

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

The UK IBD Genetics Consortium & the Wellcome Trust Case Control Consortium 2\*

**Ulcerative colitis is a common form of inflammatory bowel disease with a complex etiology. As part of the Wellcome Trust Case Control Consortium 2, we performed a genome-wide association scan for ulcerative colitis in 2,361 cases and 5,417 controls. Loci showing evidence of association at  $P < 1 \times 10^{-5}$  were followed up by genotyping in an independent set of 2,321 cases and 4,818 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 (*CDH1* and *CDH3*;  $P = 2.8 \times 10^{-8}$ ) and 7q31 (*LAMB1*;  $P = 3.0 \times 10^{-8}$ ). Of note, *CDH1* has recently been associated with susceptibility to colorectal cancer, an established complication of longstanding ulcerative colitis. The new associations suggest that changes in the integrity of the intestinal epithelial barrier may contribute to the pathogenesis of ulcerative colitis.**

Genetic epidemiological data clearly implicate inherited susceptibility in the pathogenesis of ulcerative colitis and Crohn's disease, which represent the two most common forms of inflammatory bowel disease (IBD) and together affect at least 1 in 250 individuals in 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 affected individuals, 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 Crohn's disease and ulcerative colitis are related polygenic diseases. This has recently been corroborated by the results of genetic association studies, which have highlighted both disease-specific loci and other loci that are shared between ulcerative colitis and Crohn's disease. For example, whereas genetically determined defects in the handling of intracellular bacteria (related to the gene *NOD2*, and the autophagy genes *ATG16L1*

and *IRGM*) are specific to Crohn's disease, genes encoding multiple components in the Th17 pathway (*IL23R*, *IL12B*, *JAK2* and *STAT3*) are associated with both Crohn's disease and ulcerative colitis<sup>4–12</sup>.

Until recently, most attention had focused on Crohn's disease, with genome-wide association (GWA) studies and subsequent meta-analysis yielding more than 30 confirmed Crohn's disease-susceptibility loci<sup>4,6,7,10–12</sup>. In addition to the longstanding known association in the major histocompatibility complex (MHC)<sup>13</sup>, the first GWA scans in ulcerative colitis reported associations at *IL23R*, *IL10* and loci on chromosomes 1p36 and 12q15 that 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 ulcerative colitis to date. All study subjects were UK residents of European ancestry; clinical data are presented in **Table 1**. Affected individuals (cases) and controls were genotyped on the Affymetrix 6.0 array. After application of quality control filters (see Online Methods), we analyzed GWA data from 2,361 individuals with ulcerative colitis and 5,417 controls (**Fig. 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 loci 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 2,321 ulcerative colitis cases and 4,818 controls. Three new loci showed evidence for association at  $P < 5 \times 10^{-8}$  in the combined panel, with three additional new loci showing nominal ( $P < 0.05$ ) replication (**Table 2** and **Fig. 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 hot spot on chromosome 20q13 containing the 3' untranslated region of just one gene, *HNF4A*. The SNP rs6017342 itself maps 5 kb distal to the 3' untranslated region and is located within an expressed sequence tag DB076868, which has been detected in only a single testis cDNA library and does not

\*A full list of authors and affiliations appears at the end of the paper.

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

	GWAS	Replication cohort
<b>Cases</b>	2,361	2,321
<b>Age at diagnosis<sup>a</sup></b>		
Early onset (<18 years)	5.9% (112)	6.5% (130)
Not early onset (>18 years)	94.1% (1,783)	93.5% (1,861)
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% (1,432)	88.0% (1,660)
<b>Colorectal cancer<sup>e</sup></b>		
	1.00% (23)	0.87% (20)
<b>Controls</b>		
Total	5,417	4,818
UKBS	2,675	–
1958 Birth Cohort	2,742	1,952
PoBI	–	2,866

<sup>a</sup>Data available for 80% (GWAS) and 86% (Replication) of cases. <sup>b</sup>Data available for 84% (GWAS) and 81% (Replication) of cases. <sup>c</sup>Data available for 64% (GWAS) and 79% (Replication) of cases. <sup>d</sup>Data available for 72% (GWAS) and 81% (Replication) of cases. <sup>e</sup>Data available for 93% (GWAS) and 99% (Replication) of cases.

encode any extended 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. Because rs6017342 is located within a recombination hot spot, there are few known SNPs in strong linkage disequilibrium (LD) ( $r^2 > 0.5$ ) with it; there are no additional SNPs in LD with rs6017342 on the Affymetrix chip used in this study or on the Illumina chips used in previous studies. As the evidence for association with ulcerative colitis rests on this single SNP, we subjected these data to careful scrutiny; genotype cluster plots for this SNP showed clear resolution of the three genotype classes (**Supplementary Fig. 1**), with 99.3% completeness of genotypes within this dataset.

Rare *HNF4A* mutations account for approximately 4% of individuals in the UK with 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  $\beta$ -cell dysfunction and sensitivity to sulfonylureas. Common variants of *HNF4A* influence predisposition to type 2 diabetes (rs2144908)<sup>17</sup> and dyslipidemia (rs1800961)<sup>18</sup>. The ulcerative colitis-associated SNP, rs6017342, is not in LD with either of these two common variants, nor did it show association in our study of Crohn's disease ( $P = 0.92$ )<sup>4</sup>.

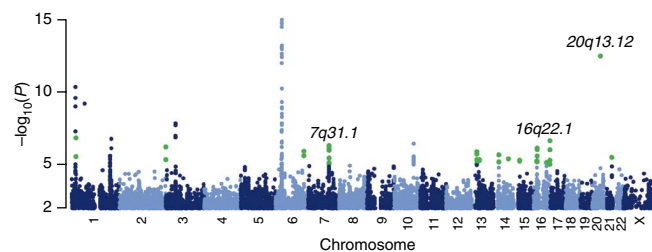
*HNF4A* encodes the transcription factor hepatocyte nuclear factor 4 $\alpha$ , which regulates the expression of multiple components within all three key compartments of the cell-cell junction, namely, the adherens junction, the tight junction and the desmosome<sup>19</sup>. Such

cell-cell junctions are fundamental to epithelial organization and barrier function. *HNF4 $\alpha$*  also has 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 fetal colon die perinatally. Histological analysis of colonic tissue recovered during late development (embryonic day (E) 18.5) 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 mouse intestinal inflammation, Ahn and colleagues circumnavigated the embryonic lethality of *Hnf4a*<sup>-/-</sup> mice by generating a conditional model of intestinal *Hnf4a* deletion<sup>21</sup>. These *Hnf4a*<sup>ΔEPC</sup> mice (*loxP*-flanked *Hnf4a* driven by the *Vill1* (villin 1) promoter) developed increased epithelial permeability and markedly more severe colitis following dextran sodium sulfate (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 ulcerative colitis<sup>21</sup>, a finding that 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 that is bounded by recombination hot spots spans 411 kb and encodes several genes. Among the strongest candidates for ulcerative colitis susceptibility is *CDH1*, which encodes E-cadherin. This transmembrane glycoprotein is one of the main components of the adherens junction and is 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 *CDH1* is known to be reduced in areas of active ulcerative colitis<sup>22</sup>.

Given the well-recognized association between ulcerative colitis and colorectal cancer<sup>2</sup>, the observation of correlated association signals at the *CDH1* locus in both diseases is noteworthy. In particular, variants in LD ( $r^2 = 0.5$ ) with the SNP most strongly associated with ulcerative colitis 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 strongly associated SNP in the colorectal cancer study is also associated with ulcerative colitis ( $P = 8 \times 10^{-4}$ ). This locus did not show association with Crohn's disease in a large international GWA meta-analysis ( $P = 0.549$ )<sup>6</sup> (**Supplementary Table 2**). However, evidence for association of *CDH1* with Crohn's disease was reported recently in the Canadian population using a candidate gene approach<sup>24</sup>, and the presence of Crohn's disease-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>.

The evidence that *HNF4 $\alpha$*  and E-cadherin cooperate to maintain epithelial barrier integrity in the intestine is of great potential relevance. In experiments focused on the liver, *Hnf4a* knockout mice



**Figure 1** Plot of genome-wide association results.  $-\log_{10}(P)$  values are from the 1-d.f. trend test. Alternating chromosomes are 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.

**Table 2** New hits from the GWAS

SNP	Chr.	LD region (Mb) <sup>a</sup>	Gene of interest (#) <sup>b</sup>	$P_{\text{scan}}$	$P_{\text{repl}}$	$P_{\text{comb}}$	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> (1)	$3.2 \times 10^{-13}$	$7.1 \times 10^{-6}$	$8.5 \times 10^{-17}$	C	0.5168	1.17 (1.09–1.26)
rs7524102 <sup>d</sup>	1p36.12	22.54–22.61	None	$1.4 \times 10^{-7}$	0.05	$3.1 \times 10^{-7}$	A	0.8264	1.10 (1.00–1.21)
rs9548988	13q13.3	39.36–39.56	None	$5.0 \times 10^{-6}$	0.0061	$2.7 \times 10^{-7}$	T	0.4594	1.10 (1.03–1.19)

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

<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. <sup>d</sup>Replication genotyping at this locus is for SNP rs12568930, which is a proxy ( $r^2=1$ ) for rs7524102.

failed to express E-cadherin<sup>19</sup>, whereas in a human intestinal cell line, E-cadherin-dependent cell-cell contact was found to be critical in determining the amount and binding activity of nuclear HNF4 $\alpha$ <sup>25</sup>. This in turn affected the expression of several genes, including *APOA4* (ref. 25), which encodes an anti-inflammatory protein known to inhibit experimental colitis<sup>26</sup>.

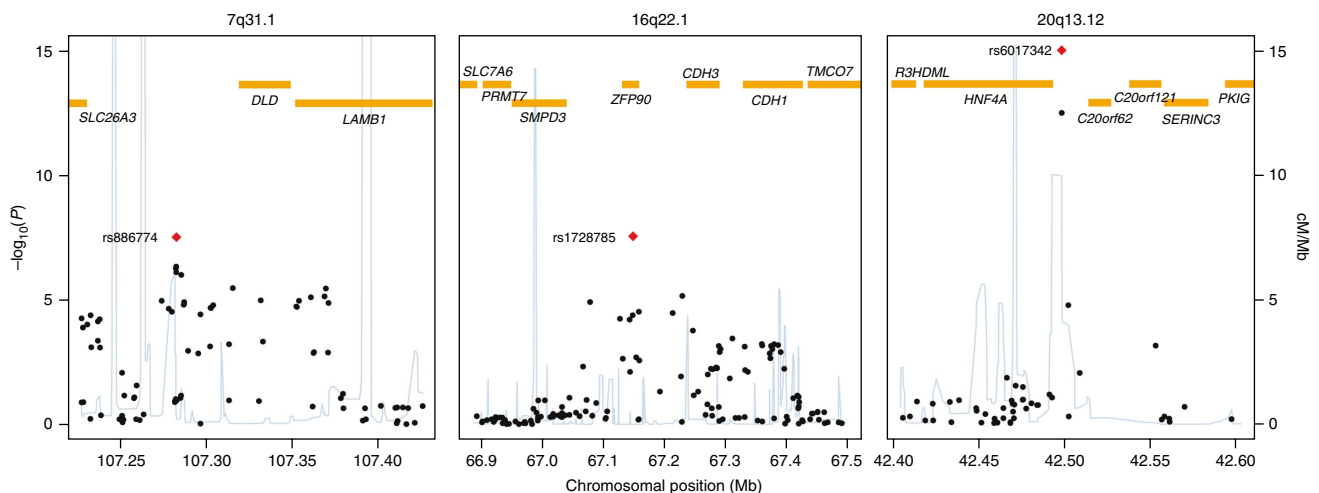
The third newly confirmed ulcerative colitis susceptibility locus was a region on chromosome 7q31, previously suggested by a recent North American GWA scan<sup>15</sup>. In the current study, the peak association was seen at rs886774 (GWA scan  $P = 4.8 \times 10^{-7}$ ; combined GWA and replication  $P = 3.0 \times 10^{-8}$ ). 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 laminin-1, laminin-2 and laminin-10. Laminins are expressed in the intestinal basement membrane and play a key role in anchoring the single-layered epithelium; expression of laminins is known to be downregulated in ulcerative colitis<sup>27</sup>. rs886774 was not associated with Crohn's disease in the meta-analysis<sup>6</sup> (Supplementary Table 2).

Two other loci previously implicated in ulcerative colitis-related phenotypes showed strong (but not genome-wide significant) association here. 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 near (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 newly discovered loci described above, our GWAS detected strong association at established ulcerative colitis 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 ulcerative colitis loci reported recently in genome-wide scans: the *IL10* locus<sup>14</sup> and the *OTUD3-PLA2G2E* locus<sup>15</sup>, on chromosomes 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 ulcerative colitis. Moderate support was obtained for some loci previously reported to be associated with ulcerative colitis, including *ECM1*, *CARD9* (ref. 31), 1q32-*KIF21B* and 9p24-*JAK2*, but weaker support was obtained for other loci, such as *IL2-IL21* (ref. 32), *IL12B* and 12q15 (Table 3). Some of the ulcerative colitis loci are clearly associated with Crohn's disease, whereas others are not associated 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. We have identified three new susceptibility loci for ulcerative colitis and provide the first genetic link between ulcerative colitis and colorectal cancer. The strongest new association intervals that we detected contain, 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 ulcerative colitis 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 significance 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.



**Figure 2** Regional association plots for the three newly discovered loci.  $-\log_{10}(P)$  values from the 1-d.f. trend test from three new loci, along with local recombination rates estimated from HapMap data. Combined  $P$  values for replicated SNPs are indicated with a red diamond.

**Table 3 GWAS signals from loci previously reported to be associated with ulcerative colitis**

SNP	Chr.	Pos	Gene	Scan <i>P</i>	Ref.
rs6426833	1p36.13	20,044,447	<i>OTUD3-PLA2G2E</i>	$2.1 \times 10^{-11}$	15
rs11209026	1p31.3	67,478,546	<i>IL23R</i>	$3.0 \times 10^{-10}$	15
rs3024493	1q32.1	205,010,591	<i>IL10</i>	$8.0 \times 10^{-8}$	14
rs10021288	4q27	123,224,984	<i>IL2-IL21</i>	0.0033	32
rs9268877	6p21.32	32,539,125	MHC	$3.9 \times 10^{-23}$	8
rs12815372	12q15	66,765,480	<i>IL26</i>	0.00070	15
rs311497	20q13.33	61,691,693	<i>TNFRSF6B</i>	0.0018	30
rs2094871	21q22.2	39,382,729	<i>PSMG1</i>	$1.6 \times 10^{-6}$	30
rs7511649	1q21.2	148,537,415	<i>ECM1</i>	0.00015	8
rs7554511	1q32.1	199,144,185	<i>KIF21B</i>	$1.2 \times 10^{-6}$	5
rs12612347	2q35	218,765,583	<i>ARPC2</i>	0.024	14
rs9858542	3p21.31	49,676,987	<i>MST1</i>	$7.0 \times 10^{-9}$	8
rs1368438	5q33.3	158,639,883	<i>IL12B</i>	0.0039	8
rs12529198	6p25.1	5,096,246	<i>LYRM4</i>	0.13	5
rs6908425	6p22.3	20,836,710	<i>CDKAL1</i>	0.0044	5
rs10974914	9p24.1	5,004,332	<i>JAK2</i>	$1.5 \times 10^{-5}$	5
rs10781500	9q34.3	138,389,159	<i>CARD9</i>	$7.0 \times 10^{-6}$	31
rs17582416	10p11.21	35,327,656	<i>CCNY</i>	0.022	5
rs10995271	10q21.2	64,108,492	none	0.32	8,9
rs6584283	10q24.2	101,280,291	<i>NKX2-3</i>	$1.7 \times 10^{-7}$	8
rs916977	15q13.1	26,186,959	<i>HERC2</i>	0.26	9
rs744166	17q21.2	37,767,727	<i>STAT3</i>	0.0025	5,9
rs2542151	18p11.21	12,769,947	<i>PTPN2</i>	0.0010	9

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

## METHODS

Methods and any associated references are available in the online version of the paper at <http://www.nature.com/naturegenetics/>.

Note: Supplementary information is available on the Nature Genetics website.

## ACKNOWLEDGMENTS

The principal funding for this study was provided by the Wellcome Trust, as part of the Wellcome Trust Case Control Consortium 2 project (083948/Z/07/Z). 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 S. Bertrand, J. Bryant, S.L. Clark, J.S. Conquer, T. Dibling, J.C. Eldred, S. Gamble, C. Hind, A. Wilk, C.R. Stribling and S. Taylor of the Wellcome Trust Sanger Institute's Sample and Genotyping Facilities 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. P. Donnelly was supported in part by a Wolfson–Royal Society Merit Award. We also acknowledge support from the UK Medical Research Council (R.C.T., grant G0601387), the Special Trustees of Moorfields National Health Service (NHS) Foundation Trust, and 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, the Central Manchester NHS Foundation Trust in partnership with the University of Manchester, and Moorfields Eye Hospital in partnership with University College London Institute of Ophthalmology. 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 W. Bodmer and B. Winney for use of the People of the British Isles DNA collection, which was funded by the Wellcome Trust.

## AUTHOR CONTRIBUTIONS

J.C.L., C.W.L., N.J.P., A.P., E.W., K.P., H.Z., H.D., E.R.N., D.M., K.B., T.E. and L.C. were involved in establishing DNA collections and/or assembling phenotypic data. C.W.L., A.P., E.W., D.M., H.D., A.J.L., C.M., J.D.S., D.P.J., C.E., T.A., J.C.M., J. Satsangi and M.P. recruited patients. W.G.N., C.E., T.A., J.C.M., J.D.S., M.P. and C.G.M. supervised clinical and laboratory work. The WTCCC2 DNA, genotyping, data QC and

informatics group executed GWAS sample handling, genotyping and quality control. The WTCCC2 data and analysis group, J.C.B. and C.A.A. performed statistical analyses. J.C.B., J.C.L., C.W.L., N.J.P., C.C.A.S., C.A.A., T.A., P. Donnelly, J. Satsangi, M.P. and C.G.M. contributed to writing the manuscript. The WTCCC2 management committee conceived and oversaw the design and execution of the GWAS. WTCCC2 group memberships are specified in the full author list.

Published online at <http://www.nature.com/naturegenetics/>.

Reprints and permissions information is available online at <http://npg.nature.com/reprintsandpermissions/>.

- Rubin, G.P., Hungin, A.P., Kelly, P.J. & Ling, J. Inflammatory bowel disease: epidemiology and management in an English general practice population. *Aliment. Pharmacol. Ther.* **14**, 1553–1559 (2000).
- Eaden, J.A., Abrams, K.R. & Mayberry, J.F. The risk of colorectal cancer in ulcerative colitis: a meta-analysis. *Gut* **48**, 526–535 (2001).
- Xavier, R.J. & Podolsky, D.K. Unravelling the pathogenesis of inflammatory bowel disease. *Nature* **448**, 427–434 (2007).
- Wellcome Trust Case Control Consortium. Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls. *Nature* **447**, 661–678 (2007).
- Anderson, C.A. *et al.* Investigation of Crohn's disease risk loci in ulcerative colitis further defines their molecular relationship. *Gastroenterology* **136**, 523–529 (2009).
- Barrett, J.C. *et al.* Genome-wide association defines more than 30 distinct susceptibility loci for Crohn's disease. *Nat. Genet.* **40**, 955–962 (2008).
- Duerr, R.H. *et al.* A genome-wide association study identifies *IL23R* as an inflammatory bowel disease gene. *Science* **314**, 1461–1463 (2006).
- Fisher, S.A. *et al.* Genetic determinants of ulcerative colitis include the *ECM1* locus and five loci implicated in Crohn's disease. *Nat. Genet.* **40**, 710–712 (2008).
- Franke, A. *et al.* Replication of signals from recent studies of Crohn's disease identifies previously unknown disease loci for ulcerative colitis. *Nat. Genet.* **40**, 713–715 (2008).
- Hampe, J. *et al.* A genome-wide association scan of nonsynonymous SNPs identifies a susceptibility variant for Crohn's disease in *ATG16L1*. *Nat. Genet.* **39**, 207–211 (2007).
- 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.* **3**, e58 (2007).
- Parkes, M. *et al.* Sequence variants in the autophagy gene *IRGM* and multiple other replicating loci contribute to Crohn's disease susceptibility. *Nat. Genet.* **39**, 830–832 (2007).

13. Satsangi, J. *et al.* Contribution of genes of the major histocompatibility complex to susceptibility and disease phenotype in inflammatory bowel disease. *Lancet* **347**, 1212–1217 (1996).
14. Franke, A. *et al.* Sequence variants in *IL10*, *ARPC2* and multiple other loci contribute to ulcerative colitis susceptibility. *Nat. Genet.* **40**, 1319–1323 (2008).
15. Silverberg, M.S. *et al.* Ulcerative colitis-risk loci on chromosomes 1p36 and 12q15 found by genome-wide association study. *Nat. Genet.* **41**, 216–220 (2009).
16. Yamagata, K. *et al.* Mutations in the hepatocyte nuclear factor-4 $\alpha$  gene in maturity-onset diabetes of the young (MODY1). *Nature* **384**, 458–460 (1996).
17. Barroso, I. *et al.* Population-specific risk of type 2 diabetes conferred by *HNF4A* P2 promoter variants: a lesson for replication studies. *Diabetes* **57**, 3161–3165 (2008).
18. Kathiresan, S. *et al.* Common variants at 30 loci contribute to polygenic dyslipidemia. *Nat. Genet.* **41**, 56–65 (2009).
19. Battle, M.A. *et al.* Hepatocyte nuclear factor 4 $\alpha$  orchestrates expression of cell adhesion proteins during the epithelial transformation of the developing liver. *Proc. Natl. Acad. Sci. USA* **103**, 8419–8424 (2006).
20. Garrison, W.D. *et al.* Hepatocyte nuclear factor 4 $\alpha$  is essential for embryonic development of the mouse colon. *Gastroenterology* **130**, 1207–1220 (2006).
21. Ahn, S.H. *et al.* Hepatocyte nuclear factor 4 $\alpha$  in the intestinal epithelial cells protects against inflammatory bowel disease. *Inflamm. Bowel Dis.* **14**, 908–920 (2008).
22. Karayiannakis, A.J. *et al.* Expression of catenins and E-cadherin during epithelial restitution in inflammatory bowel disease. *J. Pathol.* **185**, 413–418 (1998).
23. Houlston, R.S. *et al.* Meta-analysis of genome-wide association data identifies four new susceptibility loci for colorectal cancer. *Nat. Genet.* **40**, 1426–1435 (2008).
24. Muise, A.M. *et al.* Polymorphisms in E-cadherin (*CDH1*) result in a mis-localised cytoplasmic protein that is associated with Crohn's disease. *Gut* **58**, 1121–1127 (2009).
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.* **281**, 3560–3568 (2006).
26. Vowinkel, T. *et al.* Apolipoprotein A-IV inhibits experimental colitis. *J. Clin. Invest.* **114**, 260–269 (2004).
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.* **15**, 39–48 (2000).
28. Styrkarsdottir, U. *et al.* Multiple genetic loci for bone mineral density and fractures. *N. Engl. J. Med.* **358**, 2355–2365 (2008).
29. Liu, Y. *et al.* A genome-wide association study of psoriasis and psoriatic arthritis identifies new disease loci. *PLoS Genet.* **4**, e1000041 (2008).
30. Kugathasan, S. *et al.* Loci on 20q13 and 21q22 are associated with pediatric-onset inflammatory bowel disease. *Nat. Genet.* **40**, 1211–1215 (2008).
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.* **82**, 1202–1210 (2008).
32. Festen, E.A. *et al.* Genetic variants in the region harbouring *IL2/IL21* associated with ulcerative colitis. *Gut* **58**, 799–804 (2009).

### The UK IBD Genetics Consortium

Jeffrey C Barrett<sup>1</sup>, James C Lee<sup>2</sup>, Charles W Lees<sup>3</sup>, Natalie J 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 J Lobo<sup>10</sup>, Craig Mowat<sup>11</sup>, Jeremy D Sanderson<sup>7</sup>, Derek P Jewell<sup>12</sup>, William G 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>5,15</sup>, Leena Peltonen (Deputy Chair)<sup>1</sup>, Jenefer M Blackwell<sup>16</sup>, Elvira Bramon<sup>17</sup>, Matthew A Brown<sup>18</sup>, Juan P Casas<sup>19</sup>, Aiden Corvin<sup>20</sup>, Nicholas Craddock<sup>21</sup>, Panos Deloukas<sup>1</sup>, Audrey Duncanson<sup>22</sup>, Janusz Jankowski<sup>23</sup>, Hugh S Markus<sup>24</sup>, Christopher G Mathew<sup>4</sup>, Mark I McCarthy<sup>25</sup>, Colin N A Palmer<sup>26</sup>, Robert Plomin<sup>27</sup>, Anna Rautanen<sup>5</sup>, Stephen J Sawcer<sup>28</sup>, Nilesh Samani<sup>29</sup>, Richard C Trembath<sup>4</sup>, Ananth C Viswanathan<sup>30,31</sup> & Nicholas Wood<sup>32</sup>

**Data and Analysis Group:** Chris C A Spencer<sup>5</sup>, Jeffrey C Barrett<sup>1</sup>, Céline Bellenguez<sup>5</sup>, Daniel Davison<sup>15</sup>, Colin Freeman<sup>5</sup>, Amy Strange<sup>5</sup> & Peter Donnelly<sup>5,15</sup>

**DNA, Genotyping, Data QC and Informatics Group:** Cordelia Langford<sup>1</sup>, Sarah E Hunt<sup>1</sup>, Sarah Edkins<sup>1</sup>, Rhian Gwilliam<sup>1</sup>, Hannah Blackburn<sup>1</sup>, Suzannah J Bumpstead<sup>1</sup>, Serge Dronov<sup>1</sup>, Matthew Gillman<sup>1</sup>, Emma Gray<sup>1</sup>, Naomi Hammond<sup>1</sup>, Alagurevathi Jayakumar<sup>1</sup>, Owen T McCann<sup>1</sup>, Jennifer Liddle<sup>1</sup>, Marc L Perez<sup>1</sup>, Simon C Potter<sup>1</sup>, Radhi Ravindraraaj<sup>1</sup>, Michelle Ricketts<sup>1</sup>, Matthew Waller<sup>1</sup>, Paul Weston<sup>1</sup>, Sara Widaa<sup>1</sup>, Pamela Whittaker<sup>1</sup>, Panos Deloukas<sup>1</sup> & Leena Peltonen<sup>1</sup>

**Publications Committee:** Christopher G Mathew (Chair)<sup>4</sup>, Jenefer M Blackwell<sup>16</sup>, Matthew A Brown<sup>18</sup>, Aiden Corvin<sup>20</sup>, Mark I McCarthy<sup>25</sup> & Chris C A Spencer<sup>5</sup>

**UK Blood Services Controls:** Antony P Attwood<sup>1,33</sup>, Jonathan Stephens<sup>33</sup>, Jennifer Sambrook<sup>33</sup> & Willem H Ouwehand<sup>1,33</sup>

**1958 Birth Cohort Controls:** Wendy L McArdle<sup>34</sup>, Susan M Ring<sup>35</sup> & David P Strachan<sup>36</sup>

<sup>1</sup>Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge, UK. <sup>2</sup>Gastroenterology Research Unit, Addenbrooke's Hospital, Hills Road, Cambridge, UK. <sup>3</sup>Gastrointestinal Unit, Molecular Medicine Centre, University of Edinburgh, Western General Hospital, Crewe Road, Edinburgh, UK. <sup>4</sup>Department of Medical and Molecular Genetics, King's College London School of Medicine, Guy's Hospital, London, UK. <sup>5</sup>Wellcome Trust Centre for Human Genetics, Roosevelt Drive, Oxford, UK. <sup>6</sup>Peninsula College of Medicine and Dentistry, Barrack Road, Exeter, UK. <sup>7</sup>Department of Gastroenterology, Guy's & St. Thomas' NHS Foundation Trust, St Thomas' Hospital, London, 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, UK. <sup>9</sup>Department of Gastroenterology, James Cook University Hospital, South Tees Hospitals NHS Trust, Middlesbrough, UK. <sup>10</sup>Division of Molecular and Genetic Medicine, University of Sheffield Medical School, Royal Hallamshire Hospital, Sheffield, UK. <sup>11</sup>Department of General Internal Medicine, Ninewells Hospital and Medical School, Ninewells Avenue, Dundee, UK. <sup>12</sup>Gastroenterology Unit, Gibson Laboratories, Radcliffe Infirmary, Woodstock Road, Oxford, UK. <sup>13</sup>Endoscopy Regional Training Unit, Torbay Hospital, Torbay, UK. <sup>14</sup>Institute of Human Genetics, Newcastle University, Newcastle upon Tyne, UK. <sup>15</sup>Department of Statistics, University of Oxford, Oxford, UK. <sup>16</sup>Genetics and Infection Laboratory, Cambridge Institute of Medical Research, Addenbrooke's Hospital, Cambridge, UK. <sup>17</sup>Division of Psychological Medicine and Psychiatry, Biomedical Research Centre for Mental Health at the Institute of Psychiatry, King's College London and The South London and Maudsley NHS Foundation Trust, Denmark Hill, London, UK. <sup>18</sup>Diamantina Institute of Cancer, Immunology and Metabolic Medicine, Princess Alexandra Hospital, University of Queensland, Brisbane, Queensland, Australia. <sup>19</sup>Department of Epidemiology and Population Health, London School of Hygiene and Tropical Medicine, London, UK. <sup>20</sup>Neuropsychiatric Genetics Research Group, Institute of Molecular Medicine, Trinity College Dublin, Dublin, Ireland. <sup>21</sup>Department of Psychological Medicine, Cardiff University School of Medicine, Heath Park, Cardiff, UK. <sup>22</sup>Molecular and Physiological Sciences, The Wellcome Trust, London, UK. <sup>23</sup>Centre for Gastroenterology, Bart's and the London School of Medicine and Dentistry, London, UK. <sup>24</sup>Clinical Neurosciences, St. George's University of London, London, UK. <sup>25</sup>Oxford Centre for Diabetes, Endocrinology and Metabolism (ICDEM), Churchill Hospital, Oxford, UK. <sup>26</sup>Biomedical Research Centre, Ninewells Hospital and Medical School, Dundee, UK. <sup>27</sup>Social, Genetic and Developmental Psychiatry Centre, King's College London Institute of Psychiatry, Denmark Hill, London, UK. <sup>28</sup>University of Cambridge Department of Clinical Neurosciences, Addenbrooke's Hospital, Cambridge, UK. <sup>29</sup>Department of Cardiovascular Science, University of Leicester, Glenfield Hospital, Leicester, UK. <sup>30</sup>Glaucoma Research Unit, Moorfields Eye Hospital NHS Foundation Trust, London, UK. <sup>31</sup>Department of Genetics, University College London Institute of Ophthalmology, London, UK. <sup>32</sup>Department of Molecular Neuroscience, Institute of Neurology, Queen Square, London, UK. <sup>33</sup>Department of Haematology, University of Cambridge and NHS Blood and Transplant, Long Road, Cambridge, UK. <sup>34</sup>Avon Longitudinal Study of Parents and Children (ALSPAC) DNA Bank, Department of Social Medicine, University of Bristol, Bristol, UK. <sup>35</sup>ALSPAC Laboratory, Department of Social Medicine, Clifton, Bristol, UK. <sup>36</sup>Division of Community Health Sciences, St. George's Hospital, London, UK. Correspondence should be addressed to J.C.B. (barrett@sanger.ac.uk) or C.C.A.S. (chris.spencer@well.ox.ac.uk).

## ONLINE METHODS

**Subjects.** *Cases:* A total of 5,319 unrelated individuals of European ancestry with a diagnosis of ulcerative colitis established using standard endoscopic, radiological and histological criteria were recruited from ten centers within the UK (Cambridge, Oxford, London, Newcastle, Sheffield, Edinburgh, Dundee, Manchester, Torbay and Exeter) and defined as cases. All study participants 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 before sample collection (Cambridge, Oxford, London, Newcastle, Sheffield, Edinburgh, Dundee, Manchester, Torbay and Exeter Local Research Ethics Committees). After applying quality control measures (see below and **Supplementary Table 3**), we analyzed a total of 4,682 samples, which were divided between the discovery panel (2,361 samples) and replication panel (2,321 samples).

*Controls:* A total of 10,235 population control DNA samples from three sources passed our quality control filters (see below). 5,417 samples of the WTCCC2 common control set were used for the GWA experiment. This comprised 2,675 healthy blood donors recruited from the United Kingdom Blood Service (UKBS) and 2,742 samples from the 1958 Birth Cohort (1958BC) obtained from Epstein Barr virus-transformed cell lines from individuals born in England, Wales and Scotland during 1 week in 1958. The 4,818 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 obtained from individuals with self-reported British ancestry.

A summary of cases 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 and 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 standards if the original DNA concentration was  $\geq 50$  ng/ $\mu$ l, 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 quality control, raw intensities (from the .CEL files) were renormalized within collections using CelQuantileNorm. These normalized intensities were used to call genotypes with an updated version of the Chiamo software 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 used thresholds for assigning 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 quality control filters described below, this threshold led to a study-wide level of missing data of 0.20%.

An overlapping set of 4,830 controls was 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 quality control 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 ref. 33) by making genotype calls at different confidence thresholds and then plotting the fraction of calls made against concordance with the Illumina genotypes (**Supplementary Fig. 2**).

The general trend was 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 strongly 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 quality control measures (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 we excluded SNPs for which 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 the data typically violate assumptions underpinning standard tests for association, the payback in terms of increased accuracy for these tests can be substantial.

To try to obtain the maximally powerful set of samples and SNPs, we attempted to refine some standard quality control practices. For all individuals, we explicitly modeled the data as a mixture of 'normal' and 'outlier' individuals for each of ancestry, missing data and heterozygosity, and sex assignment (unpublished). 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 among 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 allowed more refined assessment of the relatedness between individuals than genome-wide sharing statistics (for example, parent-child relationships can be distinguished from those of siblings). We obtained a set of individuals with identity by descent  $< 5\%$  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 Fig. 3**). The measure is calculated automatically by the program SNPtest. SNPs were removed if this information measure was below 0.98 or if the estimated minor allele frequency 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 we excluded 16 from further analysis as likely genotyping errors.

**Supplementary Figure 4** provides quantile-quantile plots for the post-quality control comparison of our two control collections and for association statistics based on the post-quality control trend test comparing cases and the combined control set. Both visual inspection and the inflation statistics for each ( $\lambda = 1.037$  and  $\lambda = 1.079$ , respectively) suggest that the quality control 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 SNPTTEST and PLINK<sup>34</sup>. We also performed 2-d.f. genotypic tests to verify that none of our associations show significant deviation from a multiplicative model, and we performed

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://outmod.edbonsai.sourceforge.net>; Chiamo, <http://www.stats.ox.ac.uk/~marchini/>

[software/gwas/chiamo.html](http://www.stats.ox.ac.uk/~marchini/software/gwas/chiamo.html); SNPtest, <http://www.stats.ox.ac.uk/~marchini/software/gwas/snptest.html>; PoBI <http://www.peopleofthebritishisles.org>.

33. Korn, J.M. *et al.* Integrated genotype calling and association analysis of SNPs, common copy number polymorphisms and rare CNVs. *Nat. Genet.* **40**, 1253–1260 (2008).
34. Purcell, S. *et al.* PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).

