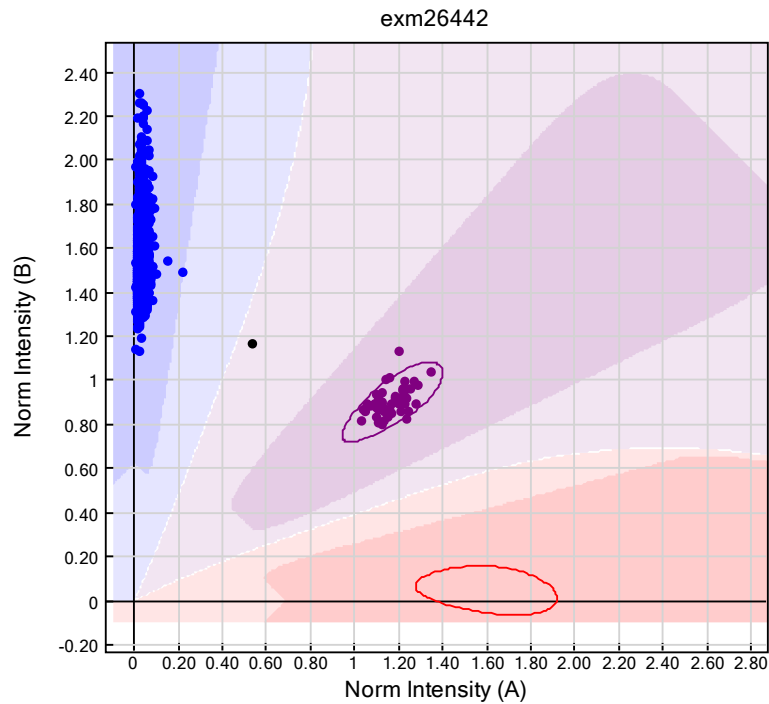
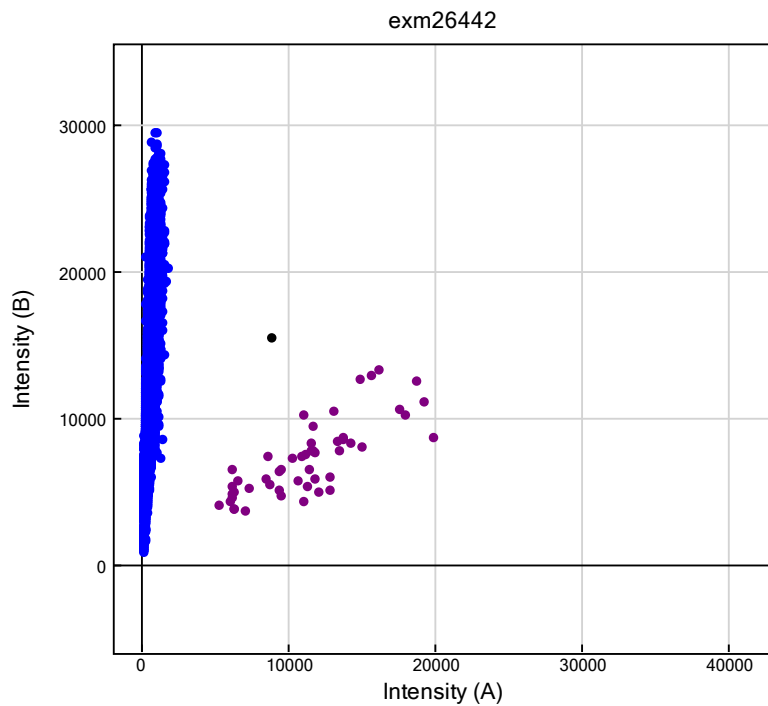


1 Supplementary Figures



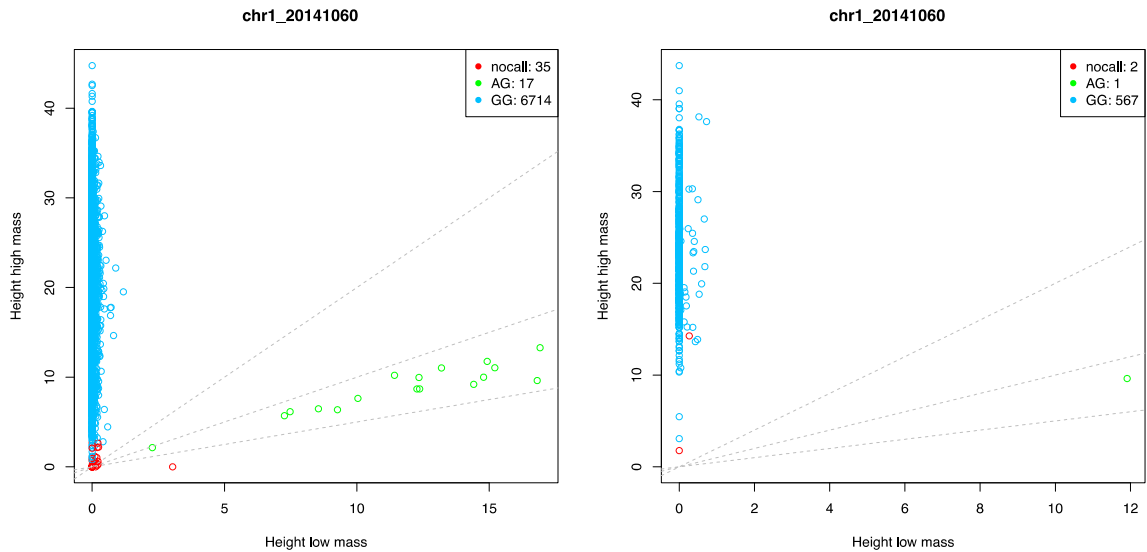
2



3

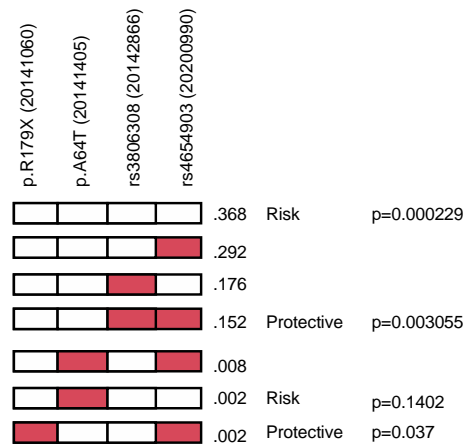
4 Supplementary figure 1. Cluster plots for rs36095412 (exm26442) from the Illumina
5 HumanExome BeadChip array. Cartesian coordinates display the cluster using
6 intensity values A and B representing the two possible alleles for this SNP. The top
7 and bottom panels represent normalized and raw values, respectively. Red and blue
8 shaded regions represent the two homozygous clusters; purple shaded region
9 represents the heterozygous cluster.

10



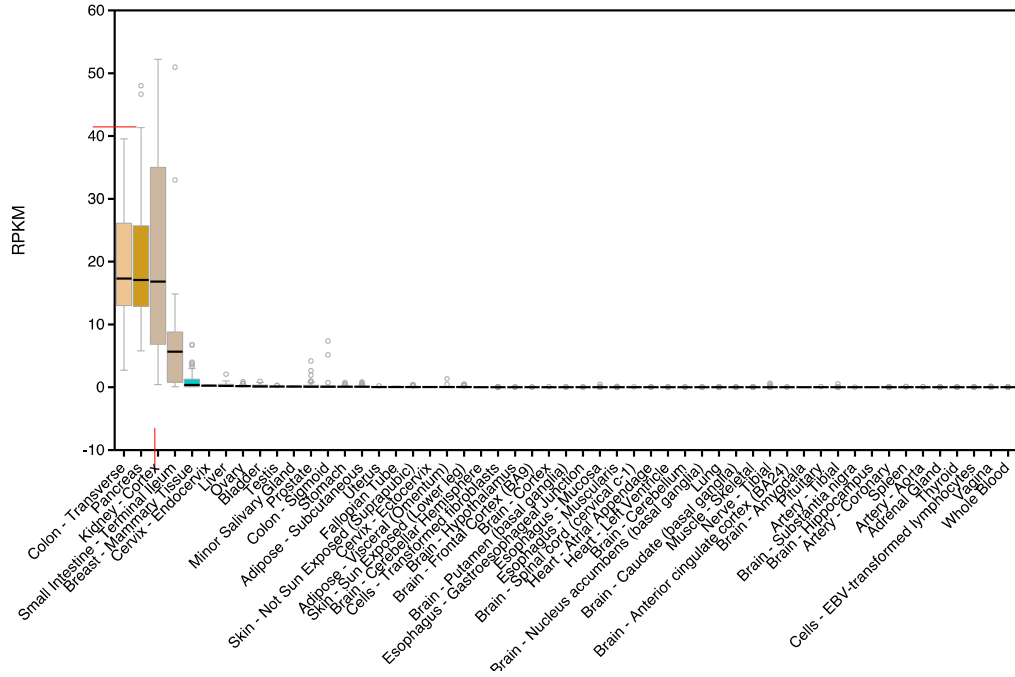
11
12
13
14
15
16
17
18
19
20

Supplementary figure 2. Cluster plots for rs36095412 (chr1_20141060) from the Sequenom genotyping. Using R and the raw data (heights of MASSSpec intensities): we obtained *skew*, where we divide allele with higher intensity by sum of height of both alleles; and *yield*, where we divide 1- height of unextended primer by sum of intensities of both alleles and unextended primer. *Yield* indicates quality of signal where a value below 0.5 indicates poor quality. *Height* refers to height of signal in the MassSpec.

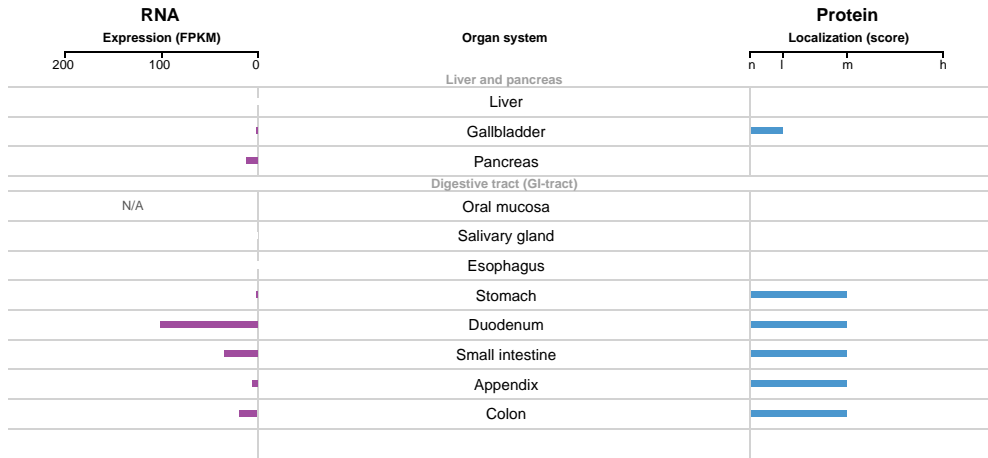


21
22
23
24
25
26
27
28
29
30
31
32
33

Supplementary figure 3. Haplotype blocks of RNF186 associated variants including: the low-frequency coding variants (1) p.R179X and (2) p.A64T; and the common associated variants (3) rs3806308 and (4) rs4654903. Haplotype-based case-control association analysis was conducted using PLINK 1.07 using a subset of individuals with array¹ and targeted genotyping data (red indicates the non-reference allele).



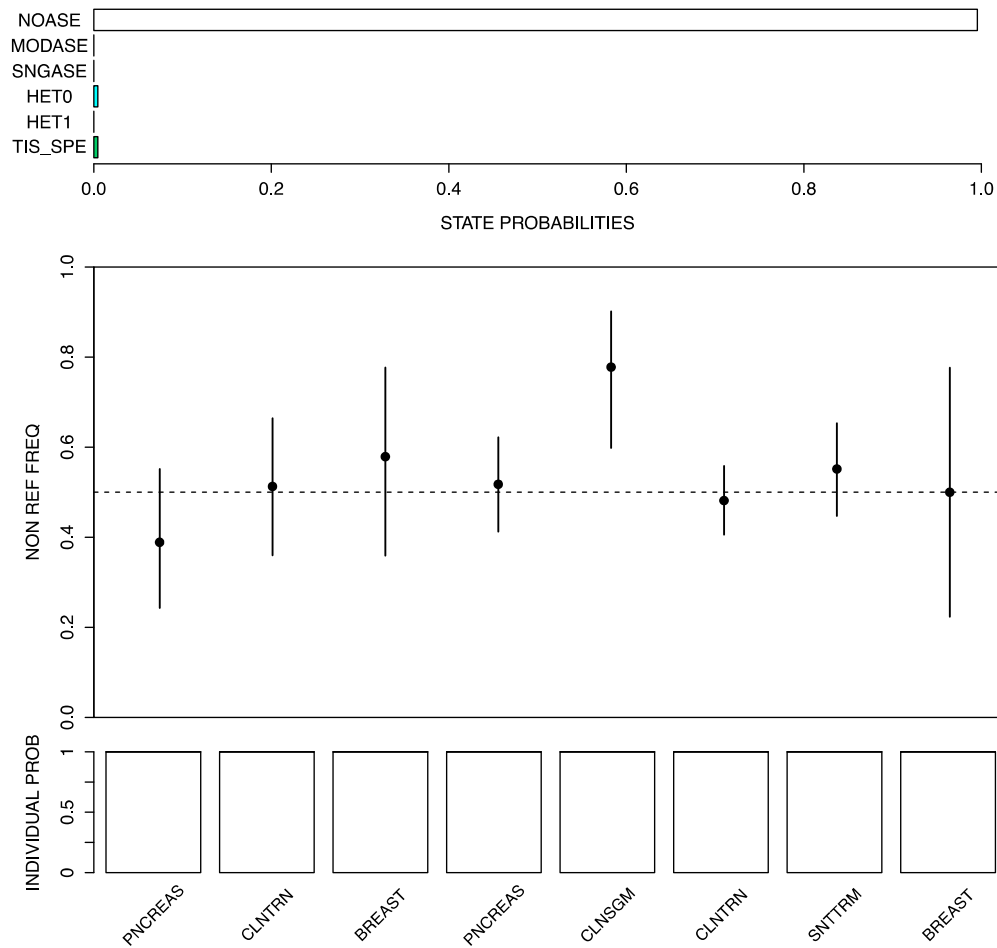
35
36
37
38



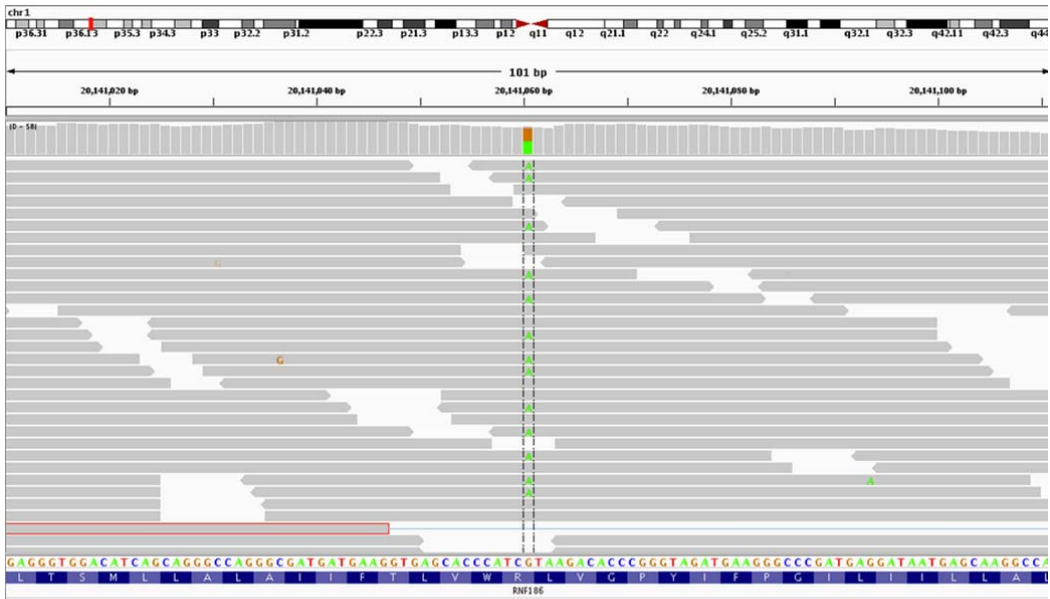
39
40
41
42
43
44
45
46
47
48

Supplementary figure 4. Tissue-wide RNA expression profile of *RNF186* in the Genotype-Tissue Expression² project (<http://gtexportal.org/home/gene/RNF186>). Tissues are sorted in decreasing order by the median RPKM value from GTEx Analysis Release V4. Protein expression profile of RNF186 protein in the human protein atlas³ shows “medium” localization score in the digestive tract.

1_20141060_G_A_b37

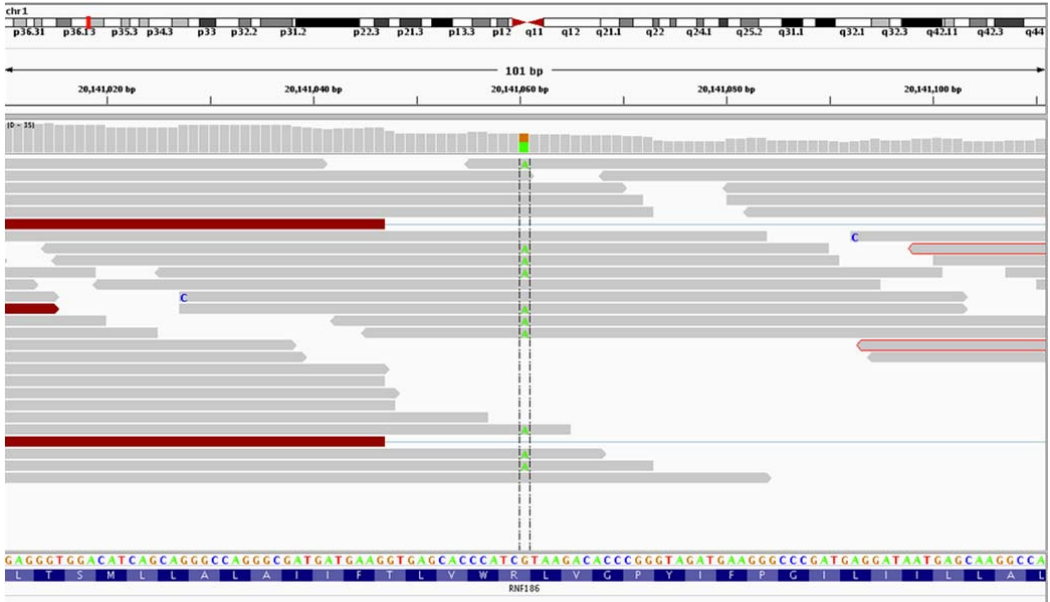


49
 50 Supplementary figure 5. Allele-specific expression data for rs36095412 (p.R179X) in
 51 GTEx. Top panel shows the posterior probabilities for six states as defined in Pirinen
 52 et al⁴. The multi-tissue classification state ‘NOASE’ (no ASE effect across all tissues)
 53 has posterior probability greater than 0.9. Middle panel shows the point estimates of
 54 the non-reference allele frequency among RNA-seq reads across eight observations
 55 (in five different tissue types named at the bottom) together with their 95% credible
 56 intervals. Bottom panel shows the posterior probability of the group indication for
 57 each tissue type, where white, gray, and black denote groups no ASE, moderate ASE,
 58 and strong ASE, respectively.
 59



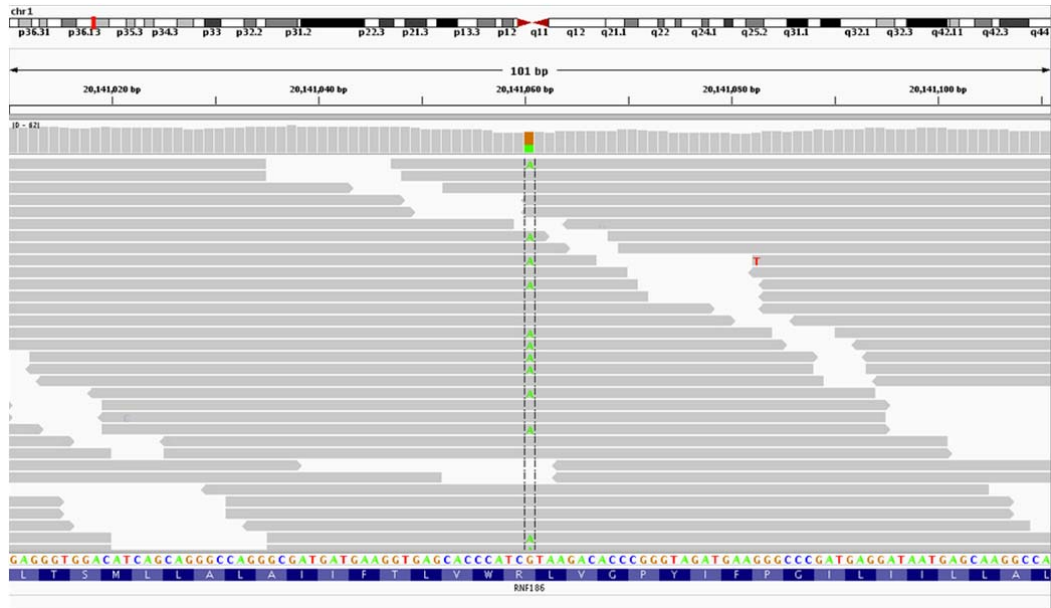
60
61

Colon transverse



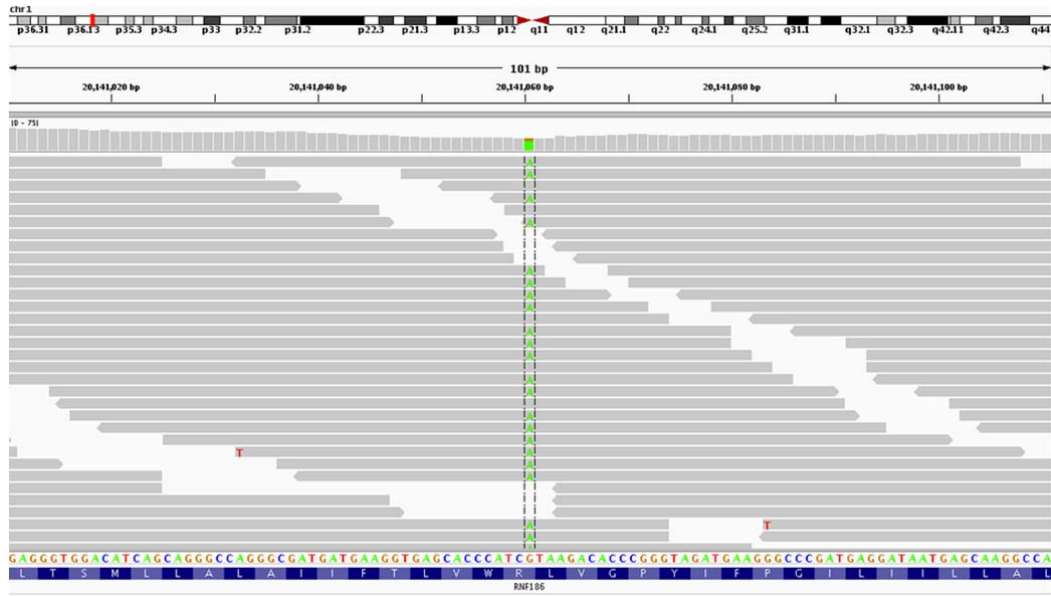
62
63
64
65

Breast mammary tissue



66
67

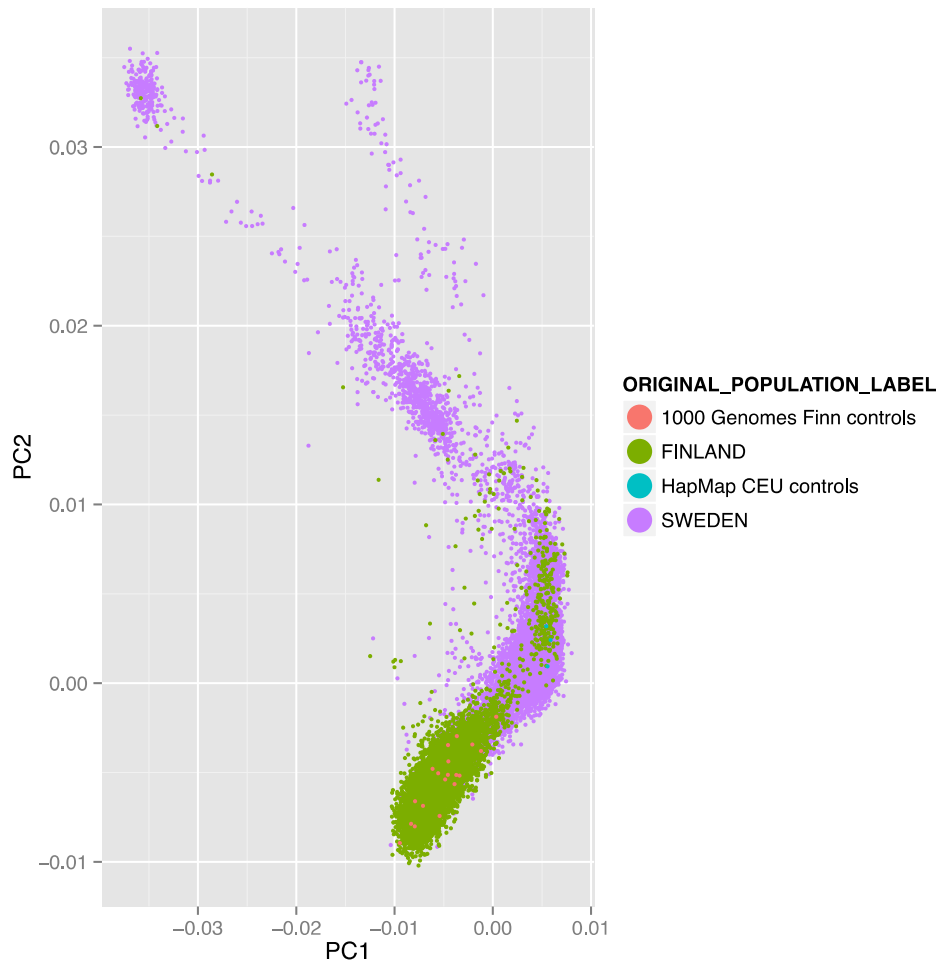
Pancreas



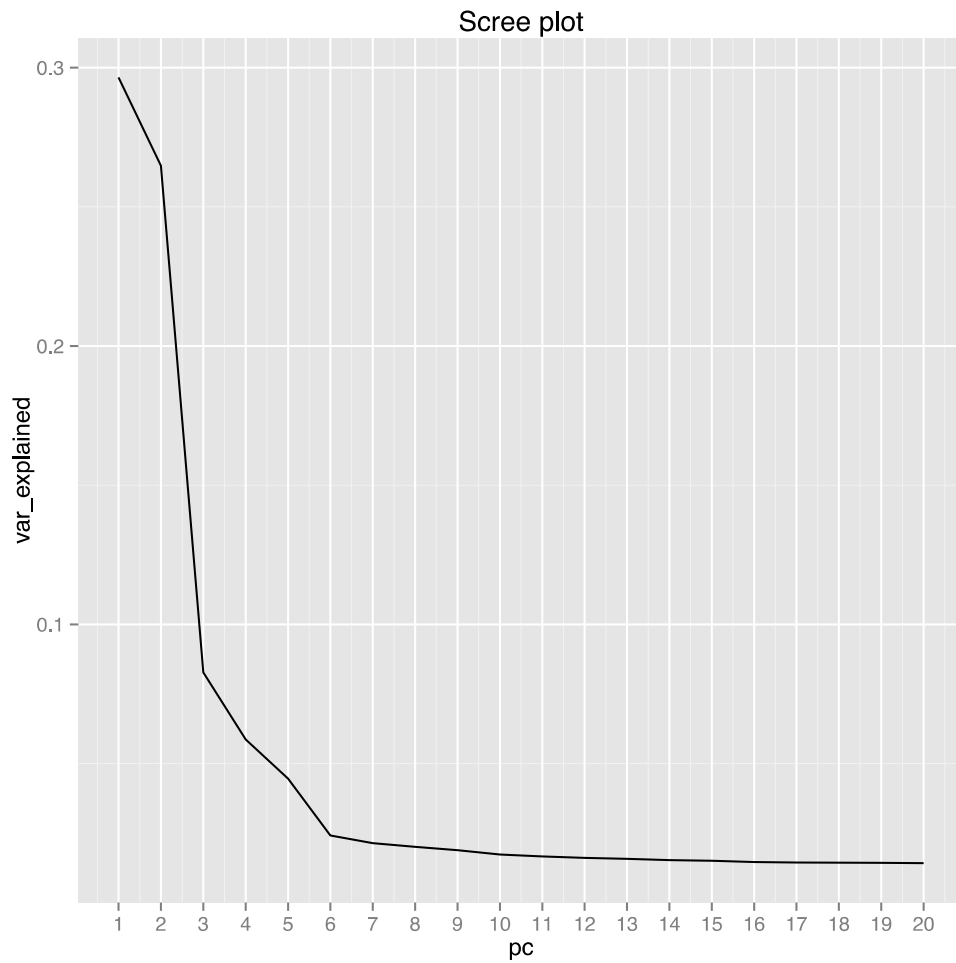
68
69
70
71

Colon sigmoid

Supplementary figure 6. IGV snapshot of RNA-seq reads from tissues with expression (RPKM) greater than zero in p.R179X carriers from the GTEx project.



72
73 Supplementary figure 7. First two principal components showing genetic differences
74 among 27885 jointly called individuals for classifying Finnish individuals among
75 Swedish samples. Original population labels are given in different colors (Finnish,
76 Swedish). Additionally, well-characterized control individuals were added to joint
77 calling (Finnish individuals from 1000 genomes project and HapMap CEU
78 individuals) and are shown in the figure as separate colors.
79



80
81 Supplementary figure 8. Proportion of genetic variance explained by first 20 principal
82 components among 27885 sequenced individuals. In x-axis principal components are
83 ordered by decreasing variance explained and y-axis gives variance explained by each
84 principal component.

85
86
87
88
89
90
91
92
93

94 **Supplementary Tables**

95 Supplementary table 1. Prioritizing protective protein truncating variants identified in
 96 the targeted sequencing data set (CMH test). For each variant we present the analysis
 97 using the indexed association in Huang et al.¹ For each data set the single variant
 98 association analysis output from PLINK/SEQ is shown
 99 (<https://atgu.mgh.harvard.edu/plinkseq/assoc.shtml#single>). CONMETA – consensus
 100 sample-variant meta-information; ALT – alternate allele(s), comma delimited; MAF –
 101 minor allele frequency; HWE – *P*-value from Hardy-Weinberg disequilibrium test
 102 (exact test); MINA – number of minor alleles in cases; MINU – number of minor
 103 alleles in controls; OBSA – number of non-null genotypes in cases; OBSU – number
 104 of non-null genotypes in controls; REFA – number of reference homozygotes in
 105 cases; HETA – number of heterozygotes in cases; HOMA – number of alternate
 106 homozygotes in cases; REFU, HETU, HOMU – same as aforementioned description,
 107 but for controls; *P* – *p*-value for single site association (allelic, two-sided test); OR –
 108 Allelic odds ratio. Finnish exome data is also shown.

109
 110 Supplementary table 2. Conditional analysis of p.R179X variant and the index
 111 common variant association rs4654903 reported in Silverberg et al (2009).⁵
 112 Conditional analysis results for the samples in the Iceland replication data set with
 113 whole genome sequencing data are shown.
 114

	rs36095412-A stop gained R179X (MAF 0.78%)		rs4654903-A Intergenic MAF(45.5%)		Phenotype (N)
Analysis	Pvalue	Effect	Pvalue	Effect	
Unadjusted	5.0 x 10 ⁻⁴	0.30	3.8x10 ⁻⁷	1.24	Ulcerative Colitis (N Cases=1,453 ;Ctrls=264,744)
Adjusted for the other marker	8.4 x 10 ⁻⁴	0.31	9.0x10 ⁻⁷	1.23	

115 $r^2 = 0.004$; $D' = 0.83$

116

117 Supplementary table 3. Association of p.R179X in *RNF186* with crohn's disease

118

119

Study	Data type	CD		Controls		Control MAF	<i>P</i>	OR
		179X	R179	179X	R179			
GWASseq	Sequence (targeted)	4	2404	6	1828	0.33%		
Finland	Sequence (exome)	0	476	23	16223	0.14%		
Screen	-	4	2880	29	18051		0.19	0.36
US+Canada	Exome Chip	16	9962	21	12883	0.16%		
Sweden	Exome Chip	6	1108	45	10813	0.41%		
Belgium	Genotyping	3	3189	0	1764	0.00%		
Germany	Genotyping	4	2822	7	4399	0.16%		
Dutch	Genotyping	6	2306	8	4164	0.19%		
Italy	Genotyping	3	2249	2	1914	0.10%		
Replication							0.56	1.16 (0.76-1.76)
Combined (screen +replication)							0.94	1.04 (0.70-1.54)

120

121

CD, crohn's disease; OR, odds ratio; *P*, *p*-value. Screen + replication *P* value is computed using Mantel-Haenszel chi-squared test with continuity

122

123 **Supplementary Notes**

124

125 **Supplementary Note 1: Pre-processing.** The sequence reads are first mapped to the
126 reference to produce a file in SAM/BAM format sorted by coordinate. Duplicate reads
127 are marked – these reads are not informative and are not used as additional evidence
128 for or against a putative variant. Next, local realignment is performed around indels.
129 This identifies the most consistent placement of the reads relative to potential indels
130 in order to clean up artifacts introduced in the original mapping step. Finally, base
131 quality scores are recalibrated in order to produce more accurate per-base estimates of
132 error emitted by the sequencing machines.

133 **Supplementary Note 2: Variant Discovery.** Once the data has been pre-processed as
134 described above, it is put through the variant discovery process, i.e. the identification
135 of sites where the data displays variation relative to the reference genome, and
136 calculation of genotypes for each sample at that site. The variant discovery process is
137 decomposed into separate steps: variant calling (performed per-sample), joint
138 genotyping (performed per-cohort) and variant filtering (also performed per-cohort).
139 The first two steps are designed to maximize sensitivity, while the filtering step aims
140 to deliver a level of specificity that can be customized for each project.

141 Variant calling is done by running the HaplotypeCaller in GVCF mode on each
142 sample's BAM file(s) to create single-sample gVCFs. If there are more than a few
143 hundred samples, batches of ~200 gVCFs are merged hierarchically into a single
144 gVCF to make the next step more tractable. Joint genotyping is then performed on the
145 gVCFs of all available samples together in order to create a set of raw SNP and indel
146 calls. Finally, variant recalibration is performed in order to assign a well-calibrated
147 probability to each variant call in a raw call set, and to apply filters that produce a
148 subset of calls with the desired balance of specificity and sensitivity.

149 **Supplementary Note 3: Identification of Finnish samples.** Initial data set consisted
150 of 27885 jointly called individuals from Finnish and Swedish cohorts. GATK PASS
151 SNPs were extracted that satisfied the following conditions: minor allele frequency >
152 0.05, HWE-p-value < 1e-6, missing genotypes <= 0.03 (after setting GQ<20 to
153 missing). Remaining variants were LD pruned so that they were approximately
154 independent ($R^2 < 0.1$ within 500kb). Remaining SNPs were used for PCA-analysis.

155

156 Majority of the Finnish and Swedish samples clustered in to clear separate clusters
157 based on PC1 and PC2 (Supplementary figure 7) as they explained over 50% of the
158 variance after which there was a clear drop in variance explained (PC3 explained
159 only 8%, Supplementary figure 8). For objectively classifying Finns we used 100 fold
160 cross validation in logistic regression framework to estimate weights for the PC1 and
161 PC2. In each round of cross validation PC1 and PC2 were regressed on population
162 label and the Z-scores were stored. PC score was calculated by dividing the mean of
163 the Z-scores for that PC by the standard deviation in cross validation replicated. Final
164 weight was obtained by dividing a PC score by the sum of all PC scores. The weights
165 were used to calculate weighted mahalanobis distance of each sample to the centroid
166 of each population learning samples. Probability of sample coming from each

167 population was calculated by squaring the mahalanobis distance and getting
168 cumulative density at that value from chisquare distribution with two degrees of
169 freedom.

170

171 **Supplementary References**

172

- 173 1. Huang, H. *et al.* Association mapping of inflammatory bowel disease loci
174 to single variant resolution. *bioRxiv* (2015).
- 175 2. Consortium, G.T. Human genomics. The Genotype-Tissue Expression
176 (GTEx) pilot analysis: multitissue gene regulation in humans. *Science* **348**,
177 648-60 (2015).
- 178 3. Uhlen, M. *et al.* Proteomics. Tissue-based map of the human proteome.
179 *Science* **347**, 1260419 (2015).
- 180 4. Pirinen, M. *et al.* Assessing allele-specific expression across multiple
181 tissues from RNA-seq read data. *Bioinformatics* **31**, 2497-504 (2015).
- 182 5. Silverberg, M.S. *et al.* Ulcerative colitis-risk loci on chromosomes 1p36
183 and 12q15 found by genome-wide association study. *Nat Genet* **41**, 216-
184 20 (2009).

185