole genome amplification due to limited volumes and amounts of genomic DNA. We used the Repli-g Mini kit (Qiagen) which implements a multiple displacement amplification reaction to generate up to 10 ug of DNA per 50 uL reaction from a starting amount of at least 10 ng of genomic template. Genomic DNA from French cases and UK controls was used for genotyping.

Genotyping

We examined 70 variants, in cases and controls, in the following cancer candidate genes: *APC*, *AXIN1*, *AXIN2*, *BRCA1*, *BRCA2*, *CDH1*, *CHEK2*, *CTNNB1*, *EPHB2*, *EXO1*, *MLH1*, *MLH3*, *MSH2*, *MTHFR*, *PMS2*, *SMAD4* and *TP53*, which were selected based on their involvement in familial cancers and the presence of somatic mutations in cancer. Of these, sixteen variants were genotyped in a subset of only 227 controls. The complete list of variants analyzed is given in Supplementary Table 1.

Four variants were genotyped using restriction fragment length polymorphism (RFLP) analysis (*CDH1*-2, *CHEK2*-1, *CTNNB1*-1, *MSH2*-8). Variants *APC*-10 and *APC*-11 (i.e. I1307K and E1317Q) were typed using allele-specific PCR⁷. For details of primers, enzymes and fragment sizes for these see Supplementary Table 2. Primers and conditions for all other variants are available upon request. Genotyping of the remaining variants was done using the Sequenom MassArray technology, namely matrix-assisted laser desorption/ionization time-of-flight (MALDI-TOF) mass spectrometry and the iPLEX Gold assay (Sequenom Inc., San Diego, CA, USA)²⁰.

Statistical analysis

Hardy-Weinberg equilibrium was assessed using an exact test implemented in the program PLINK v.1.07²¹. Case-control association analyses were also conducted with PLINK. Two-sided p-values were calculated using Fisher's exact test and those below 0.05 (with no multiple comparison correction) were considered statistically significant for an initial analysis. Combined ORs were estimated using the Mantel-Haenszel test^{22,23}.

Functional in silico analysis

We used the web-based programs PolyPhen-2 and SNPs&GO to predict the effect of non-synonymous variants on protein function^{24,25}. FastSNP and F-SNP were similarly used for non-coding variants^{26,27}.

RESULTS

Populations

Clinical characteristics of the patients and controls are shown in Table 1.

Variant selection

Among the 70 variants examined, 24 (34%) were monomorphic in both UK cases and controls and therefore not useful for analysis (Supplementary Table 1). Of the remaining 46, ten were monomorphic in cases only and five were monomorphic only in controls. Thirty-one variants were considered rare having a control population MAF below 1%, seven were low frequency variants (i.e. MAF between 1 and 5%) and eight were common polymorphisms (i.e. MAF >5%). If we define variant class based on the combined MAF in cases and controls as recently suggested^{28,29} only one variant (APC-17) changes categories, going from being a low frequency variant to a rare one.

No variant was out of Hardy-Weinberg equilibrium in the control population at a Bonferroni-corrected p-value of 0.001 (0.05/46) or less. Three variants were in Hardy-Weinberg disequilibrium in controls (CDH1-1 TP53-1 and BRCA1-6, $p < 0.05$) and three in patients (EPHB2-3, EXO1-12 and BRCA1-22, $p < 0.05$) if a correction for multiple testing was not applied.

Association analysis

UK cases vs controls

When comparing UK cases with controls four rare variants were found to have a significantly higher MAF among the patients (EXO1-12, MLH1-1, CTNNB1-1 and BRCA2-

37, $p < 0.05$; Table 2). Variant EX01-12 was more frequent in individuals with cancer than in those with adenomas, as opposed to MLH1-1, CTNNB1-1 and BRCA2-37, which were only present in the multiple adenoma cases. (Table 2). Results close to significance were also seen for rare variant EPHB2-3 and common variant CDH1-2 ($p = 0.07$ for both), although the CDH1-2 A allele appears to protect against disease. When analyzing carrier frequencies instead of allele frequencies, only BRCA2-37 was significant, with an OR of 4.1 (1.2-14.3, $p = 0.05$, Table 3). Pooling together all the rare variants showed that the proportion of patients carrying rare variants is higher than the proportion of control carriers, regardless of whether the full set or a subset of controls was used (Table 3). The combined OR, obtained by merging OR1 (effect of variants typed in the full set of controls) and OR2 (effect of variants typed in the smaller set of controls) with the Mantel-Haenszel test, was 1.2 (95% CI, 0.8-1.8, $p = 0.42$). This effect became much stronger when only variants with a MAF lower than 0.5% were tested (combined OR 1.8, 95% CI, 1.0-3.1, $p = 0.05$; Table 3). On the other hand, the analysis of pooled low frequency variants showed a protective but non-significant effect (combined OR 0.8, 95% CI, 0.5-1.1, $p = 0.18$; Supplementary Table 3). When APC-17 is switched categories, from low frequency to rare variant, results change slightly. For rare variants with $MAF < 1\%$, combined OR = 1.1 (95% CI, 0.8-1.7, $p = 0.54$); for low frequency variants, combined OR = 0.8 (95% CI, 0.5-1.2, $p = 0.24$). Also, two variants (MLH3-1 and CHEK2-1) do not make the 0.5% cutoff when assessed from the frequency in the combined set of cases and controls. Taking them out of the analysis of variants with $MAF < 0.5\%$ yields a combined OR of 1.6 (95% CI, 0.8-2.9, $p = 0.17$).

UK multiple adenoma patients vs early onset CRC patients

Analysis by disease group (i.e. multiple adenoma vs early onset patients) of all variants with frequencies below 0.5% revealed an increase in susceptibility to disease for carriers of rare variants, especially among multiple adenoma patients (combined OR 1.9;

95% CI, 1.0-3.5; $p=0.05$; Table 4) Individually, the carrier frequency for BRCA2-37 in multiple polyp cases was significantly higher than that of controls, while this variant was absent among early onset patients. MLH3-1, on the other hand, showed a significantly higher carrier frequency in the early onset group as compared to controls, whereas BRCA2-27 was only detected in individuals with early onset disease (Table 4). Overall, out of the 31 variants with $MAF < 1\%$, 14 were present in multiple adenoma patients only, while 3 (4, if counting APC-17) were present only in early onset cases. The difference is significant whether APC-17 is included or not ($p < 0.05$), although the smaller number of early onset CRC patients might be introducing bias.

There were significant allele frequency differences between individuals with multiple adenomas and those with early onset CRC in two variants in *MLH3* (one rare [MLH3-1] and one low frequency [MLH3-5]) and in a common variant in *CDH1* (CDH1-5) (Table 2). In these three instances the allele frequency among early onset patients was higher than among individuals with multiple adenomatous polyps.

Comparison between UK and French samples

Twenty-two of the variants genotyped in UK patients (16 rare, 4 low frequency and 2 common variants) were also examined in French subjects affected by either multiple polyps or early onset CRC, recruited using the same set of criteria employed in the UK (Table 5). Three rare variants (MSH2-8, APC-10 and BRCA2-48) were absent from both, the UK and the French population. Eight variants that were detected in UK cases were not found in French patients (EPHB2-1, EPHB2-4, EPHB2-7, EXO1-4, CTNNB1-1, BRCA2-35, BRCA2-37, CHEK2-1). On the other hand, no variant identified in French patients was missing from the UK sample (cases and/or controls). There was no significant difference between French and UK cases with respect to the overall number of rare variants with $MAF < 1\%$, but UK patients show an excess of rare variants with $MAF < 0.5\%$ compared with this set of French patients ($p=0.02$, Supplementary Table 4

In silico analysis of functional effects

We investigated the putative effects of non-synonymous variants with the programs PolyPhen-2 and SNPs&GO. Based on PolyPhen-2, 11 variants were classified as probably damaging, 6 as possibly damaging and 22 as benign, whilst according to SNP&GO there were 18 disease variants and 21 neutral variants. Although these numbers seem fairly similar there were several disagreements between programs with respect to the prediction of particular variants (Table 2). Among the rare variants with $MAF < 1\%$, there were nine probably damaging, five possibly damaging and 14 benign, or 13 disease and 15 neutral variants. When only those variants with $MAF < 0.5\%$ are considered the ratio of damaging (probably + possibly)/disease to benign/neutral variants increases from ~50 to 60%. All of the low frequency variants, in contrast, were predicted to be benign/neutral. In addition to the missense variants, there was one synonymous (SMAD4-1) and one deleterious coding (CHEK2-1) rare variants, and five non-coding variants (one rare and one common variant in the promoter and two common intronic variants in *CDH1*, and one low frequency variant in the 3'UTR of *APC*). Promoter variant CDH1-1 is predicted to eliminate a S8 transcription factor binding site, while the CDH1-2 promoter variant A allele has been found to decrease transcriptional efficiency by 68% with respect to the C allele, also probably by altering transcription factor binding sites ³⁰. Using the programs FastSNP and F-SNP a variant in *CDH1* intron 1 was determined to potentially affect a splicing site while a variant in intron 4 of the same gene showed a low risk of being an intronic enhancer (Table 2).

DISCUSSION

We have examined 55 rare variants, seven low frequency variants and eight polymorphisms in a sample of UK colorectal cancer and multiple adenoma cases and controls. Two of the four rare variants that were individually significantly associated with

disease (i.e. *MLH1*-1 and *CTNNB1*-1) had already been identified in the same set of individuals with multiple polyps, although not then found to be individually significant⁷. In this study we showed that these *MLH1* and *CTNNB1* variants were not present in a different and much larger UK control population, which explains the present case-control significant difference, and were also absent from a sample of early onset CRC UK patients. Given that these two variants have not been found in our set of French patients, and that having ~30% fewer French cases may not in itself fully explain the UK-French apparent difference, they may represent UK founder effects, as previously suggested⁴. However, replication of our findings in another UK sample of multiple adenoma patients, as well as functional studies, are necessary to establish their importance as CRC risk factors. The remaining two individually significant variants (*EXO1*-12 and *BRCA2*-37) have not been associated with CRC before. *BRCA2*-37 was classified as not clinically significant by the Breast Cancer Information Core database (in early 2010) and predicted to be benign by PolyPhen-2, yet it was recently found to be overrepresented among subjects with familial prostate cancer³¹ and SNP&GO considered it to be disease associated. As mentioned above, the fact that it was not found among French patients could indicate a restricted distribution of this variant. Variant *EPHB2*-3, which was detected in a Finnish individual with rectal and prostate cancer in an earlier study³², showed a nearly significant result. All associated variants code for non-synonymous amino acid changes. However, *CTNNB1*-1, *BRCA2*-37 and *EPHB2*-3 were predicted to be benign by PolyPhen-2, whereas *MLH1*-1 and *EXO1*-12 were considered probably damaging. SNPs&GO, on the other hand, predicted *CTNNB1*-1 and *BRCA2*-37 to be disease-associated and the remaining three variants to be neutral. Recently, SNPs&GO has been found to be more accurate than PolyPhen and other similar programs³³. However, even though it identifies *CTNNB1*-1 and *BRCA2*-37 as potentially pathogenic, it misses *MLH1*-1 and *EXO1*-12. Also, *APC*-11, demonstrably pathogenic⁵, was not identified as such by any of these computational

methods. These discrepancies indicate that the use of *in silico* methods to evaluate the effects of non-synonymous rare variants is not yet sufficiently reliable to be confident of their predictions. This is especially important when using them to predict which variants to focus on.

The grouping of all rare variants in the association analysis (23 or 8, depending on the control set used) yielded a combined OR of ~ 1.2 , which suggested there was no strong evidence of an effect on CRC. However, pooling all variants with a MAF lower than 0.5% considerably bolstered the association, taking the OR to ~ 1.8 . Notably, even though several of the variants included in the analysis are, on the basis of the *in silico* analysis and the examination of other parameters of pathogenicity³⁴, considered to be benign, neutral, not clinically significant or of unknown significance, there is nevertheless an elevated risk from their combined action. The conclusion is that these variants may well be pathologically relevant, but that the *in silico* approaches are not yet adequate to detect this. The low frequency variants (MAF between 1% and 5%) do not appear to influence susceptibility to CRC, as we described earlier for *CCND1*³⁵. It is clear that further research is needed to evaluate more fully the role of low frequency variants in cancer³⁶. Defining rare variants using a threshold based on the combined set of cases and controls, as compared to just the controls, only altered the classification of three variants in our study and so did not appreciably affect the results. This was to be expected since we had a substantially larger number of controls than cases.

Extensively studied common variants MTHFR A222V (rs1801133) and TP53 R72P (rs1042522) did not show significant frequency differences between cases and controls. Conversely, CDH1 -284C/A (rs16260) exhibited a lower frequency of the A allele in patients than in controls, revealing weak statistical evidence of a protective effect of this polymorphism on colorectal disease (0.25 vs. 0.31, $p=0.07$). This is in agreement with previous findings on CRC where the C allele increases risk^{14,15}, while in gastric, prostate

and breast cancer the A allele tends to be the risk allele³⁷⁻³⁹. However, our study is underpowered for the detection of effects from common variants.

The analysis by disease group showed that, even though the collection of rare variants in each set of patients carries a higher risk of disease, our findings are mostly driven by the effects on individuals with multiple adenomas. BRCA2-37, the rare variant with the strongest effect in this study, was, for example, found only in patients with multiple adenomas. Although the sample size for the early onset group was limited, our results clearly suggest that the genetic influence on colorectal cancer may mostly be seen in individuals with multiple adenomas, as compared to early onset cases. This parallels what is found in the clear cut familial cases of inherited colorectal cancer. The extra layer of activity needed to go from polyp to cancer leads to an additional amount of variation that may be 'less genetically determined' and so obscure the underlying genetic susceptibility due to the multiple adenomas. There were, however, no significant differences in carrier frequencies between the two groups of patients despite the fact that over half of the variants with MAF<1% were found only in the multiple adenoma group. This, again, is probably due to the relatively smaller size of the early onset group of patients. However, the allelic frequencies of two missense variants (one rare and one low frequency) in *MLH3* and one intronic common variant in *CDH1* differed significantly between multiple adenoma and early onset CRC cases, with the latter exhibiting higher frequencies of these variants. This suggests that different sets of rare variants are quite likely to be involved in different pathologies, but to detect their effect would require larger numbers of patients than we were able to study.

The association p-values reported in this study have not been corrected for multiple hypotheses testing. Taking into account the number of variants analysed, a Bonferroni correction would take the significance threshold to 0.001. Nonetheless, we believe that because there is an a priori case for each candidate variant to be potentially

functional, such a correction would be unsuitably stringent. The lack of French controls precluded a similar association study from being carried out with French samples as was done for the UK samples. Using UK controls would be inappropriate because of population stratification within Europe, especially for analysis of rare variants as they are likely to be population specific. The presence of such founder effects is appreciably suggested by the fact that the variants MLH1-1 and CTNNB1-1, which are very clearly associated with multiple adenomas in UK cases, were not found in the French multiple adenoma cases. Further analysis of such differences requires larger numbers of French cases and appropriately selected French controls.

In summary, because rare variants appear to be associated with higher ORs than common variants, a relatively small study like ours can uncover the effects of candidate variants with low population frequencies on complex diseases such as, in this case, colorectal cancer. We have also shown that variants with frequencies below 0.5% appear to have the biggest effects regardless of the *in silico* prediction of their function. The role of the individual variant BRCA2-37 (V2728I) on the development of multiple adenomatous polyps deserves further examination. In general, the multiple adenoma phenotype seems to be more susceptible to genetic influence than early onset CRC, but a larger early onset patient sample would be necessary to confirm this finding. We have found some differences between UK and French patients in terms of the distribution of rare variants that justify closer inspection as population stratification within Europe can lead to spurious association results.

To conclude, we have confirmed that rare variants are important risk factors in CRC and as such, should be systematically assayed alongside common variation in the search for the genetic basis of complex diseases taking great care to match cases with appropriate controls.

ACKNOWLEDGMENTS

We would like to thank all individuals who contributed samples for this study.

We are also grateful to Shazad Asraf, Oliver Sieber, and Haitao Wang for making samples available.

We thank the patients who participated in VICTOR, without whose help, the study would not have been possible. The trial was enabled by an educational study grant from Merck & Co. Inc. who also supplied the study drug. VICTOR was endorsed by the UK National Cancer Research Institute (NCRI) and Cancer Research Campaign Clinical Trials Awards and Advisory Committee (CTAAC) allowing support from the UK National Cancer Research Network. The University of Birmingham was the centre for randomisation and for statistics initially before these moved to the Universities of Oxford and Warwick respectively. We also thank VICTOR trial staff: A Stanley, J Stokes, F Duchesne, K Reed, S Pendlebury, R Duvvuri, I Kennedy, G Davis, J Birtwistle, L Blair, A Bange, K Bell, L Dunham, R Deacon, A Mellor, B Mikolajczyk, M Davoudianfar and A Ferch; and VICTOR laboratory support staff: K Elliot, M Presz, H Wang; Clinical Coordinator: D Rea; The VICTOR Trial Advisory Group and the VICTOR Data and Safety Monitoring Committee and all nurses and medical staff who participated in this trial.

CB was Cancer Research UK Julia Bodmer Fellow.

This work was made possible by a Cancer Research UK programme grant, a Wellcome Trust grant and a Serroussi Foundation grant to WFB.

URLs:

- People of the British Isles study (PoBI): www.peopleofthebritishisles.org

- Breast Cancer Information Core (BIC) dataset:
<http://research.nhgri.nih.gov/bic/index.shtml>
- PolyPhen v.2: <http://genetics.bwh.harvard.edu/pph2>
- PLINK : <http://pngu.mgh.harvard.edu/~purcell/plink/>
- SNPs&GO : <http://snps-and-go.biocomp.unibo.it/snps-and-go/>
- FastSNP:
http://fastsnp.ibms.sinica.edu.tw/pages/input_CandidateGeneSearch.jsp
- F-SNP: <http://compbio.cs.queensu.ca/F-SNP>

REFERENCES

1. Fearnhead NS, Wilding JL, Bodmer WF. Genetics of colorectal cancer: hereditary aspects and overview of colorectal tumorigenesis. *British Medical Bulletin* 2002;64:27–43.
2. Bodmer WF, Bailey CJ, Bodmer J, Bussey HJ, Ellis A, Gorman P, Lucibello FC, Murday VA, Rider SH, Scambler P. Localization of the gene for familial adenomatous polyposis on chromosome 5. *Nature* 1987;328:614–6.
3. Al-Tassan N, Chmiel NH, Maynard J, Fleming N, Livingston AL, Williams GT, Hodges AK, Davies DR, David SS, Sampson JR, Cheadle JP. Inherited variants of MYH associated with somatic G:C-->T:A mutations in colorectal tumors. *Nature Genetics* 2002;30:227–32.
4. Bodmer W, Bonilla C. Common and rare variants in multifactorial susceptibility to common diseases. *Nature Genetics* 2008;40:695–701.
5. Frayling IM, Beck NE, Ilyas M, Dove-Edwin I, Goodman P, Pack K, Bell JA, Williams CB, Hodgson SV, Thomas HJW, Talbot IC, Bodmer WF, et al. The APC variants I1307K and E1317Q are associated with colorectal tumors, but not always with a family history. *Proceedings of the National Academy of Sciences* 1998;95:10722–7.
6. Bodmer W. Familial adenomatous polyposis (FAP) and its gene, APC. *Cytogenetic and Genome Research* 1999;86:99–104.
7. Fearnhead NS, Wilding JL, Winney B, Tonks S, Bartlett S, Bicknell DC, Tomlinson IPM, Mortensen NJM, Bodmer WF. Multiple rare variants in different genes account for multifactorial inherited susceptibility to colorectal adenomas. *Proceedings of the National Academy of Sciences* 2004;101:15992–7.

8. Allan JM, Shorto J, Adlard J, Bury J, Coggins R, George R, Katory M, Quirke P, Richman S, Scott D, Scott K, Seymour M, et al. MLH1 -93G>A promoter polymorphism and risk of mismatch repair deficient colorectal cancer. *International Journal of Cancer* 2008;123:2456–9.
9. Tulupova E, Kumar R, Hanova M, Slyskova J, Pardini B, Naccarati A, Polakova V, Vodickova L, Novotny J, Hemminki K, Halamkova J, Vodicka P. Do polymorphisms and haplotypes of mismatch repair genes modulate risk of sporadic colorectal cancer? *Mutation Research* 2008;648:40–5.
10. Wasielewski M, Vasen H, Wijnen J, Hooning M, Dooijes D, Tops C, Klijn JGM, Meijers-Heijboer H, Schutte M. CHEK2 1100delC is a susceptibility allele for HNPCC-related colorectal cancer. *Clinical Cancer Research* 2008;14:4989–94.
11. Winney B, Boumertit A, Day T, Davison D, Echeta C, Evseeva I, Hutnik K, Leslie S, Nicodemus K, Royrvik EC, Tonks S, Yang X, et al. People of the British Isles: preliminary analysis of genotypes and surnames in a UK-control population. *European Journal of Human Genetics* 2011;:1–8.
12. Lamlum H, Al Tassan N, Jaeger E, Frayling I, Sieber O, Reza FB, Eckert M, Rowan A, Barclay E, Atkin W, Williams C, Gilbert J, et al. Germline APC variants in patients with multiple colorectal adenomas, with evidence for the particular importance of E1317Q. *Human Molecular Genetics* 2000;9:2215–22.
13. Lefevre JH, Rodrigue CM, Mourra N, Bennis M, Flejou J-F, Parc R, Tired E, Gespach C, Parc YR. Implication of MYH in colorectal polyposis. *Annals of Surgery* 2006;244:874–9; discussion 879–80.
14. Grünhage F, Jungck M, Lamberti C, Berg C, Becker U, Schulte-Witte H, Plassmann D, Rahner N, Aretz S, Friedrichs N, Buettner R, Sauerbruch T, et al. Association of familial colorectal cancer with variants in the E-cadherin (CDH1) and cyclin D1 (CCND1) genes. *International Journal of Colorectal Disease* 2008;23:147–54.
15. Pittman AM, Twiss P, Broderick P, Lubbe S, Chandler I, Penegar S, Houlston RS. The CDH1-160C>A polymorphism is a risk factor for colorectal cancer. *International Journal of Cancer* 2009;125:1622–5.
16. Levine AJ, Figueiredo JC, Lee W, Poynter JN, Conti D, Duggan DJ, Campbell PT, Newcomb P, Martinez ME, Hopper JL, Le Marchand L, Baron JA, et al. Genetic variability in the MTHFR gene and colorectal cancer risk using the colorectal cancer family registry. *Cancer Epidemiology, Biomarkers & Prevention* 2010;19:89–100.
17. Le Marchand L, Wilkens LR, Kolonel LN, Henderson BE. The MTHFR C677T polymorphism and colorectal cancer: the multiethnic cohort study. *Cancer Epidemiology, Biomarkers & Prevention* 2005;14:1198–203.
18. Polakova V, Pardini B, Naccarati A, Landi S, Slyskova J, Novotny J, Vodickova L, Bermejo JL, Hanova M, Smerhovsky Z, Tulupova E, Kumar R, et al. Genotype and haplotype analysis of cell cycle genes in sporadic colorectal cancer in the Czech Republic. *Human Mutation* 2009;30:661–8.

19. Ahn SJ, Costa J, Rettig Emanuel J. PicoGreen quantitation of DNA: effective evaluation of samples pre- or post-PCR. *Nucleic Acids Research* 1996;24:2623–5.
20. Sauer S, Gut I. Genotyping single-nucleotide polymorphisms by matrix-assisted laser-desorption/ionization time-of-flight mass spectrometry. *Journal of Chromatography B* 2002;782:73–87.
21. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, Maller J, Sklar P, de Bakker PIW, Daly MJ, Sham PC. PLINK: a tool set for whole-genome association and population-based linkage analyses. *American Journal of Human Genetics* 2007;81:559–75.
22. Mantel N, Haenszel W. Statistical aspects of the analysis of data from retrospective studies of disease. *Journal of the National Cancer Institute* 1959;22:719–48.
23. Svejgaard A, Jersild C, Nielsen LS, Bodmer WF. HL-A antigens and disease. Statistical and genetical considerations. *Tissue Antigens* 1974;4:95–105.
24. Adzhubei IA, Schmidt S, Peshkin L, Ramensky VE, Gerasimova A, Bork P, Kondrashov AS, Sunyaev SR. A method and server for predicting damaging missense mutations. *Nature Methods* 2010;7:248–9.
25. Calabrese R, Capriotti E, Fariselli P, Martelli PL, Casadio R. Functional annotations improve the predictive score of human disease-related mutations in proteins. *Human Mutation* 2009;30:1237–44.
26. Yuan H-Y, Chiou J-J, Tseng W-H, Liu C-H, Liu C-K, Lin Y-J, Wang H-H, Yao A, Chen Y-T, Hsu C-N. FASTSNP: an always up-to-date and extendable service for SNP function analysis and prioritization. *Nucleic Acids Research* 2006;34:W635–41.
27. Lee PH, Shatkay H. F-SNP: computationally predicted functional SNPs for disease association studies. *Nucleic Acids Research* 2008;36:D820–4.
28. Lemire M. Defining rare variants by their frequencies in controls may increase type I error. *Nature Genetics* 2011;43:391–2.
29. Pearson RD. Bias due to selection of rare variants using frequency in controls. *Nature Genetics* 2011;43:392–3.
30. Li LC, Chui RM, Sasaki M, Nakajima K, Perinchery G, Au HC, Nojima D, Carroll P, Dahiya R. A single nucleotide polymorphism in the E-cadherin gene promoter alters transcriptional activities. *Cancer Research* 2000;60:873–6.
31. Luedeke M, Linnert CM, Hofer MD, Surowy HM, Rinckleb AE, Hoegel J, Kuefer R, Rubin MA, Vogel W, Maier C. Predisposition for TMPRSS2-ERG fusion in prostate cancer by variants in DNA repair genes. *Cancer Epidemiology, Biomarkers & Prevention* 2009;18:3030–5.
32. Kokko A, Laiho P, Lehtonen R, Korja S, Carvajal-Carmona LG, Järvinen H, Mecklin J-P, Eng C, Schleutker J, Tomlinson IPM, Vahteristo P, Aaltonen LA. EPHB2 germline

- variants in patients with colorectal cancer or hyperplastic polyposis. *BMC Cancer* 2006;6:145.
33. Thusberg J, Olatubosun A, Vihinen M. Performance of mutation pathogenicity prediction methods on missense variants. *Human Mutation* 2011;32:358–68.
 34. Spurdle AB. Clinical relevance of rare germline sequence variants in cancer genes: evolution and application of classification models. *Current Opinion in Genetics & Development* 2010;
 35. Bonilla C, Lefèvre JH, Winney B, Johnstone E, Tonks S, Colas C, Day T, Hutnik K, Boumertit A, Midgley R, Kerr D, Parc Y, et al. Cyclin D1 rare variants in UK multiple adenoma and early-onset colorectal cancer patients. *Journal of Human Genetics* 2011;
 36. Zhu Q, Ge D, Maia JM, Zhu M, Petrovski S, Dickson SP, Heinzen EL, Shianna KV, Goldstein DB. A genome-wide comparison of the functional properties of rare and common genetic variants in humans. *American Journal of Human Genetics* 2011;88:458–68.
 37. Humar B, Graziano F, Cascinu S, Catalano V, Ruzzo AM, Magnani M, Toro T, Burchill T, Futschik ME, Merriman T, Guilford P. Association of CDH1 haplotypes with susceptibility to sporadic diffuse gastric cancer. *Oncogene* 2002;21:8192–5.
 38. Qiu L-X, Li R-T, Zhang J-B, Zhong W-Z, Bai J-L, Liu B-R, Zheng M-H, Qian X-P. The E-cadherin (CDH1)--160 C/A polymorphism and prostate cancer risk: a meta-analysis. *European Journal of Human Genetics* 2009;17:244–9.
 39. Beeghly-Fadiel A, Lu W, Gao Y-T, Long J, Deming SL, Cai Q, Zheng Y, Shu X-O, Zheng W. E-cadherin polymorphisms and breast cancer susceptibility: a report from the Shanghai Breast Cancer Study. *Breast Cancer Research and Treatment* 2009;

Table 1. Case and control sample description.

	N	mean age (years)	male:female	mean no. of polyps
UK multiple adenomas	112	59 ^a	68:20 ^b	11 ^b
UK early onset	70	42	38:31 ^c	n/a
French multiple adenomas	75	51 ^d	44:31	26 ^e
French early onset	56	40	24:32	n/a
PoBI controls	866	62	478:382 ^f	n/a

Missing data for: ^a33, ^b24, ^c1, ^d3, ^e14, ^f6 individuals. n/a=not applicable.

Table 2. Variants analyzed in UK cases and controls

id	variant	dbSNP	major/ minor allele	MAF cases	MAF controls	p-value ^a	MAF multiple adenomas	MAF early onset	p- value ^b	PolyPhen-2/FastSNP	SNPs&GO
MTHFR-1	A222V	rs1801133	C/T	0.338	0.336	0.95	0.337	0.342	0.94	probably damaging	disease
EPHB2-1	R80H	n/a	G/A	0.003	0.001	0.19	0.005	0.000	0.53	probably damaging	disease
EPHB2-3	I361V	rs56180036	A/G	0.007	0.001	0.07	0.010	0.000	0.37	benign	neutral
EPHB2-4	R568W	n/a	C/T	0.003	0.001	0.19	0.000	0.012	0.11	probably damaging	disease
EPHB2-7	M883V	n/a	A/G	0.003	0.000	0.22	0.005	0.000	0.53	possibly damaging	neutral
EXO1-2	E109K	n/a	G/A	0.000	0.001	0.66	0.000	0.000	n/a	possibly damaging	neutral
EXO1-12	D249N	rs61750993	G/A	0.020	0.007	0.03	0.018	0.024	0.74	probably damaging	neutral
EXO1-4	L410R	n/a	T/G	0.004	0.000	0.18	0.005	0.000	0.54	probably damaging	neutral
EXO1-10	G759E	rs4150001	G/A	0.003	0.009	0.38	0.000	0.013	0.11	benign	neutral
MLH1-1	G22A	rs41295280	G/C	0.003	0.000	0.03	0.005	0.000	0.53	probably damaging	neutral
CTNNB1-1	N287S	rs35288908	A/G	0.003	0.000	0.05	0.005	0.000	0.42	benign	disease
APC-7	L1129S	n/a	T/C	0.000	0.003	0.37	0.000	0.000	n/a	possibly damaging	disease
APC-11	E1317Q	rs1801166	G/C	0.011	0.007	0.46	0.014	0.007	0.59	benign	neutral
APC-15	G2502S	rs2229995	G/A	0.010	0.021	0.19	0.009	0.012	0.83	benign	neutral
APC-16	R2505Q	n/a	G/A	0.000	0.002	0.43	0.000	0.000	n/a	probably damaging	neutral
APC-17	S2621C	rs72541816	C/G	0.003	0.011	0.25	0.000	0.012	0.12	benign	neutral
APC-20	8636 C/A	n/a	C/A	0.024	0.019	0.66	0.024	0.025	0.96	n/a	n/a
PMS2-1	T511A	rs2228007	A/G	0.023	0.029	0.60	0.023	0.024	0.55	benign	neutral
PMS2-2	T597S	rs1805318	A/T	0.034	0.020	0.23	0.043	0.012	0.19	benign	neutral
PMS2-3	M622I	rs1805324	G/A	0.017	0.022	0.59	0.014	0.024	0.98	benign	neutral
BRCA2-7	N372H	rs144848	T/G	0.320	0.286	0.23	0.329	0.298	0.60	benign	disease
BRCA2-8	S384F	rs41293475	C/T	0.000	0.003	0.38	0.000	0.000	n/a	possibly damaging	disease
BRCA2-27	R2034C	rs1799954	C/T	0.004	0.004	0.90	0.000	0.012	0.12	benign	disease
BRCA2-35	D2665G	rs28897745	A/G	0.003	0.000	0.22	0.005	0.000	0.53	probably damaging	disease

BRCA2-37	V2728I	rs28897749	G/A	0.014	0.003	0.02	0.019	0.000	0.21	benign	disease
MLH3-1	V741F	rs28756990	G/T	0.014	0.004	0.16	0.005	0.036	0.04	benign	neutral
MLH3-2	M809V	rs61752722	A/G	0.003	0.009	0.37	0.005	0.000	0.53	benign	neutral
MLH3-5	S845G	rs28756992	A/G	0.014	0.033	0.09	0.005	0.036	0.04	benign	neutral
AXIN1-4	D495E	n/a	C/G	0.004	0.009	0.44	0.007	0.000	0.37	benign	neutral
AXIN1-6	R841Q	rs34015754	G/A	0.007	0.007	0.92	0.010	0.000	0.45	probably damaging	disease
CDH1-1	-1128	rs13335980	A/T	0.000	0.001	0.53	0.000	0.000	n/a	-S8 binding site	n/a
CDH1-2	-284	rs16260	C/A	0.246	0.308	0.07	0.218	0.291	0.12	affects TF binding sites	n/a
CDH1-3	IVS1+6	rs3743674	T/C	0.102	0.121	0.37	0.116	0.071	0.27	splicing site (medium-high risk)	n/a
CDH1-5	IVS4+10	rs33963999	G/C	0.070	0.080	0.57	0.051	0.119	0.04	intronic enhancer (very low-low risk)	n/a
CDH1-7	A592T	rs35187787	G/A	0.003	0.007	0.45	0.005	0.000	0.53	benign	neutral
CDH1-8	T599S	n/a	C/G	0.000	0.001	0.65	0.000	0.000	n/a	benign	neutral
CDH1-10	A634V	n/a	C/T	0.000	0.001	0.52	0.000	0.000	n/a	possibly damaging	disease
TP53-1	R72P	rs1042522	G/C	0.255	0.248	0.80	0.270	0.220	0.38	benign	disease
BRCA1-6	Q356R	rs1799950	A/G	0.042	0.060	0.22	0.043	0.038	0.53	possibly damaging	disease
BRCA1-8	R496H	rs28897677	G/A	0.000	0.001	0.65	0.000	0.000	n/a	benign	disease
BRCA1-16	T826K	rs28897683	C/A	0.000	0.001	0.65	0.000	0.000	n/a	probably damaging	disease
BRCA1-22	E1038G	rs16941	A/G	0.338	0.342	0.90	0.346	0.317	0.64	probably damaging	disease
BRCA1-28	S1512I	rs1800744	G/T	0.003	0.004	0.82	0.005	0.000	0.86	benign	disease
AXIN2-1	N412S	rs115931022	A/G	0.000	0.007	0.19	0.000	0.000	n/a	benign	neutral
SMAD4-1	A118A	n/a	G/A	0.003	0.005	0.75	0.005	0.000	0.53	synonymous	n/a
CHEK2-1	1100 delC	n/a	C/del	0.008	0.003	0.31	0.009	0.007	0.88	p.T367fsX15	n/a

n/a : not available.

MAF: minor allele frequency.

number in bold : $p \leq 0.05$.

^ap-value for the comparison case vs control

^bp-value for the comparison multiple adenoma vs early onset CRC

Table 3. Rare variant counts in UK cases and controls.

rare variant	cases ^a	controls ^a	p-value
EPHB2-1	1/149	1/778	0.30
EPHB2-3	1/145	2/746	0.41
EPHB2-4	1/149	1/775	0.30
EPHB2-7	1/149	0/224	0.40
EXO1-2	0/150	1/751	1.00
EXO1-4	1/125	0/226	0.36
EXO1-10	1/145	4/225	0.65
EXO1-12	5/150	10/745	0.15
MLH1-1	1/147	0/740	0.17
CTNNB1-1	1/174	0/702	0.20
APC-7	0/149	4/743	1.00
APC-11	4/176	11/745	0.50
APC-16	0/150	3/729	1.00
BRCA2-8	0/147	4/748	1.00
BRCA2-27	1/142	6/740	1.00
BRCA2-35	1/149	0/224	0.40
BRCA2-37	4/148	5/744	0.05
MLH3-1	4/144	2/225	0.21
MLH3-2	1/148	4/227	0.65
AXIN1-4	1/117	14/749	0.71
AXIN1-6	2/146	11/746	1.00
CDH1-1	0/150	2/748	1.00
CDH1-7	1/147	11/747	0.70
CDH1-8	0/151	1/742	1.00
CDH1-10	0/150	2/731	1.00
BRCA1-8	0/151	1/735	1.00
BRCA1-16	0/149	1/744	1.00
BRCA1-28	1/149	2/225	1.00
AXIN2-1	0/128	10/737	0.37
SMAD4-1	1/148	7/739	1.00
CHEK2-1	3/178	1/179	0.37
MAF < 1% (n=31)			
total carriers/total non-carriers 1 ^b	24/146	108/739	
OR 1 (95% CI) ^c	1.13 (0.70, 1.81)		0.63
total carriers/total non-carriers 2 ^b	13/146	13/217	
OR 2 (95% CI) ^c	1.49 (0.68, 3.24)		0.33
combined OR (95% CI)	1.21 (0.80, 1.82)		0.42
MAF < 0.5% (n=23)			
total carriers/total non-carriers 1 ^b	11/149	41/740	
OR 1 (95% CI) ^c	1.33 (0.67, 2.65)		0.41
total carriers/total non-carriers 2 ^b	11/146	5/215	
OR 2 (95% CI) ^c	3.24 (1.10, 9.52)		0.04

combined OR (95% CI)	1.77 (0.97, 3.08)		0.05
----------------------	-------------------	--	-------------

number in bold : $p \leq 0.05$.

^anumber of individuals with variant/total number of individuals typed.

^bnumber of non-carriers corresponds to the harmonic mean of individuals without the rare variant for each variant typed.

^cOR 1 was calculated using the larger set of controls whereas OR 2 was calculated using a subset of 227 controls.

Table 4. Rare variant counts in UK multiple adenoma, early onset CRC and control subjects for variants with MAF<0.5% in controls.

rare variant	multiple adenomas ^a	early onset ^a	controls ^a	p-value ^b	p-value ^c
EPHB2-1	1/107	0/42	1/778	0.23	1.00
EPHB2-3	1/103	0/42	2/746	0.32	1.00
EPHB2-4	0/107	1/42	1/775	1.00	0.10
EPHB2-7	1/107	0/42	0/224	0.32	1.00
EXO1-2	0/108	0/42	1/751	1.00	1.00
EXO1-4	1/91	0/34	0/226	0.29	1.00
MLH1-1	1/105	0/42	0/740	0.12	1.00
CTNNB1-1	1/106	0/68	0/702	0.13	1.00
APC-7	0/107	0/42	4/743	1.00	1.00
APC-16	0/108	0/42	3/729	1.00	1.00
BRCA2-8	0/105	0/42	4/748	1.00	1.00
BRCA2-27	0/101	1/41	6/740	1.00	0.05
BRCA2-35	1/107	0/42	0/224	0.32	1.00
BRCA2-37	4/106	0/42	5/744	0.02	1.00
MLH3-1	1/102	3/42	2/225	1.00	0.03
CDH1-1	0/108	0/42	2/748	1.00	1.00
CDH1-8	0/109	0/42	1/742	1.00	1.00
CDH1-10	0/109	0/41	2/731	1.00	1.00
BRCA1-8	0/109	0/42	1/744	1.00	1.00
BRCA1-16	0/107	0/42	1/734	1.00	1.00
BRCA1-28	1/107	0/42	2/225	1.00	1.00
SMAD4-1	1/106	0/42	7/739	1.00	1.00
CHEK2-1	2/111	1/67	1/179	0.56	0.47
adenomas vs controls					
total carriers/ total non-carriers 1 ^d	9/106		41/740		
OR 1 (95% CI) ^e	1.53 (0.72, 3.24)			0.26	
total carriers/	7/103		5/215		

total non-carriers 2 ^d					
OR 2 (95% CI) ^e	2.92 (0.91, 9.43)			0.06	
combined OR (95% CI)	1.87 (0.98, 3.48)			0.05	
early onset vs controls					
total carriers/ total non-carriers 1 ^d		2/43	41/740		
OR 1 (95% CI) ^e		0.84 (0.20, 3.59)		0.81	
total carriers/ total non-carriers 2 ^d		4/42	5/215		
OR 2 (95% CI) ^e		4.10 (1.06, 15.89)		0.03	
combined OR (95% CI)		1.72 (0.73, 5.27)		0.25	

number in bold : $p \leq 0.05$.

^anumber of individuals with variant/total number of individuals typed.

^bp-value for the comparison between multiple adenoma patients and controls.

^cp-value for the comparison between early onset CRC patients and controls.

^dnumber of non-carriers corresponds to the harmonic mean of individuals without the rare variant for each variant typed.

^eOR 1 was calculated using the larger set of controls whereas OR 2 was calculated using a subset of 227 controls.

Table 5. Variants genotyped in French patients.

id	variant	major/minor allele	French cases ^a (N=131)	UK cases ^a (N=182)	class
EPHB2-1	R80H	G/A	0.000	0.003	rare variant
EPHB2-4	R569W	C/T	0.000	0.003	rare variant
EPHB2-7	M883V	A/G	0.000	0.003	rare variant
EXO1-4	L410R	T/G	0.000	0.004	rare variant
EXO1-10	G759E	G/A	0.015	0.003	rare variant
MSH2-8	E808X	G/T	0.000	0.000	rare variant
CTNNB1-1	N287S	A/G	0.000	0.003	rare variant
APC-10	I1307K	T/A	0.000	0.000	rare variant
APC-11	E1317Q	G/C	0.004	0.011	rare variant
APC-17	S2621C	C/G	0.008	0.003	low freq variant
APC-20	8636 C/A	C/A	0.051	0.024	low freq variant
PMS2-2	T597S	A/T	0.008	0.034	low freq variant
BRCA2-7	N372H	T/G	0.273	0.320	polymorphism
BRCA2-35	D2665G	A/G	0.000	0.003	rare variant
BRCA2-37	V2728I	G/A	0.000	0.014	rare variant
BRCA2-48	P3194Q	C/A	0.000	0.000	rare variant
MLH3-1	V741F	G/T	0.016	0.014	rare variant

MLH3-2	M809V	A/G	0.012	0.003	rare variant
MLH3-5	S845G	A/G	0.042	0.014	low freq variant
CDH1-2	-284	C/A	0.315	0.246	polymorphism
BRCA1-28	S1512I	G/A	0.004	0.003	rare variant
CHEK2-1	1100 delC	C/del	0.000	0.008	rare variant

^aMinor allele frequency is shown.

Supplementary Table 4. Rare variant counts in UK and French patients.

rare variant id	UK cases ^a	French cases ^a
EPHB2-1	1/149	0/130
EPHB2-4	1/149	0/124
EPHB2-7	1/149	0/129
EXO1-4	1/125	0/130
EXO1-10	1/145	4/130
CTNNB1-1	1/174	0/130
APC-11	4/176	1/131
BRCA2-35	1/149	0/130
BRCA2-37	4/148	0/131
MLH3-2	1/148	3/129
MLH3-1	4/144	4/127
BRCA1-28	1/149	1/129
CHEK2-1	3/178	0/131
MAF < 1% (n=13)		
total carriers/total non-carriers ^b	24/149	13/128
OR (95% CI)	1.59 (0.71, 3.02)	0.20
MAF < 0.5% (n=10)		
total carriers/total non-carriers ^b	18/148	5/128
OR (95% CI) ^e	3.14 (1.13, 8.69)	0.02

number in bold : $p \leq 0.05$.

^anumber of individuals with variant/total number of individuals typed.

^bnumber of non-carriers corresponds to the harmonic mean of individuals without the rare variant for each variant typed.