

Published in final edited form as:

Nat Genet. 2009 December ; 41(12): 1330–1334. doi:10.1038/ng.483.

## Genome-wide association study of ulcerative colitis identifies three new susceptibility loci, including the *HNF4A* region

The UK IBD Genetics Consortium and The Wellcome Trust Case Control Consortium 2

### Abstract

Ulcerative colitis (UC) is a common form of inflammatory bowel disease with a complex aetiology. As part of the Wellcome Trust Case Control Consortium 2, we performed a genome-wide association scan for UC in 2361 cases and 5417 controls. Loci showing evidence of association at  $P < 1 \times 10^{-5}$  were followed up by genotyping in an independent set of 2321 cases and 4818 controls. We find genome-wide significant evidence of association at three new loci, each containing at least one biologically relevant candidate gene, on chromosomes 20q13 (*HNF4A*;  $P = 3.2 \times 10^{-17}$ ), 16q22 (*CDHI* and *CDH3*;  $P = 2.8 \times 10^{-8}$ ) and 7q31 (*LAMB1*;  $3.0 \times 10^{-8}$ ). Of note, *CDHI* has recently been associated with susceptibility to colorectal cancer, which is an established complication of longstanding UC. The new associations suggest that changes in the integrity of the intestinal epithelial barrier may contribute to the pathogenesis of UC.

Genetic epidemiological data clearly implicate inherited susceptibility in the pathogenesis of UC and Crohn's disease (CD), which represent the two common forms of inflammatory bowel disease (IBD) and together affect at least 1 in 250 of the Northern European population.<sup>1</sup> Notwithstanding recent therapeutic advances, disease-related morbidity in ulcerative colitis continues to be high. Recognized complications of severe disease refractory to medical therapy include colectomy, often as an emergency, in 15–20% of patients, as well as colorectal cancer.<sup>2</sup>

Substantial progress has been made in understanding IBD pathogenesis in recent years. In genetically susceptible individuals it appears that a dysregulated mucosal immune response to commensal enteric bacteria predisposes to chronic, relapsing intestinal inflammation which is the hallmark of IBD.<sup>3</sup> Clinical features combined with epidemiological evidence have long suggested that CD and UC are related polygenic diseases. This has recently been corroborated by the results of genetic association studies, which have highlighted both disease-specific loci and others which are shared between UC and CD. For example, while genetically-determined defects in the handling of intracellular bacteria (*NOD2* and the autophagy genes *ATG16L1* and *IRGM*) are specific to CD, multiple components in the Th17 pathway (*IL23R*, *IL12B*, *JAK2*, *STAT3*) are associated with both CD and UC.<sup>4–12</sup>

Until recently most attention had focused on CD, with genome-wide association (GWA) studies and subsequent meta-analysis yielding more than 30 confirmed CD susceptibility loci.<sup>4, 6, 7, 10–12</sup> In addition to the longstanding known association in the MHC,<sup>13</sup> the first

Correspondence to: JCB (barrett@sanger.ac.uk), CCAS (chris.spencer@well.ox.ac.uk).

#### Author contributions

JL, CL, NP, AP, EW, KP, HZ, HD, ERN, DM, KB, TE, LC were involved in establishing DNA collections, and/or assembling phenotypic data; CL, AP, EW, DM, HD, AL, CM, JeS, DPJ, CE, TA, JCM, JS, MP recruited patients; WN, CE, TA, JCM, JS, MP, CGM supervised clinical and laboratory work; WTCCC2 DNA, Genotyping, Data QC and Informatics group executed GWAS sample handling, genotyping and QC; WTCCC2 Data and Analysis group, JCB, CAA performed statistical analyses; JCB, JL, CL, CCAS, CAA, TA, PD, JS, MP, CGM contributed to writing the manuscript. WTCCC2 Management Committee conceived and oversaw the design and execution of the GWAS. WTCCC2 group memberships are specified in the full author list.

GWA scans in UC reported associations at *IL23R*, *IL10* and loci on chromosomes 1p36 and 12q15 which meet accepted genome-wide significance thresholds.<sup>14, 15</sup>

As part of the Wellcome Trust Case Control Consortium 2 (WTCCC2) study of 15 complex disorders and traits, we report here the results of the largest GWA scan in UC to date. All study subjects were UK residents of white, European ancestry; clinical data are presented in Table 1. Cases and controls were genotyped on the Affymetrix 6.0 array. After application of quality control filters (see Methods), we analysed GWA data from 2361 individuals with UC and 5417 controls (Figure 1). An initial analysis revealed 24 distinct loci (comprising 156 SNPs) which showed evidence of association at  $P < 1 \times 10^{-5}$ . Sixteen of these had not been previously reported, and were followed up by genotyping the most strongly associated SNP from each locus using the Sequenom iPLEX platform in an independent panel of 2321 UC cases and 4818 controls. Three new loci showed evidence for association at  $P < 5 \times 10^{-8}$  in the combined panel, with three further new loci showing nominal ( $P < 0.05$ ) replication (Table 2 and Figure 2). We describe these loci below and highlight the most plausible candidate gene for each, recognizing that fine mapping and functional studies are required to define causal variants and identify the gene from which each signal arises. A list of all loci for which replication was attempted is shown in Supplementary Table 1.

The most significant new association was seen at rs6017342 (GWA scan  $P = 3.2 \times 10^{-13}$ ; combined GWA and replication  $P = 8.5 \times 10^{-17}$ ), which maps within a recombination hotspot on chromosome 20q13 containing the 3' untranslated region (UTR) of just one gene, *HNF4A*. The SNP rs6017342 itself maps 5kb distal to the 3' UTR. Although within an expressed sequence tag DB076868, this has been detected in just a single testis cDNA library and does not encode a significant open reading frame. The region contains two small blocks of sequence that are conserved in mammals and may include regulatory sequences affecting the expression of surrounding genes. Since rs6017342 is located within a recombination hotspot, there are few known SNPs in strong linkage disequilibrium ( $r^2 > 0.5$ ) with it; there are none on the Affymetrix chip used in this study or on the Illumina chips used in previous studies. As the evidence for this association rests on this single SNP, we subjected these data to careful scrutiny; genotype cluster plots for this SNP showed clear resolution of the 3 genotype classes (Supplementary Figure 1), with 99.3% completeness of genotypes within this dataset.

Rare *HNF4A* mutations account for approximately 4% of UK cases of maturity-onset diabetes of the young (MODY),<sup>16</sup> a monogenic form of diabetes mellitus characterized by autosomal dominant inheritance, young age of onset, pancreatic b-cell dysfunction and sensitivity to sulphonylureas. Common variants of *HNF4A* influence predisposition to Type II diabetes (rs2144908)<sup>17</sup> and dyslipidaemia (rs1800961).<sup>18</sup> The UC associated SNP, rs6017342 is not in LD with either of these 2 common variants, nor did it not show association in our study of CD ( $P=0.92$ ).<sup>3</sup>

*HNF4A* encodes the transcription factor hepatocyte nuclear factor 4  $\alpha$  which regulates the expression of multiple components within all three key compartments of the cell-cell junction, namely the adherens junction, the tight junction and the desmosome.<sup>19</sup> Such cell-cell junctions are fundamental to epithelial organization and barrier function. HNF4 $\alpha$  also plays a key role in the development of the embryonic mammalian gastrointestinal tract. Previous studies demonstrated that mice with targeted deletion of *HNF4a* in epithelial cells of the foetal colon die perinatally. Histological analysis of colonic tissue recovered during late development (E18.5) demonstrated absent crypt formation, reduced epithelial cell proliferation and defective goblet cell maturation.<sup>20</sup> In order to explore the role of *HNF4a* in murine intestinal inflammation, Ahn and colleagues circumnavigated the embryonic lethality of *Hnf4a*<sup>-/-</sup> mice by generating a conditional model of intestinal *Hnf4a* deletion.

21 These *Hnf4a*<sup>ΔIEpC</sup> mice (floxed *Hnf4a* driven by the *villin* promoter) developed increased epithelial permeability and a markedly more severe colitis following dextran sodium sulphate (DSS) challenge, than their wild-type littermates.<sup>21</sup> The same investigators provided preliminary evidence for dysregulated *HNF4A* gene expression in the intestinal epithelium in Crohn's disease and in ulcerative colitis,<sup>21</sup> a finding which now merits detailed re-exploration.

Significant association was also seen for a locus on chromosome 16q22, with the strongest signal at rs1728785 (GWA scan  $P = 1.8 \times 10^{-5}$ ; combined GWA and replication  $P = 2.8 \times 10^{-8}$ ). The interval bounded by recombination hotspots spans 411 kb and encodes several genes. Among the strongest candidates for UC susceptibility is *CDHI* which encodes E-cadherin. This transmembrane glycoprotein is one of the main components of the adherens junction and a key mediator of intercellular adhesion in the intestinal epithelium. It also plays a key role in epithelial restitution and repair following mucosal damage and expression of *CDHI* is known to be significantly reduced in areas of active UC.<sup>22</sup>

Given the well-recognised association between UC and colorectal cancer,<sup>2</sup> the observation of correlated association signals at the *CDHI* locus in both diseases is striking. Thus variants in LD ( $r^2 = 0.5$ ) with the most strongly UC associated SNP in our study were recently identified in a GWA scan meta-analysis to be associated with colorectal cancer susceptibility<sup>23</sup>; conversely, we find that a perfect proxy for the most associated SNP in the colorectal cancer study is also associated with UC ( $P = 8 \times 10^{-4}$ ). This locus did not show association with CD in a large international GWA meta-analysis of CD ( $P = 0.549$ )<sup>6</sup> (Supplementary Table 2). However, evidence for association of *CDHI* with CD was reported recently in the Canadian population using a candidate gene approach,<sup>24</sup> and the CD associated SNPs resulted in a truncated E-cadherin protein *in vitro* which accumulated in the cytoplasm and led to disorganized epithelial architecture.<sup>24</sup>

Of great potential relevance is the evidence that *HNF4A* and E-cadherin co-operate to maintain epithelial barrier integrity in the intestine. In experiments focused on the liver, *HNF4a* knockout mice failed to express E-cadherin,<sup>19</sup> while in the gut E-cadherin dependent cell-cell contact was found to be critical in determining the amount and binding activity of nuclear *HNF4a*. This in turn affected the expression of several genes including *ApoA-IV*,<sup>25</sup> an anti-inflammatory protein known to inhibit experimental colitis.<sup>26</sup>

The third newly confirmed UC susceptibility locus was a region on chromosome 7q31, previously suggested by a recent North American GWA scan.<sup>14</sup> In the current study the peak association was seen at rs886774 (GWA scan  $P = 4.8 \times 10^{-7}$ ; combined GWA and replication  $P = 3.0 \times 10^{-8}$ ). A strong positional candidate gene at this locus is *LAMB1*, encoding the laminin beta 1 subunit. Laminins are heterotrimers; the beta-1 light chain is present in laminins-1 -2 and -10. Laminins are expressed in the intestinal basement membrane, and play a key role in anchoring the single-layered epithelium; expression is known to be down-regulated in UC.<sup>27</sup> rs886774 was not associated with CD in the meta-analysis.<sup>5</sup> (Supplementary Table 2)

Two other loci previously implicated in UC-related phenotypes showed strong (but not genome-wide significant) association with UC. These comprise a SNP previously associated with osteoporosis<sup>28</sup> (rs7524102 on chromosome 1p36, combined GWA and replication  $P = 3.1 \times 10^{-7}$ ) and a SNP nearby (though not in LD with) a marker known to be associated with psoriasis<sup>29</sup> (rs9548988 on 13q.13, combined GWA and replication  $P = 2.7 \times 10^{-7}$ ).

In addition to the novel loci described above, our GWAS detected strong association at established UC loci such as the MHC, *IL23R*, 3p21/*MST1* and *NKX2-3* (one tailed  $P$  values in the direction of the previously reported association in Table 3). We also provide robust

confirmation of two UC loci reported recently in genome-wide scans, the *IL10* locus<sup>11</sup> and the *OTUD3/PLA2G2E* locus<sup>12</sup> on chromosome 1q31 and 1p36 respectively. Also of interest is our finding that the *PSMG1* locus on chromosome 21, which has previously been associated with pediatric-onset IBD,<sup>30</sup> is likely to contribute specifically to disease susceptibility in UC. Variable degrees of support were obtained for some previously reported UC loci, including *ECM1*, *CARD9*,<sup>31</sup> *KIF21B*/chromosome 1q32, and *JAK2*/chromosome 9p24, but weaker support for other loci such as *IL2/IL21*,<sup>32</sup> *IL12B* and 12q15 (Table 3). Some of the UC loci are clearly associated with CD, while others are not, or have not been tested (Supplementary Table 2). We also tested for epistatic interaction among all pairwise combinations of these loci (both previously described and new) but found none.

This is the first report of a new series of GWA scans undertaken by the WTCCC2 consortium. We have identified three new susceptibility loci for UC, and provide the first genetic link between UC and colorectal cancer. Each of the strongest new association intervals that we have identified contains respectively *HNF4A*, *CDH1* and *LAMB1* as the most plausible positional candidate genes, thus providing further evidence for the re-emerging concept that altered epithelial barrier function may be a key factor in UC pathogenesis.<sup>8</sup> Indeed, this is the first time that variants within genetic loci encoding such epithelial barrier genes have shown association with IBD at stringent genome-wide significant thresholds. Fine mapping and functional studies are clearly required to investigate this connection further, but our study provides strong scientific justification for the exploration of new therapeutic targets relevant to epithelial barrier function.

## METHODS

### Subjects

**Cases**—A total of 5319 unrelated patients of white, European, non-Jewish ancestry with a diagnosis of ulcerative colitis established using standard endoscopic, radiological and histological criteria, were recruited from ten centres within the United Kingdom (Cambridge, Oxford, London, Newcastle, Sheffield, Edinburgh, Dundee, Manchester, Torbay and Exeter, Supplementary Table 4). All patients provided written consent and either a sample of blood or saliva, from which DNA was extracted according to standard protocols. Research Ethics Committee approval was obtained prior to sample collection (Cambridge, Oxford, London, Newcastle, Sheffield, Edinburgh, Dundee, Manchester, Torbay and Exeter Local Research Ethics Committees). After QC (see below), we analyzed a total of 4682 samples, which were divided between the discovery panel (2361 samples) and replication panel (2321 samples).

**Controls**—A total of 10,235 control DNA samples from 3 sources passed our QC filters (see below). 5417 samples of the WTCCC2 common control set were used for the GWA experiment. This comprised 2675 healthy blood donors recruited from the United Kingdom Blood Service (UKBS), and 2742 samples from the 1958 Birth Cohort (1958BC) obtained from EBV-transformed cell lines from individuals born in England, Wales and Scotland during one week in 1958. The 4818 samples used as controls for the replication cohort were recruited from the Wellcome Trust-funded People of the British Isles (PoBI) DNA collection, obtained from rural populations throughout the British Isles, and from a further independent set of DNA samples obtained from 1958BC. All of the control samples used were from individuals with self-reported Caucasian ethnicity.

A summary of patients and controls is shown in Table 1 and Supplementary Table 4.

**DNA sample preparation:** Genomic DNA for all cases was shipped to the Sanger Institute, Cambridge. DNA quality plus subject identity were validated using the Sequenom iPLEX

assay designed to genotype 4 gender SNPs and 26 SNPs present on the Affymetrix array. DNA concentrations were quantified using a PicoGreen assay (Invitrogen) and an aliquot assayed by agarose gel electrophoresis. A DNA sample was considered to pass quality control if the original DNA concentration was  $\geq 50$  ng/ul, the DNA was not degraded, the gender assignment from the iPLEX assay matched that provided in the patient data manifest and genotypes were obtained for over 65% of the SNPs on the iPLEX.

### GWA Genotyping

Samples were genotyped at Affymetrix's service laboratory on the Genome-Wide Human SNP Array 6.0. For all samples passing Affymetrix's laboratory QC, raw intensities (from the .CEL files) were renormalized within collections using *CelQuantileNorm* (see <http://outmodedbonsai.sourceforge.net/>). These normalized intensities were used to call genotypes with an updated version of the *Chiamo* software (see [www.stats.ox.ac.uk/~marchini/software/gwas/chiamo.html](http://www.stats.ox.ac.uk/~marchini/software/gwas/chiamo.html)), adapted for Affymetrix 6.0 SNP data. The *Chiamo* algorithm simultaneously calls genotypes for individuals in several collections; here it was applied to 15,068 individuals from five collections genotyped as part of the WTCCC2. *Chiamo* generates posterior probabilities for each of the three possible genotypes plus a fourth class of outliers. Our analyses use thresholded genotypes: for each individual, if one genotype had posterior probability greater than 0.9, this was set as the genotype for that individual, otherwise the genotype was set to be missing. After applying the QC filters described below, this threshold led to a study-wide level of missing data of 0.20%.

An overlapping set of 4830 controls were also genotyped on the Illumina 1.2M chip as part of a separate WTCCC2 project, and the 50,000 SNPs which are shared between that platform and the Affymetrix 6.0 (used in this study) were used to evaluate genotype accuracy. For the same QC thresholds and similar levels of missing data, discordance between *Chiamo* and *Illuminus*, which we regard as an upper bound on genotyping error rate, was 0.05857% for 1958BC and 0.07476% for UKBS.

We compared *Chiamo* to *Birdsuite* (the default Affymetrix calling algorithm applied on a plate-by-plate basis as recommended in 33) by making genotype calls at different confidence thresholds, and then plotting the fraction of calls made against concordance with the Illumina genotypes (Supplementary Figure 2). The general trend is that, when matched for the proportion of missing data, *Chiamo* has slightly higher concordance than *Birdsuite*. We are therefore confident that *Chiamo* is an acceptable alternative to *Birdsuite*.

### Replication Genotyping

In the replication stage, genotyping was carried out at the Sanger Institute using the Sequenom iPLEX Gold assay. For one locus, the most associated SNP could not be genotyped with this technology, so a perfect ( $r^2 = 1$  in all HapMap populations) proxy was used instead. 19 SNPs (including 3 gender markers) were typed in a multiplex reaction; 15 passed experimental QC (one SNP with Hardy-Weinberg P value  $< 1 \times 10^{-6}$  was discarded). Samples with  $> 20\%$  missing genotypes ( $n=300$ ) were excluded; these samples are not included in the tallies in Table 1.

### Quality Control

**Samples**—As is now standard practice for GWAS studies, we excluded sets of individuals whose genome-wide patterns of diversity are outliers compared to the bulk of those in the study, and SNPs where there is evidence that genotype calls do not provide precise estimates of genotype frequencies. Ignoring individuals and SNPs in this way throws away data gained at some expense, but because they typically violate assumptions underpinning

standard tests for association, the payback in terms of increased accuracy for these tests can be substantial.

In order to try to obtain the maximally powerful set of samples and SNPs we attempted to refine some standard QC practices. For all individuals we explicitly model the data as a mixture of “normal” and “outlier” individuals for each of ancestry, missing data and heterozygosity, and sex assignment.<sup>34</sup> We fit each model in a Bayesian framework, and exclude individuals whose posterior probability of belonging to the outlier class was above 0.5. This approach replaces (and we believe improves upon) the traditional concept of fixed exclusion thresholds for parameters such as call rate, heterozygosity and ancestry. In total 413 case individuals and 567 control individuals were excluded from the analyses (Supplementary Table 3).

To assess relatedness amongst study individuals we compared each individual with the 100 individuals they were most closely related to (on the basis of genome-wide levels of allele sharing) and used a hidden Markov Model (HMM) to decide, at each position in their genome, whether the two individuals shared 0, 1, or 2 chromosomes identical by descent. This allows more refined assessment of the relatedness between individuals than do genome-wide sharing statistics (for example, parent-child relationships can be distinguished from siblings). We obtained a set of individuals with IBD < 5% by iteratively removing the member of each pair of putatively related individuals with more missing genotypes.

**SNPs**—For each SNP we considered a measure of the (Fisher) information carried by the genotype calls for the underlying allele frequency. Informally, this will decrease as the number of individuals with low posterior probabilities for the most likely call increases, and it can be thought of as a more refined measure of both missing data levels and minor allele frequency (Supplementary Figure 3). The measure is calculated automatically by the program SNPtest ([www.stats.ox.ac.uk/~marchini/software/gwas/snptest.html](http://www.stats.ox.ac.uk/~marchini/software/gwas/snptest.html)). SNPs were removed if this information measure was below 0.98, or if the estimated MAF was below 0.01% (both calculated on the combined case-control data). 14.7% of SNPs were removed by these criteria. Again, this approach appears to offer advantages over conventional SNP filters, in excluding fewer SNPs for the same level of improved data quality. Because associated SNPs are expected to be enriched in the tiny fraction of poorly performing markers on these chips, we subsequently examined 155 cluster plots for SNPs with  $p < 1 \times 10^{-5}$ , and excluded 16 from further analysis as likely genotyping errors.

Supplementary Figure 4 provides QQ plots for the post-QC comparison of our two control collections, and for association statistics based on the post-QC trend test comparing cases and the combined control set. Both visual inspection, and the inflation statistic for each ( $\lambda = 1.037$  and  $\lambda = 1.079$  respectively), suggest that the QC filtered data provides a good basis for association analyses.

## Statistical Methods

We report p-values from 1-d.f. Cochran-Armitage tests for trend as implemented in the software SNPTEST and PLINK.<sup>35</sup> We also performed 2-d.f. genotypic tests to verify that none of our associations show significant deviation from a multiplicative model, and two marker logistic regressions to test for epistasis between associated markers. Effect size estimates are based on replication samples only, and represent per-allele increase of risk in a multiplicative model.

## URLs

Affymetrix, [http://www.affymetrix.com/Auth/support/downloads/manuals/genotyping\\_console\\_manual.pdf](http://www.affymetrix.com/Auth/support/downloads/manuals/genotyping_console_manual.pdf); CelQuantileNorm, <http://outmodedbonsai.sourceforge.net>; Chiamo, [www.stats.ox.ac.uk/~marchini/software/gwas/chiamo.html](http://www.stats.ox.ac.uk/~marchini/software/gwas/chiamo.html); SNPtest, [www.stats.ox.ac.uk/~marchini/software/gwas/snpctest.html](http://www.stats.ox.ac.uk/~marchini/software/gwas/snpctest.html) [www.peopleofthebritishisles.org](http://www.peopleofthebritishisles.org)

## Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

## Acknowledgments

The principal funding for this study was provided by the Wellcome Trust, as part of the Wellcome Trust Case Control Consortium 2 project. We thank all subjects who contributed samples, and consultants and nursing staff across the UK who helped with recruitment of study subjects. We also thank Sami Bertrand, Jackie Bryant, Sarah L. Clark, Jen S. Conquer, Thomas Dibling, Stephen Gamble, Clifford Hind, Alicja Wilk, Claire R. Stribling, Sam Taylor, Julia C. Wyatt of the Wellcome Trust Sanger Institute's DNA Logistics and Genotyping Facility for technical assistance. Case collections were supported by the National Association for Colitis and Crohn's disease (NACC), the Wellcome Trust, the Medical Research Council UK, the Guy's and St Thomas' Charity, the Clinical Research Facility at the Peninsular College of Medicine and Dentistry, Exeter, the Torbay Hospital Medical Fund and the Evelyn Trust. We also acknowledge support from the Department of Health via the National Institute for Health Research (NIHR) comprehensive Biomedical Research Centre awards to Guy's & St Thomas' NHS Foundation Trust in partnership with King's College London, the Cambridge University Hospitals NHS Foundation Trust in partnership with the University of Cambridge School of Clinical Medicine and the Central Manchester Foundation Trust in partnership with the University of Manchester. We acknowledge use of the British 1958 Birth Cohort DNA collection, funded by the Medical Research Council grant G0000934 and the Wellcome Trust grant 068545/Z/02, and thank Professor Walter Bodmer and Dr Bruce Winney for use of the People of the British Isles DNA collection which was funded by the Wellcome Trust.

## Complete List of Authors

†The complete list of authors who contributed to this study is as follows:

### : The UK IBD Genetics Consortium

Jeffrey C Barrett<sup>1</sup>, James Lee<sup>2</sup>, Charlie Lees<sup>3</sup>, Natalie Prescott<sup>4</sup>, Carl A Anderson<sup>1,5</sup>, Anne Phillips<sup>3</sup>, Emma Wesley<sup>6</sup>, Kirstie Parnell<sup>6</sup>, Hu Zhang<sup>2</sup>, Hazel Drummond<sup>3</sup>, Elaine R Nimmo<sup>3</sup>, Dunecan Massey<sup>2</sup>, Kasia Blaszczyk<sup>4</sup>, Timothy Elliott<sup>7</sup>, Lynn Cotterill<sup>8</sup>, Helen Dallal<sup>9</sup>, Alan Lobo<sup>10</sup>, Craig Mowat<sup>11</sup>, Jeremy Sanderson<sup>7</sup>, Derek P Jewell<sup>12</sup>, William

<sup>1</sup>Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge CB10 1SA, UK

<sup>2</sup>Gastroenterology Research Unit, Addenbrooke's Hospital, Hills Road, Cambridge CB2 2QQ, UK

<sup>3</sup>Gastrointestinal Unit, Molecular Medicine Centre, University of Edinburgh, Western General Hospital, Crewe Road, Edinburgh EH4 2XU

<sup>4</sup>Department of Medical and Molecular Genetics, King's College London School of Medicine, Floor 8 Tower Wing, Guy's Hospital, London SE1 9RT, UK

<sup>5</sup>Wellcome Trust Centre for Human Genetics, Roosevelt Drive, Oxford OX3 7BN, UK

<sup>6</sup>Peninsula College of Medicine and Dentistry, Barrack Road, Exeter EX2 5DW, UK

<sup>7</sup>Dept Gastroenterology, Guy's & St Thomas' NHS Foundation Trust, St Thomas' Hospital, London SE1 9RT, UK

<sup>8</sup>Department of Medical Genetics, Manchester Academic Health Science Centre (MAHSC), University of Manchester and NIHR Biomedical Research Centre, Central Manchester NHS Foundation Trust, Manchester M13 0JH, UK

<sup>9</sup>Department of Gastroenterology, James Cook University Hospital, South Tees Hospitals NHS Trust, Marton Road, Middlesbrough TS4 3BW, UK

<sup>10</sup>Division of Molecular and Genetic Medicine, University of Sheffield Medical School, Royal Hallamshire Hospital, Sheffield S10 2JF, UK

<sup>11</sup>Department of General Internal Medicine, Ninewells Hospital and Medical School, Ninewells Avenue, Dundee DD1 9SY, UK

<sup>12</sup>Gastroenterology Unit, Gibson Laboratories, Radcliffe Infirmary, Woodstock Road, Oxford OX2 6HE, UK

Newman<sup>8</sup>, Cathryn Edwards<sup>13</sup>, Tariq Ahmad<sup>6</sup>, John C Mansfield<sup>14</sup>, Jack Satsangi<sup>3</sup>, Miles Parkes<sup>2</sup>, Christopher G Mathew<sup>4</sup>

## The Wellcome Trust Case Control Consortium 2

### Management Committee

Peter Donnelly (Chair)<sup>1,2</sup>, Leena Peltonen (Deputy Chair)<sup>3</sup>, Elvira Bramon<sup>4</sup>, Matthew Brown<sup>5</sup>, Juan Casas<sup>6</sup>, Aiden Corvin<sup>7</sup>, Nicholas Craddock<sup>8</sup>, Panos Deloukas<sup>3</sup>, Janus Jankowski<sup>9</sup>, Hugh Markus<sup>10</sup>, Christopher G Mathew<sup>11</sup>, Mark McCarthy<sup>12</sup>, Colin Palmer<sup>13</sup>, Robert Plomin<sup>14</sup>, Stephen Sawcer<sup>15</sup>, Richard C Trembath<sup>11</sup>, Ananth Viswanathan<sup>16</sup>, Nick Wood<sup>17</sup>

### Data and Analysis Group

Chris C A Spencer<sup>1</sup>, Jeffrey C Barrett<sup>3</sup>, Celine Bellenguez<sup>1</sup>, Daniel Davison<sup>2</sup>, Colin Freeman<sup>1</sup>, Amy Strange<sup>1</sup>, Peter Donnelly<sup>1,2</sup>

### DNA, Genotyping, Data QC and Informatics Group

Cordelia Langford<sup>3</sup>, Sarah E Hunt<sup>3</sup>, Sarah Edkins<sup>3</sup>, Rhian Gwilliam<sup>3</sup>, Hannah Blackburn<sup>3</sup>, Suzannah J. Bumpstead<sup>3</sup>, Serge Dronov<sup>3</sup>, Matthew Gillman<sup>3</sup>, Emma Gray<sup>3</sup>, Naomi Hammond<sup>3</sup>, Alagurevathi Jayakumar<sup>3</sup>, Owen T McCann<sup>3</sup>, Jennifer Liddle<sup>3</sup>, Marc L Perez<sup>3</sup>, Simon Potter<sup>3</sup>, Radhi Ravindrarah<sup>3</sup>, Michelle Ricketts<sup>3</sup>, 9 Matthew Waller<sup>3</sup>, Paul Weston<sup>3</sup>, Sara Widaa<sup>3</sup>, Pamela Whittaker<sup>3</sup>, Panos Deloukas<sup>3</sup>, Leena Peltonen<sup>3</sup>

### Publications Committee

Christopher Mathew (Chair)<sup>11</sup>, Jenefer Blackwell<sup>18</sup>, Matthew Brown<sup>5</sup>, Aiden Corvin<sup>7</sup>, Mark I McCarthy<sup>12</sup>, Chris C A Spencer<sup>1</sup>

### UK Blood Services Controls

<sup>13</sup>Endoscopy Regional Training Unit, Torbay Hospital, Torbay TQ2 7AA, UK

<sup>14</sup>Institute of Human Genetics, Newcastle University, Newcastle upon Tyne NE1 3BZ, UK

<sup>1</sup>Wellcome Trust Centre for Human Genetics, Roosevelt Drive, Oxford OX3 7LJ, UK

<sup>2</sup>Dept Statistics, University of Oxford, Oxford OX1 3TG, UK

<sup>3</sup>Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge CB10 1SA, UK

<sup>4</sup>Dept Psychological Medicine, King's College London Institute of Psychiatry Denmark Hill, London SE5 8AF, UK

<sup>5</sup>Nuffield Department of Orthopaedics, Rheumatology and Musculoskeletal Sciences, University of Oxford, Nuffield Orthopaedic Centre, Oxford OX3 7LD, UK

<sup>6</sup>Dept Epidemiology and Population Health, London School of Hygiene and Tropical Medicine, London WC1E 7HT, UK

<sup>7</sup>Neuropsychiatric Genetics Research Group, Institute of Molecular Medicine, Trinity College Dublin, Dublin 2, Eire

<sup>8</sup>Dept Psychological Medicine, Cardiff University School of Medicine, Heath Park, Cardiff CF14 4XN, UK

<sup>9</sup>Centre for Gastroenterology, Bart's and the London School of Medicine and Dentistry, London E1 2AT, UK

<sup>10</sup>Division of Cardiac and Vascular Sciences, Dept Clinical Neurosciences, St George's Hospital, London SW17 0RE, UK

<sup>11</sup>Dept Medical and Molecular Genetics, King's College London School of Medicine, Guy's Hospital, London SE1 9RT, UK

<sup>12</sup>Oxford Centre for Diabetes, Endocrinology and Metabolism (ICDEM), Churchill Hospital, Oxford OX3 7LJ, UK

<sup>13</sup>Biomedical Research Centre, Ninewells Hospital and Medical School, Dundee DD1 9SY, UK

<sup>14</sup>Social, Genetic and Developmental Psychiatry Centre, King's College London Institute of Psychiatry, Denmark Hill, London SE5 8AF, UK

<sup>15</sup>Dept Clinical Neurosciences, University of Cambridge, Addenbrooke's Hospital, Cambridge CB2 2QQ, UK

<sup>16</sup>Glaucoma Research Unit, Moorfields Eye Hospital NHS Foundation Trust, London EC1V 2PD, UK

<sup>17</sup>Dept Molecular Neuroscience, Institute of Neurology, Queen Square, London WC1N 3BG, UK

<sup>18</sup>Genetics and Infection Laboratory, Cambridge Institute of Medical Research, Addenbrooke's Hospital, Cambridge CB2 0XY, UK

Antony P Attwood<sup>3,19</sup>, Jonathan Stephens<sup>19</sup>, Jennifer Sambrook<sup>19</sup>, Willem H Ouwehand<sup>3,19</sup>

1958 Birth Cohort Controls

Wendy L McArdle<sup>20</sup>, Susan M Ring<sup>21</sup>, David P Strachan<sup>22</sup>

## Reference List

1. Rubin GP, Hungin AP, Kelly PJ, Ling J. Inflammatory bowel disease: epidemiology and management in an English general practice population. *Aliment. Pharmacol. Ther.* 2000; 14:1553–1559. [PubMed: 11121902]
2. Eaden JA, Abrams KR, Mayberry JF. The risk of colorectal cancer in ulcerative colitis: a meta-analysis. *Gut.* 2001; 48:526–535. [PubMed: 11247898]
3. Xavier RJ, Podolsky DK. Unravelling the pathogenesis of inflammatory bowel disease. *Nature.* 2007; 448:427–434. [PubMed: 17653185]
4. Wellcome Trust Case Control Consortium. Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls. *Nature.* 2007; 447:661–678. [PubMed: 17554300]
5. Anderson CA, et al. Investigation of Crohn's disease risk loci in ulcerative colitis further defines their molecular relationship. *Gastroenterology.* 2009; 136:523–529. [PubMed: 19068216]
6. Barrett JC, et al. Genome-wide association defines more than 30 distinct susceptibility loci for Crohn's disease. *Nat. Genet.* 2008; 40:955–962. [PubMed: 18587394]
7. Duerr RH, et al. A genome-wide association study identifies IL23R as an inflammatory bowel disease gene. *Science.* 2006; 314:1461–1463. [PubMed: 17068223]
8. Fisher SA, et al. Genetic determinants of ulcerative colitis include the ECM1 locus and five loci implicated in Crohn's disease. *Nat. Genet.* 2008; 40:710–712. [PubMed: 18438406]
9. Franke A, et al. Replication of signals from recent studies of Crohn's disease identifies previously unknown disease loci for ulcerative colitis. *Nat. Genet.* 2008; 40:713–715. [PubMed: 18438405]
10. Hampe J, et al. A genome-wide association scan of nonsynonymous SNPs identifies a susceptibility variant for Crohn's disease in ATG16L1. *Nat. Genet.* 2007; 39:207–211. [PubMed: 17200669]
11. Libioulle C, et al. Novel Crohn's disease locus identified by genome-wide association maps to a gene desert on 5p13.1 and modulates expression of PTGER4. *PLoS. Genet.* 2007; 3:e58. [PubMed: 17447842]
12. Parkes M, et al. Sequence variants in the autophagy gene IRGM and multiple other replicating loci contribute to Crohn's disease susceptibility. *Nat. Genet.* 2007; 39:830–832. [PubMed: 17554261]
13. Satsangi J, et al. Contribution of genes of the major histocompatibility complex to susceptibility and disease phenotype in inflammatory bowel disease. *Lancet.* 1996; 347:1212–1217. [PubMed: 8622450]
14. Franke A, et al. Sequence variants in IL10, ARPC2 and multiple other loci contribute to ulcerative colitis susceptibility. *Nat. Genet.* 2008; 40:1319–1323. [PubMed: 18836448]
15. Silverberg MS, et al. Ulcerative colitis-risk loci on chromosomes 1p36 and 12q15 found by genome-wide association study. *Nat. Genet.* 2009; 41:216–220. [PubMed: 19122664]
16. Yamagata K, et al. Mutations in the hepatocyte nuclear factor-4alpha gene in maturity-onset diabetes of the young (MODY1). *Nature.* 1996; 384:458–460. [PubMed: 8945471]
17. Barroso I, et al. Population-specific risk of type 2 diabetes conferred by HNF4A P2 promoter variants: a lesson for replication studies. *Diabetes.* 2008; 57:3161–3165. [PubMed: 18728231]

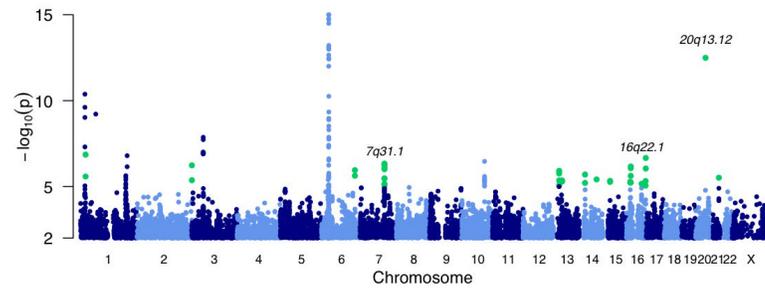
<sup>19</sup>Dept Haematology, University of Cambridge and National Health Service Blood and Transplant, Long Road, Cambridge CB2 2PT, UK

<sup>20</sup>ALSPAC DNA Bank, Dept Social Medicine, University of Bristol, 24 Tyndall Avenue, Bristol BS8 1TQ, UK

<sup>21</sup>ALSPAC Laboratory, Dept Social Medicine, Clifton, Bristol BS8 2BN, UK

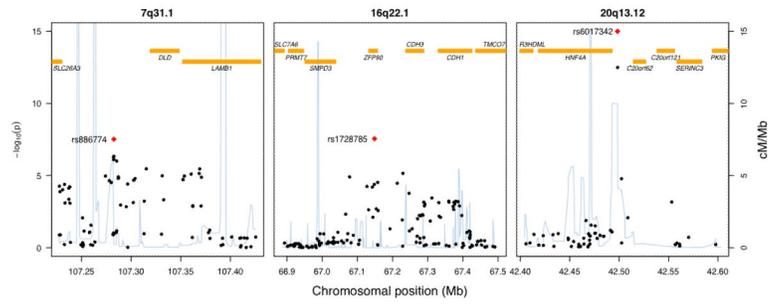
<sup>22</sup>Division of Community Health Sciences, St George's Hospital, London SW17 0RE, UK.

18. Kathiresan S, et al. Common variants at 30 loci contribute to polygenic dyslipidemia. *Nat. Genet.* 2009; 41:56–65. [PubMed: 19060906]
19. Battle MA, et al. Hepatocyte nuclear factor 4alpha orchestrates expression of cell adhesion proteins during the epithelial transformation of the developing liver. *Proc. Natl. Acad. Sci. U. S. A.* 2006; 103:8419–8424. [PubMed: 16714383]
20. Garrison WD, et al. Hepatocyte nuclear factor 4alpha is essential for embryonic development of the mouse colon. *Gastroenterology.* 2006; 130:1207–1220. [PubMed: 16618389]
21. Ahn SH, et al. Hepatocyte nuclear factor 4alpha in the intestinal epithelial cells protects against inflammatory bowel disease. *Inflamm. Bowel Dis.* 2008; 14:908–920. [PubMed: 18338782]
22. Karayiannakis AJ, et al. Expression of catenins and E-cadherin during epithelial restitution in inflammatory bowel disease. *J. Pathol.* 1998; 185:413–418. [PubMed: 9828841]
23. Houlston RS, et al. Meta-analysis of genome-wide association data identifies four new susceptibility loci for colorectal cancer. *Nat. Genet.* 2008; 40:1426–1435. [PubMed: 19011631]
24. Muise AM, et al. Polymorphisms in E-cadherin (CDH1) result in a mis-localised cytoplasmic protein that is associated with Crohn's disease. *Gut.* 2009; 58:1121–1127. [PubMed: 19398441]
25. Peignon G, et al. E-cadherin-dependent transcriptional control of apolipoprotein A-IV gene expression in intestinal epithelial cells: a role for the hepatic nuclear factor 4. *J. Biol. Chem.* 2006; 281:3560–3568. [PubMed: 16338932]
26. Vowinkel T, et al. Apolipoprotein A-IV inhibits experimental colitis. *J. Clin. Invest.* 2004; 114:260–269. [PubMed: 15254593]
27. Schmehl K, Florian S, Jacobasch G, Salomon A, Korber J. Deficiency of epithelial basement membrane laminin in ulcerative colitis affected human colonic mucosa. *Int. J. Colorectal Dis.* 2000; 15:39–48. [PubMed: 10766090]
28. Styrkarsdottir U, et al. Multiple genetic loci for bone mineral density and fractures. *N. Engl. J. Med.* 2008; 358:2355–2365. [PubMed: 18445777]
29. Liu Y, et al. A genome-wide association study of psoriasis and psoriatic arthritis identifies new disease loci. *PLoS Genet.* 2008; 4:e1000041. [PubMed: 18369459]
30. Kugathasan S, et al. Loci on 20q13 and 21q22 are associated with pediatric-onset inflammatory bowel disease. *Nat. Genet.* 2008; 40:1211–1215. [PubMed: 18758464]
31. Zhernakova A, et al. Genetic analysis of innate immunity in Crohn's disease and ulcerative colitis identifies two susceptibility loci harboring CARD9 and IL18RAP. *Am. J. Hum. Genet.* 2008; 82:1202–1210. [PubMed: 18439550]
32. Festen EA, et al. Genetic variants in the region harbouring IL2/IL21 associated with ulcerative colitis. *Gut.* 2009; 58:799–804. [PubMed: 19201773]
33. Korn JM, et al. Integrated genotype calling and association analysis of SNPs, common copy number polymorphisms and rare CNVs. *Nat. Genet.* 2008; 40:1253–1260. [PubMed: 18776909]
34. Spencer CCA. A simple clustering approach to pre-analysis exclusion of individuals from GWAS. In preparation. 2009
35. Purcell S, et al. PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am. J. Hum. Genet.* 2007; 81:559–575. [PubMed: 17701901]



**Figure 1.**

$-\log_{10}(P)$  values from the 1 d.f. trend test. Alternating chromosomes shown in shades of blue. SNPs with  $P < 1 \times 10^{-5}$  which had not been previously reported are highlighted in green. The three new loci identified in this study are noted.



**Figure 2.**

$-\log_{10}(P)$  values from the 1 d.f. trend test from three new loci, along with local recombination rate estimated from HapMap data. Combined  $P$  values for replicated SNPs are indicated with a red diamond.

**Table 1**  
**Clinical details of cases and controls**

	GWAS	Replication cohort
<b>CASES</b>	2361	2321
<b>Age at diagnosis <sup>a</sup></b>		
Early onset (<18years)	5.9% (112)	6.5% (130)
Not early onset (>18years)	94.1% (1783)	93.5% (1861)
Median	33.4	35.2
Mean	36.3	38.2
<b>Disease Extent <sup>b</sup></b>		
Proctitis	17.9% (357)	15.1% (285)
Left sided	38.8% (774)	47.5% (897)
Extensive	43.3% (864)	37.4% (707)
<b>Smoking at diagnosis <sup>c</sup></b>		
Ex-smoker	36.7% (556)	30.3% (553)
Current smoker	10.8% (163)	16.4% (300)
Never smoked	52.5% (794)	53.2% (971)
<b>Colectomy <sup>d</sup></b>		
Yes	15.7% (266)	12.0% (226)
No	84.3% (1432)	88.0% (1660)
<b>Colorectal cancer <sup>e</sup></b>		
	1.00% (23)	0.87% (20)
<b>CONTROLS</b>		
Total	5417	4818
UKBS	2675	-
1958 Birth Cohort	2742	1952
POBI	-	2866

<sup>a</sup>Data available for 80% (GWAS) 86% (Replication),

<sup>b</sup>Data available for 84% (GWAS) 81% (Replication),

<sup>c</sup>Data available for 64% (GWAS) 79% (Replication),

<sup>d</sup>Data available for 72% (GWAS) 81% (Replication),

<sup>e</sup>Data available for 93% (GWAS) 99% (Replication).

and

Page 14

Table 2

## New hits from the GWAS

Top tier reaches  $5 \times 10^{-8}$  in combined analysis, second tier hits have replication  $P < 0.05$  and require further study to completely verify.

SNP	Chr	LD region(Mb) <sup>a</sup>	Gene of interest (#) <sup>b</sup>	$P_{scan}$	$P_{repl}$	$P_{comb}$	Risk allele	RAF <sup>c</sup>	OR (95% CI)
rs886774	7q31.1	107.25-107.39	<i>LAMB1</i> (2)	$4.8 \times 10^{-7}$	0.005	$3 \times 10^{-8}$	G	0.4136	1.11 (1.03-1.19)
rs1728785	16q22.1	66.98-67.40	<i>CDH1</i> (5)	$1.8 \times 10^{-5}$	0.0004	$2.8 \times 10^{-8}$	G	0.7641	1.17 (1.07-1.27)
rs6017342	20q13.12	42.49-42.52	<i>HNF4A</i> (7)	$3.2 \times 10^{-13}$	$7.1 \times 10^{-6}$	$8.5 \times 10^{-17}$	C	0.5168	1.17 (1.09-1.26)
rs7524102*	1p36.12	22.54-22.61	(0)	$1.4 \times 10^{-7}$	0.05	$3.1 \times 10^{-7}$	A	0.8264	1.10 (1.00-1.21)
rs9548988	13q13.3	39.36-39.56	(0)	$5.0 \times 10^{-6}$	0.0061	$2.7 \times 10^{-7}$	T	0.4594	1.10 (1.03-1.19)

\* replication genotyping at this locus is for SNP rs12568930, which is an  $r^2 = 1$  proxy for rs7524102.

<sup>a</sup>LD region of 0.2 cM centered on focal SNP, in NCBI Build 36 coordinates.

<sup>b</sup>Number of genes in LD region.

<sup>c</sup>Risk allele frequency.

and

Page 15

Table 3

**GWAS signals from previously reported UC loci**

Top tier were previously reported at genome-wide significance ( $5 \times 10^{-8}$ ), bottom tier were previously reported with weaker evidence. *P* values are one-tailed in the direction of the previously reported association.

SNP	Chrom	Pos	Gene	Scan <i>P</i>	Ref
rs6426833	1p36.13	20044447	<i>OTUD3/PLA2GGE</i>	$2.1 \times 10^{-11}$	15
rs11209026	1p31.3	67478546	<i>IL23R</i>	$3.0 \times 10^{-10}$	15
rs3024493	1q32.1	205010591	<i>IL10</i>	$8.0 \times 10^{-8}$	14
rs10021288	4q27	123224984	<i>IL2/21</i>	0.0033	32
rs9268877	6p21.32	32539125	<i>MHC</i>	$3.9 \times 10^{-23}$	8
rs12815372	12q15	66765480	<i>IL26</i>	0.00070	15
rs311497	20q13.33	61691693	<i>TNFRSF6B</i>	0.0018	30
rs2094871	21q22.2	39382729	<i>PSMG1</i>	$1.6 \times 10^{-6}$	30
rs7511649	1q21.2	148537415	<i>ECMI</i>	0.00015	8
rs7554511	1q32.1	199144185	<i>KIF21B</i>	$1.2 \times 10^{-6}$	5
rs12612347	2q35	218765583	<i>ARPC2</i>	0.024	14
rs9858542	3p21.31	49676987	<i>MST1</i>	$7.0 \times 10^{-9}$	8
rs1368438	5q33.3	158639883	<i>IL12B</i>	0.0039	8
rs12529198	6p25.1	5096246	<i>LYRM4</i>	0.13	5
rs6908425	6p22.3	20836710	<i>CDKALI</i>	0.0044	5
rs10974914	9p24.1	5004332	<i>JAK2</i>	$1.5 \times 10^{-5}$	5
rs10781500	9q34.3	138389159	<i>CARD9</i>	$7.0 \times 10^{-6}$	31
rs17582416	10p11.21	35327656	<i>CCNY</i>	0.022	5
rs10995271	10q21.2	64108492	none	0.32	8, 9
rs6584283	10q24.2	101280291	<i>NKX2-3</i>	$1.7 \times 10^{-7}$	8
rs916977	15q13.1	26186959	<i>HERC2</i>	0.26	9
rs744166	17q21.2	37767727	<i>STAT3</i>	0.0025	5, 9
rs2542151	18p11.21	12769947	<i>PTPN2</i>	0.0010	9