# Supplement for:

**Genome-Wide Meta-Analysis Increases to 71 the Number of Confirmed Crohn's Disease Susceptibility Loci**

Andre Franke[1,*], Dermot P.B. McGovern[2,3*], Jeffrey C. Barrett[4,*], Kai Wang[5], Graham L. Radford-Smith[6], Tariq Ahmad[7], Charlie W. Lees[8], Tobias Balschun[9], James Lee[10], Rebecca Roberts[11], Carl A. Anderson[4], Joshua C. Bis[12], Suzanne Bumpstead[4], David Ellinghaus[1], Eleonora M. Festen[13], Michel Georges[14], Todd Green[15], Talin Haritunians[3], Luke Jostins[4], Anna Latiano[16], Christopher G. Mathew[17], Grant W. Montgomery[18], Natalie J. Prescott[17], Soumya Raychaudhuri[15], Jerome I. Rotter[3], Philip Schumm[19], Yashoda Sharma[20], Lisa A. Simms[6], Kent D. Taylor[3], David Whiteman[18], Cisca Wijmenga[13], Robert N. Baldassano[21], Murray Barclay[11], Theodore M. Bayless[22], Stephan Brand[23], Carsten Buning[24], Albert Cohen[25], Jean-Frederick Colombel[26], Mario Cottone[27], Laura Stronati[28], Ted Denson[29], Martine De Vos[30], Renata D'Inca[31], Marla Dubinsky[32], Cathryn Edwards[33], Tim Florin[34], Denis Franchimont[35], Richard Gearry[11], Jürgen Glas[23, 36, 37], Andre Van Gossum[35], Stephen L. Guthery[38], Jonas Halfvarson[39], Hein W. Verspaget[40], Jean-Pierre Hugot[41], Amir Karban[42], Debby Laukens[30], Ian Lawrance[43], Marc Lemann[44], Arie Levine[45], Cecile Libioulle[46], Edouard Louis[46], Craig Mowat[47], William Newman[48], Julián Panés[49], Anne Phillips[47], Deborah D. Proctor[20], Miguel Regueiro[50], Richard Russell[51], Paul Rutgeerts[52], Jeremy Sanderson[53], Miquel Sans[49], Frank Seibold[54], A. Hillary Steinhart[55], Pieter C.F. Stokkers[56], Leif Torkvist[57], Gerd Kullak-Ublick[58], David Wilson[59], Thomas Walters[60], Stephan R. Targan[2], Steven R. Brant[22], John D. Rioux[61], Mauro D'Amato[62], Rinse K. Weersma[63], Subra Kugathasan[64], Anne M. Griffiths[60], John C. Mansfield[65], Severine Vermeire[52], Richard H. Duerr[50,66], Mark S. Silverberg[55], Jack Satsangi[8], Stefan Schreiber[1,67], Judy H. Cho[20,68], Vito Annese[16,69], Hakon Hakonarson[5,21], Mark J. Daly[15, †], Miles Parkes[10,†]


[1-69]   see affiliation details in main manuscript

*   Shared first authorship
†   Shared senior authorship

‡   Corresponding author
    Inflammatory Bowel Disease Research Group
    Addenbrooke's Hospital
    University of Cambridge
    Cambridge CB2 2QQ
    United Kingdom
    eMail:   miles.parkes@addenbrookes.nhs.uk
    Tel.:    +44 (0) 1223-216389
    Fax:     +44 (0) 1223 596213

**Supplementary Table 1** - Index GWAS studies used for the meta-analysis, with genotyping platform and numbers of cases and controls from each centre for which (post quality control) genome-wide association study data were available.

| Index GWAS | Crohn's disease cases | Healthy controls | GWAS platform |
| --- | --- | --- | --- |
| Early Onset | 1,689 | 6,197 | Illumina HumanHap550 |
| German | 479 | 1,145 | Illumina HumanHap550 |
| USA (Cedars-Sinai) | 925 | 2,882 | Illumina HumanHap300 |
| Belgium | 537 | 913 | Illumina HumanHap300 |
| USA (NIDDK) | 956 | 982 | Illumina HumanHap300 |
| UK (WTCCC) | 1,747 | 2,937 | Affymetrix GeneChip 500 |
| **TOTAL** | 6,333 | 15,056 | |

**Supplementary Table 2** - Country of origin of samples used in the replication experiment, with employed genotyping platform.

| REPLICATION - country | Crohn's disease cases | Healthy controls | Genotyping platform |
|---|---|---|---|
| Australia | 1,357 | 1,923 | Sequenom iPlex |
| Belgium | 1,282 | 1,682 | SNPlex/Taqman |
| France | 414 trios | | SNPlex/Taqman |
| Germany | 3,808 | 2,747 | SNPlex/Taqman |
| Israel | 444 | 376 | SNPlex/Taqman |
| Italy | 921 | 899 | SNPlex/Taqman |
| Netherlands | 1,101 | 269 | SNPlex/Taqman |
| New Zealand | 514 | 457 | Sequenom iPlex |
| Spain | 325 | 987 | SNPlex/Taqman |
| Sweden | 724 | 992 | SNPlex/Taqman |
| UK | 3,243 | 2,431 | Sequenom iPlex |
| USA (Cedars-Sinai) | 1,172 | 501 | SNPlex/Taqman |
| USA (NIDDK) | 803 | 762 | Illumina Goldengate |
| TOTAL | 15,694 | 14,026 | |

**Supplementary Table 3 – Odds ratios (OR) and risk allele frequencies (RAF) for the 71 SNPs listed in Table 1 and 2.**

Column **Repl. Heterogeneity** lists the Breslow-Day heterogeneity chi-square with 11 degrees of freedom. To avoid any bias due to winner's curse we performed this test in replication datasets only (available for 35/71 loci). No significant heterogeneity was observed after correction for multiple tests. For the Belgian and Cedar_2 GWAS sample, which had high genomic inflation factors (see **Supplementary Figure 1**), PCA-corrected $P$-values are shown (**P_corr**) besides the uncorrected $P$-values (**P_uncorr**) from **Table 1 and 2.** These $P$-value columns are also highlighted by light red shading. The comparison demonstrates that the observed associations do not arise due to population stratification. **Chr.:** chromosome

This supplementary table is available for download as an Excel file.

**Supplementary Table 4 – Raw allele counts and empirical variance for the 71 SNPs listed in Table 1 and 2.**

**Chr.:** chromosome; **A/B**: counts per allele (aligned to + strand of NCBI's build 36) for healthy controls (**CTRL**) and Crohn's disease patients (**CASE**); **EMPVAR**: empirical variance. The empirical variance is calculated as:

1) Conversion of genotypes to allelic dosages, which are the weighted sum of the genotype class probabilities, e.g. if "A" is the reference allele and genotype probabilities are 0.8 AA 0.1 Aa and 0.1 aa, then the dosage is calculated as

$$2 \times 0.8 + 1 \times 0.1 + 0 \times 0.1 = 1.7$$

2) Calculation of dosages for all individuals.
3) Calculation of the mean dosage for all individuals.
4) Calculation of the square of the deviation from the mean for each individual dosage.
5) Sum the squared deviations and division by twice the number of individuals.

This supplementary table is available for download as an Excel file.

**Supplementary Table 5 – Details on 1000G coding SNPs (cSNP) that are in LD with the SNPs listed in Table 1 and 2.**

**Chr.:** chromosome; Position (**pos**) are shown according to NCBI's dbSNP build 130.

| No. | Table 1 | lead SNP | | | | coding SNP | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | bSNP ID | #HR. | os (dbSNP 130) | Risk allele - Allele frequency in control population | dbSNP ID | pos (dbSNP 130) | distance to lead SNP [kb] | RefGene | amino acid substitution |
| 1 | 6 | rs780093 | 2p23 | 27,742,603 | T - 0.418 | rs1260326 | 27,730,940 | 11.66 | GCKR | L446P |
| 2 | 7 | rs10495903 | 2p21 | 43,660,422 | T - 0.129 | rs35720761; rs7578597 | 43,373,481; 43,586,327 | 286.94; 74.10 | THADA | C1605Y; T1187A |
| 3 | 10 | rs6738825 | 2q33 | 198,605,140 | A - 0.473 | rs1064213 | 198,658,485 | 53.35 | PLCL1 | V667I |
| 4 | 14 | rs2549794 | 5q15 | 96,270,305 | C - 0.409 | rs2549782 | 96,256,756 | 13.55 | ERAP2 | K392N |
| 5 | 21 | rs4077515 | 9q34 | 138,386,317 | T - 0.411 | same | | | CARD9 | S12N |
| | | | | | | rs3812571 | 138,395,115 | 8.80 | SNAPC4 | H799Q |
| 6 | 29 | rs8005161 | 14q35 | 87,542,348 | T - 0.119 | rs1805078 | 87,520,523 | 21.83 | GALC | R184C |
| | | | | | | rs3742704 | 87,547,635 | 5.29 | GPR65 | I231L |
| 7 | 31 | rs151181 | 16p11 | 28,398,018 | G - 0.386 | rs180743 | 28,415,145 | 17.13 | APOB48R | P419A |
| | | | | | | rs181206 | 28,420,904 | 22.89 | IL27 | L119P |
| | | | | | | rs7498665 | 28,790,742 | 392.72 | SH2B1 | T484A |
| 8 | 33 | rs12720356 | 19p13 | 10,330,975 | G - 0.084 | same | | | TYK2 | I684S |
| 9 | 35 | rs281379 | 19q13 | 53,906,086 | A - 0.487 | rs601338; rs602662 | 53,898,486; 53,898,797 | 7.60; 7.29 | FUT2 | **W154X**; G258S |
| | | | | | | rs2287922 | 53,924,038 | 17.95 | RASIP1 | R601C |
| 10 | 36 | rs4809330 | 20q13 | 61,820,030 | G - 0.709 | rs3208008 | 61,796,554 | 23.48 | RTEL1 | Q1042H |
| 11 | 37 | rs181359 | 22q11 | 20,258,641 | T - 0.203 | rs2298428 | 20,312,892 | 54.25 | YDJC | A263T |
| 12 | 1 | rs11209026 | 1p31 | 67,705,958 | G - 0.932 | same | | | IL23R | R381Q |
| 13 | 2 | rs2476601 | 1p13 | 114,179,091 | G - 0.907 | same | | | PTPN22 | W602R |
| 14 | 5 | rs7554511 | 1q32 | 199,144,185 | C - 0.726 | rs296520 | 199,147,601 | 3.42 | C1orf106 | R453C |
| 15 | 6 | rs3792109 | 2q37 | 233,849,156 | A - 0.529 | rs2241880 | 233,848,107 | 1.05 | ATG16L1 | T300A |
| 16 | 7 | rs3197999 | 3p21 | 49,696,536 | A - 0.297 | same | | | MST1 | R703C |
| | | | | | | rs1050450 | 49,369,838 | 326.70 | GPX1 | P200L |
| | | | | | | rs34762726 | 49,664,214 | 32.32 | BSN | A741T |
| 17 | 9 | rs12521868 | 5q31 | 131,812,292 | T - 0.422 | rs1050152 | 131,704,219 | 108.07 | SLC22A4 | L503F |
| 18 | 13 | rs1799964 | 6p21 | 31,650,287 | C - 0.209 | rs2259435 | 31,604,894 | 45.39 | MCCD1 | E42K |
| | | | | | | rs2229094 | 31,648,535 | 1.75 | LTA | C13R |
| 19 | 24 | rs11564258 | 12q12 | 40,792,300 | A - 0.025 | | | | MUC19 | |
| 20 | 25 | rs3764147 | 13q14 | 43,355,925 | G - 0.245 | same | | | C13orf3 1 | I254V |
| 21 | 27 | rs2872507 | 17q21 | 35,294,289 | A - 0.458 | rs2345480; rs2305479 | 35,315,722; 35,315,743 | 21.43; 21.45 | GSMDL | P289S; G282R |
| | | | | | | rs11557467 | 35,282,160 | 12.13 | ZPBP2 | S173I |
| 22 | 28 | rs11871801 | 17q21 | 37,824,298 | A - 0.756 | rs665268 | 37,975,555 | 151.26 | MLX | Q233R |

**Supplementary Table 6**: Positional candidate genes mapping within regions of confirmed association for Crohn's disease. Three *in silico* techniques were used to further highlight genes of interest (see main text for more details):

(1) 1000 Genomes Project and HapMap databases were searched for coding SNPs in linkage disequilibrium with our most associated SNP.

(2) GRAIL software was used to identify connectivity between genes mapping to different loci.

(3) An eQTL database was searched to identify loci within which the focal SNP correlated with gene expression with LOD≥5.0 (for details see **Supplementary Note**).

This supplementary table is available for download as an Excel file.

**Supplementary Figure 1 -** Quantile-quantile (Q-Q) plot of the association test statistic for all eight constituent GWAS data sets and the meta analysis (meta). For a detailed description of the constituent GWAS see **Supplementary Table 1**. For the Q-Q plots, the Cedar's sample (USA ("Cedars Sinai") in **Supplementary Table 1**) was dividided into the two sub-GWAS ("cedar1" and "cedar2"). For the USA (NIDDK) sample, a separate Q-Q plot is shown for the Jewish ("niddkj") and for the non-Jewish ("niddknj") sub-GWAS. Plots were calculated for all 953,241 SNPs that passed the quality control. For better scaling the y-axis was limited to a maximum chi$^2$ of 30 and the SNPs with higher chi$^2$ values are forming a "plateau" along the top. The over-dispersion of the association test statistic was estimated to be $\lambda_{GC}$=1.27 for the meta-analysis. The shaded region is the 95% concentration band that is formed by calculating, for each order statistic, the 2.5th and 97.5th centiles of the respective distribution under the null hypothesis.

**Supplementary Figure 2 – Results of association analysis ("Manhattan Plot").** The negative common logarithm of the $P$-values for the test statistic using single-SNP Z scores of the genome-wide association study are shown according to chromosome. Only markers that passed the quality criteria were used for plotting (n=953,241). Marker positions are in NCBI's build 36 (hg18).

**Supplementary Figure 3 - Quantile-quantile (Q-Q) plot of interaction analysis for the 71 SNPs listed in Table 1 and 2**. For a detailed description of the analysis see **Online Methods**. No'deviation from the null is observed and no results were significant when considering the number of tests performed.

# CD Interaction
# Analysis

**Supplementary Figure 4 - Regional plots of Table 1 and 2.** Regional Plots of the negative decadic logarithm of the *P*-values obtained in the GWAS in a window of ~250 kb around each of the 71 SNPs displayed in **Table 1 and 2**. The SNP ID of the index SNP is given above each plot and this SNP is marked by a large blue-filled symbol. The magnitude of linkage disequilibrium with the central SNP is reflected by the fill color of the symbols using the measure $r^2$ (for color coding see legend in the upper right corner of each plot). Recombination activity (**cM/Mb**) is depicted by a blue line. Positions and gene annotations are according to NCBI's build 36 (hg18). Plots are ordered according to their order in **Table 1 and 2**.

# Supplementary Note – eQTL Analysis

*Regional annotation via eQTL analysis*

The effects on expression of neighboring genes of the 71 SNPs listed in Table 1 cpf "4"(of the main manuscript) were studied using the Dixon *et al.* transcriptome data, based on Epstein-Barr virus–transformed lymphoblastoid cell lines from 400 children, as described previously [**1**]. The Dixon"*et al.* study genotyped subjects using the Illumina HumanHap300 platform and measured gene expression levels on 54,675 transcripts representing 20,599 genes using the Affymetrix U133 Plus 2.0 GeneChip. For each SNPs listed in Table 1 and 2 of the main manuscript, if it was not genotyped in the Dixon *et al.* study, the best proxy based on CEU $r^2$ was examined instead (see **Table I**). A LOD score threshold of 5 was used to declare a significant association between SNP genotypes and gene expression levels.

To evaluate the significance of the eQTL findings with the CD-associated SNPs, we compared the observed number of *cis* eQTLs yielding LOD scores >5, and the corresponding number among a randomly selected sets of 71 frequency-matched SNPs. The coordinates of the transcripts were provided by Dixon *et al.* (NCBI's 35 assembly), so the coordinates of the SNPs that we used were also based on NCBI 35 assembly. The allele frequency distribution of the 71 SNPs are shown in **Figure I**. The 71 CD SNPs produced seven eQTLs with LOD >5. In 1000 simulated experiments, only 1/1000 produced more than seven eQTLs with LOD >5 (see **Figure II**). Therefore, there is a significant enrichment of cis-operating eQTLs among the CD-associated loci.

**Table I.** Genotype-expression correlation analysis on CD-associated SNPs to identify cis-eQTLs. For primary SNPs that were not genotyped in the Dixon *et al.* data set, a proxy SNP was determined for the analysis, and the pairwise linkage disequilibrium (LD) between the primary and proxy SNPs are shown using the measure $r^2$ (R2). The column "eSNP nearby for same transcript?" shows whether nearby SNPs with stronger eQTL effects are present for the same transcript in the "eQTL?" column.

| Chr. | Left | Right | No. P<$10^{-5}$ | eQTL? | eSNP nearby for same transcript? | eSNP tested | Same SNP or Proxy |
|------|------|-------|-----------------|-------|-----------------------------------|-------------|-------------------|
| 1 | 204,869,063 | 205,099,374 | 1 | None | | rs3024505 | same SNP |
| 1 | 153,244,553 | 154,389,854 | 83 | None | | rs3180018 | same SNP |
| 1 | 7,660,386 | 7,891,753 | 5 | None | | rs2797685 | same snp |
| 1 | 195,577,400 | 196,205,063 | 23 | None | | rs1998598 | same SNP |
| 2 | 43,301,270 | 43,795,977 | 38 | None | | rs10495903 | same SNP |
| 2 | 230,761,562 | 230,944,130 | 13 | *SP140* (LOD=8.8) | | rs13397985 | rs13397985 - R2 = 0.89 |
| 2 | 25,301,554 | 25,461,262 | 2 | None | | rs7583409 | rs7583409 - R2 = 1.0 |
| 2 | 197,850,455 | 198,667,617 | 47 | None | | rs1541953 | rs1541953 - R2 = 1.0 |
| 3 | 18,582,796 | 18,857,492 | 5 | None | | rs6792314 | rs6792314 - R2 = 0.93 |
| 5 | 173,150,790 | 173,472,805 | 3 | *CPEB4* (LOD=6.1) | rs747472 (LOD=12.0) | rs359457 | same SNP |
| 5 | 72,492,752 | 72,616,167 | 1 | None | | rs7702331 | same SNP |
| 5 | 96,105,158 | 96,450,821 | 36 | *LRAP* (LOD=47.2) | rs39602 (LOD=50.3) | rs27306 | rs27306 - R2 = 1.0 |
| 5 | 141,394,644 | 141,622,106 | 10 | None | | rs11167764 | same SNP |
| 6 | 159,260,927 | 159,464,567 | 2 | None | | rs212388 | same SNP |
| 6 | 90,863,556 | 91,139,647 | 6 | None | | rs1010474 | rs1010474 - R2 = 0.92 |
| 7 | 152,992,808 | 153,135,560 | 4 | None | | rs2098112 | same SNP |
| 8 | 129,560,876 | 129,668,460 | 11 | None | | rs6651253 | rs6651253 - R2 = 1.0 |
| 9 | 138,274,802 | 138,544,419 | 41 | *CARD9* (LOD=12.4) | rs11794847 (LOD=13.1) | rs4077515 | same SNP |
| 10 | 80,671,808 | 80,775,690 | 11 | None | | rs1250550 | same SNP |
| 10 | 59,493,900 | 59,815,801 | 34 | None | | rs1698408 | rs1698408 - R2 = 0.80 |
| 10 | 6,070,249 | 6,205,531 | 1 | None | | rs122722561 | rs122722561 - R2 = 0.93 |
| 11 | 61,283,132 | 61,441,126 | 13 | *FADS1* (LOD=5.0) | rs174578 (LOD=7.2) | rs102275 | same SNP |
| 13 | 41,724,842 | 41,999,763 | 5 | *TNFSF11* (LOD=5.9) | | rs9594759 | rs9594759 - R2 = 0.97 |
| 14 | 87,280,056 | 87,712,056 | 17 | None | | rs3742704 | rs3742704 - R2 = 1.0 |
| 14 | 68,226,845 | 68,387,815 | 12 | None | | rs194772 | rs194772 - R2 = 0.67 |
| 15 | 65,195,295 | 65,269,057 | 13 | None | | rs16950687 | rs16950687 - R2 = 0.88 |
| 16 | 28,202,322 | 28,935,308 | 33 | *EIF3S8* (LOD=11.3) | rs7189927 (LOD=11.9) | rs4788084 | rs4788084 - R2 = 0.89 |
| 19 | 10,261,304 | 10,495,264 | 9 | None | | rs12720356 | same SNP |
| 19 | 53,784,242 | 53,969,894 | 13 | None | | rs676388 | rs676388 - R2 = 1.0 |
| 22 | 37,984,993 | 38,135,698 | 8 | None | | rs2413583 | same SNP |
| 22 | 20,141,991 | 20,393,337 | 26 | None | | rs5754217 | rs5754217 - R2 = 1.0 |
| 22 | 28,232,382 | 28,998,308 | 8 | None | | rs757024 | rs757024 - R2 = 0.90 |

**Figure I.** The minor allele frequency (MAF) distribution of the 71 SNPs associated with CD in Table 1 and 2.