

1 **Leveraging the resolution of RNA-Seq markedly increases the number of**
2 **causal eQTLs and candidate genes in human autoimmune disease**

3 Mapping eQTLs in autoimmune disease using RNA-Seq

4

5 Christopher A. Odhams¹, Deborah S. Cunninghame Graham^{1,2}, Timothy J. Vyse^{1,2*}

6

7 ¹ Department of Medical & Molecular Genetics, King's College London, London, UK

8 ² Academic Department of Rheumatology, Division of Immunology, Infection and Inflammatory

9 Disease, King's College London, London, UK

10

11 * Corresponding author

12 Email: timothy.vyse@kcl.ac.uk (TJV)

13

14 **Abstract**

15 Genome-wide association studies have identified hundreds of risk loci for autoimmune disease, yet
16 only a minority (~25%) share a single genetic effect with changes to gene expression (eQTLs) in
17 primary immune cell types. RNA-Seq based quantification at whole-gene resolution, where
18 abundance is estimated by culminating expression of all transcripts or exons of the same gene, is
19 likely to account for this observed lack of colocalisation as subtle isoform switches and expression
20 variation in independent exons are concealed. We perform integrative *cis*-eQTL analysis using
21 association data from twenty autoimmune diseases (846 SNPs; 584 independent loci), with RNA-Seq
22 expression from the GEUVADIS cohort profiled at gene-, isoform-, exon-, junction-, and intron-level
23 resolution. After testing for a shared causal variant, we found exon-, and junction-level analyses
24 produced the greatest frequency of candidate-causal *cis*-eQTLs; many of which were concealed at
25 whole-gene resolution. In fact, only 9% of autoimmune loci shared a disease-relevant eQTL effect at
26 gene-level. Expression profiling at all resolutions however was necessary to capture the full array of
27 eQTL associations, and by doing so, we found 45% of loci were candidate-causal *cis*-eQTLs. Our
28 findings are provided as a web resource for the functional annotation of autoimmune disease
29 association studies (www.insidegen.com). As an example, we dissect the genetic associations of
30 Ankylosing Spondylitis as only a handful of loci have documented causative relationships with gene
31 expression. We classified fourteen of the thirty-one associated SNPs as candidate-causal *cis*-eQTLs.
32 Many of the newly implicated genes had direct relevance to inflammation through regulation of TNF
33 signalling (for example *NFATC2IP*, *PDE4A*, and *RUSC1*), and were supported by integration of
34 functional genomic data from epigenetic and chromatin interaction studies. We have provided a
35 deeper mechanistic understanding of the genetic regulation of gene expression in autoimmune disease
36 by profiling the transcriptome at multiple resolutions.

37

38 **Author Summary**

39 It is now well acknowledged that non-coding genetic variants contribute to susceptibility of
40 autoimmune disease through alteration of gene expression levels (eQTLs). Identifying the variants
41 that are causal to both disease risk and changes to expression levels has not been easy and we believe
42 this is in part due to how expression is quantified using RNA-Sequencing (RNA-Seq). Whole-gene
43 expression, where abundance is estimated by culminating expression of all transcripts or exons of the
44 same gene, is conventionally used in eQTL analysis. This low resolution may conceal subtle isoform
45 switches and expression variation in independent exons. Using isoform-, exon-, and junction-level
46 quantification can not only point to the candidate genes involved, but also the specific transcripts
47 implicated. We make use of existing RNA-Seq expression data profiled at gene-, isoform-, exon-,
48 junction-, and intron-level, and perform eQTL analysis using association data from twenty
49 autoimmune diseases. We find exon-, and junction-level thoroughly outperform gene-level analysis,
50 and by leveraging all five quantification types, we find 45% of autoimmune loci share a single genetic
51 effect with gene expression. We highlight that existing and new eQTL cohorts using RNA-Seq should
52 profile expression at multiple resolutions to maximise the ability to detect causal eQTLs and
53 candidate-genes.

54

55 **Introduction**

56 The autoimmune diseases are a family of heritable, often debilitating, complex disorders whereby
57 immune system dysfunction leads to loss of tolerance to self-antigens and chronic inflammation [1].
58 Genome-wide association studies (GWAS) have now detected hundreds of susceptibility loci
59 contributing to risk of autoimmunity [2] yet their biological interpretation still remains challenging
60 [3]. Mapping single nucleotide polymorphisms (SNPs) that influence gene expression (eQTLs) can
61 provide crucial insight into the potential candidate genes and etiological pathways connected to
62 discrete disease phenotypes [4]. For example, such analyses have implicated dysregulation of
63 autophagy in Crohn's Disease [5], the pathogenic role of CD4⁺ effector memory T-cells in
64 Rheumatoid Arthritis [6], and an overrepresentation of transcription factors in Systemic Lupus
65 Erythematosus [7].

66

67 Expression profiling in appropriate cell types and physiological conditions is necessary to capture the
68 pathologically relevant regulatory changes driving disease risk [8]. Lack of such expression data is
69 thought to explain the observed disparity of shared genetic architecture between disease association
70 and gene expression at certain autoimmune loci [9]. A much overlooked cause of this disconnect
71 however, is not only the use of microarrays to profile gene expression, but also the resolution to
72 which expression is quantified using RNA-Sequencing (RNA-Seq) [10]. Expression estimates of
73 whole-genes, individual isoforms and exons, splice-junctions, and introns are obtainable with RNA-
74 Seq [11–18]. The SNPs that affect these discrete units of expression vary strikingly in their proximity
75 to the target gene, localisation to specific epigenetic marks, and effect on translated isoforms [18]. For
76 example, in over 57% of genes with both an eQTL influencing overall gene expression and a
77 transcript ratio QTL (trQTL) affecting the ratio of each transcript to the gene total, the causal variants
78 for each effect are independent and reside in distinct regulatory elements of the genome [18].

79

80 RNA-Seq based eQTL investigations that solely rely on whole-gene expression estimates are likely to
81 mask the allelic effects on independent exons and alternatively-spliced isoforms [16–19]. This is in

82 part due to subtle isoform switches and expression variation in exons that cannot be captured at gene-
83 level [20]. Recent evidence also suggests that exon-level based strategies are more sensitive than
84 conventional gene-level approaches, and allow for detection of moderate but systematic changes in
85 gene expression that are not necessarily derived from alternative-splicing events [15,21]. Furthermore,
86 gene-level summary counts can be biased in the direction of extreme exon outliers [21]. Use of
87 isoform-, exon-, and junction-level quantification in eQTL analysis also support the potential to not
88 only point to the candidate genes involved, but also the specific transcripts or functional domains
89 affected [10,18]. This of course facilitates the design of targeted functional studies and better
90 illuminates the causative relationship between regulatory genetic variation and disease. Lastly, though
91 intron-level quantification is not often used in conventional eQTL analysis, it can still provide
92 valuable insight into the role of unannotated exons in reference gene annotations, retained introns, and
93 even intronic enhancers [22,23].

94

95 Low-resolution expression profiling with RNA-Seq will impede the subsequent identification of
96 causal eQTLs when applying genetic and epigenetic fine-mapping approaches [24]. In this
97 investigation, we aim to increase our knowledge of the regulatory mechanisms and candidate genes of
98 human autoimmune disease through integration of GWAS and RNA-Seq expression data profiled at
99 gene-, isoform-, exon-, junction-, and intron-level in lymphoblastoid cell lines (LCLs). Our findings
100 are provided as a web resource to interrogate the functional effects of autoimmune associated SNPs
101 (www.insidegen.com), and will serve as the basis for targeted follow-up investigations.

102

103 **Results**

104

105 **Detection of *cis*-eQTLs and candidate-genes of autoimmune disease using RNA-Seq**

106 Using association data from twenty human autoimmune diseases, we performed integrative *cis*-eQTL
107 analysis in lymphoblastoid cell lines (LCLs) with RNA-Seq expression data profiled at five
108 resolutions: gene-, isoform-, exon-, junction-, and intron-level. We tested for a shared causal variant
109 between disease and expression at each association. The 846 autoimmune-associated SNPs taken
110 forward for analysis are documented in S1 Table and an overview of the analysis pipeline to detect
111 candidate-causal *cis*-eQTLs and eGenes is depicted in Fig 1. Expression targets at each level of RNA-
112 Seq quantification were interrogated in *cis* (± 1 Mb) to the 846 SNPs; comprising a total of 7,969
113 genes, 28,220 isoforms, 54,043 exons, 49,909 junctions, and 35,662 introns (Fig 2A).

114

115 We found that *cis*-eQTL association analysis using exon-, junction-, and intron-level quantification
116 yielded the greatest frequency of significant ($q < 0.05$) *cis*-eQTLs and eGenes (Fig 2B). These
117 findings persisted after testing whether each statistically significant *cis*-eQTL showed strong evidence
118 for colocalisation with the genetic variant underlying the autoimmune disease association ($q < 0.05$
119 and $RTC \geq 0.95$). For clarity, we define such eQTLs as candidate causal *cis*-eQTLs and we define
120 their targets as eGenes (Fig 2C). Exon-level analysis detected the most candidate-causal *cis*-eQTLs
121 (235) and eGenes (233) out of all quantification types, followed by junction- and intron-level
122 quantification. Isoform- and gene-level analysis were thoroughly outperformed, with the latter
123 detecting only 70 candidate-causal *cis*-eQTLs and 65 eGenes. In fact, we observed gene-level
124 quantification presented the greatest dropout of significant *cis*-eQTLs that were candidate causal (Fig
125 2D). Only 23.8% of significant *cis*-eQTLs were candidate-causal at gene-level compared to 49.8% at
126 exon-level; suggesting that in the autoimmune susceptibility loci tested more strongly associated *cis*-
127 eQTLs are captured by the exon-level analysis and they are distinct from gene-level *cis*-eQTLs. Gene-
128 level analysis under estimated candidate-causal eGenes. Our findings, highlighting the need to profile
129 gene-expression at multiple resolutions, are summarised in Fig 2E.

130

131 **Profiling at all resolutions is necessary to capture the full array associated *cis*-eQTLs**

132 We pruned the 846 autoimmune associated SNPs using an r^2 cut-off of 0.8 and 100kb limit to create a
133 subset of 584 independent susceptibility loci. By combining all five resolutions of RNA-Seq, we
134 found 267 loci (45.7%) presented a shared genetic effect between disease association and gene
135 expression (Fig 3A). Strikingly, only 9.3% of associated loci shared an underlying causal variant at
136 gene-level, in contrast to the 29.1% classified at exon-level. We mapped the candidate-causal *cis*-
137 eQTLs detected by RNA-Seq back to the diseases to which they are associated (Fig 3B). On average,
138 47% of associated SNPs per disease were classified as candidate-causal *cis*-eQTLs using all five
139 RNA-Seq quantification types. Interestingly, we observed the diseases that fell most below this
140 average comprised autoimmune disorders related to the gut: celiac disease (29%), inflammatory
141 bowel disease (36%), ulcerative colitis (39%), and Crohn's disease (41%), as well as Type 1 Diabetes
142 (37%). These observations are possibly a result of the cellular expression specificity of associated
143 genes in colonic and pancreatic tissue. This conclusion is supported by the above-average frequency
144 of candidate-causal *cis*-eQTLs detected in Systemic Lupus Erythematosus (50%) and Rheumatoid
145 Arthritis (54%); diseases in which the pathogenic role of B-lymphocytes is well documented [33,34].
146 We further broke down our results per disease by RNA-Seq quantification type (Fig 3C) and in almost
147 all cases, the greatest frequency of candidate-causal *cis*-eQTLs and eGenes were captured by exon-
148 and junction-level analyses.

149

150 By separating candidate-causal *cis*-eQTL associations out by quantification type, we found over half
151 were detected by either exon- or junction-level, and considerable overlap of *cis*-eQTL associations
152 existed between both types (Fig 3D). The greatest correlation of effect sizes (r^2 : 0.88) of candidate-
153 causal *cis*-eQTLs between exon- and junction-level (S1 Fig). Strong correlation also existed between
154 the effect sizes of gene- and isoform-level candidate-causal *cis*-eQTLs as expected (r^2 : 0.83); yet
155 gene-level analysis detected only 19% of all candidate-causal associations. Gene- and isoform-level
156 analysis did however capture six and eighteen candidate-causal *cis*-eQTLs unique to their

157 quantification type respectively. Thus, our data suggest that although exon- and junction-level, and to
158 a lesser extent intron-level analysis, capture the majority of candidate-causal *cis*-eQTL associations, it
159 is necessary to profile gene-expression at all quantification types to avoid misinterpretation of the
160 functional impact of disease associated SNPs.

161

162 **Web resource for functional interpretation of association studies of autoimmune disease**

163 We provide our data as a web resource (www.insidegen.com) for researchers to lookup candidate-
164 causal *cis*-eQTLs and eGenes of autoimmune diseases detected across the five RNA-Seq
165 quantification types. Data are sub-settable and exportable by SNP ID, gene, RNA-Seq resolution,
166 genomic position, and association to specific autoimmune diseases. Full data are also made available
167 in S2 Table.

168

169 **Functional dissection of Ankylosing Spondylitis genetic associations using RNA-Seq**

170 We decided to apply the results of our integrative *cis*-eQTL analysis to functionally dissect the
171 genetic associations of ankylosing spondylitis (AS). By doing so, we highlight the necessity of
172 profiling at all resolutions of RNA-Seq to shed light on novel regulatory variants, candidate genes,
173 and molecular pathways involved in pathogenesis. AS is a heritable inflammatory arthritis with a
174 largely unexplained genetic contribution outside of the *HLA-B*27* allele (> 30 risk loci) [35,36]. Only
175 a handful of loci show causative relationships with changes in gene expression [35,36]. Candidate-
176 genes implicated by association studies however suggest discrete immunological processes such as
177 antigen presentation, lymphocyte differentiation and activation, and regulation of the TNF/NF- κ B
178 signalling pathways are involved, and of note, strong genetic overlap exists with psoriasis, psoriatic
179 arthritis, and inflammatory bowel disease; indicating the pathogenesis of these diseases are tightly
180 connected [37].

181

182 Of the 31 AS associated SNPs taken forward for functional interrogation, 14 were classified as
183 candidate-causal *cis*-eQTLs (Fig 4A; full results found at www.insidegen.com for all diseases). We

184 replicated the association of risk allele rs4129267 [C] with expression reduction of *IL6R* by junction-
185 level analysis ($\beta = -0.36$; $P = 1.14 \times 10^{-06}$) [35]. Interestingly, we found the expression of
186 neighbouring gene, *RUSC1*, is also influenced by candidate-causal *cis*-eQTL rs4129267 where the
187 risk allele was also reduced expression of *RUSC1* at exon-level ($\beta = -0.24$; $P = 1.59 \times 10^{-03}$). *RUSC1*
188 is able to polyubiquitinate *IKBKG*, a key regulator of NF- κ B [38].

189 The effect of independently associated variants within the 5q15 locus on the expression of
190 aminopeptidase genes *ERAPI* and *ERAP2* was also replicated (S3 Fig) [36]. This includes the
191 association of risk allele rs30187 [T] with increased expression of *ERAPI* ($\beta = -1.09$; $P = 1.60 \times 10^{-71}$),
192 and the striking effect of protective allele rs2910686 [T] on the near-complete loss of *ERAP2* ($\beta =$
193 -1.37 ; $P = 1.95 \times 10^{-175}$). Again however, additional genes at this locus with no previous association to
194 expression changes with regards to AS risk alleles were detected. *LNPEP* also belongs to the
195 endoplasmic reticulum aminopeptidase family and has been shown to regulate the NF- κ B pathway
196 and antigen presentation via peptide trimming [39]. Interestingly, a missense variation in this gene is
197 linked to psoriasis and is down-regulated in psoriatic lesions relative to healthy skin [40]. We found at
198 junction-level, AS risk allele rs2910686 [C] also contributes to expression reduction of *LNPEP* ($\beta =$
199 0.41 ; $P = 3.09 \times 10^{-08}$). Similarly, the risk allele rs30187 [T] correlated strongly with decreased
200 expression of *CAST* ($\beta = -0.46$; $P = 2.47 \times 10^{-10}$); encoding calpastatin, a calcium-dependent cysteine
201 protease inhibitor. Cysteine protease activity positively correlates with the severity of arthritic lesions
202 and degree of inflammation [41]. Our data support the notion of multiple functional effects at this
203 locus and suggests novel pathological mechanisms including decreased expression of the inhibitor
204 *CAST* leading to increased cysteine protease activity.

205 Other AS susceptibility loci contributing to expression modulation of multiple genes include
206 rs9901869 for *TBKBPI* and *ITGB3*, and rs75301646 for *NFATC2IP* and *TUFM*. Candidate genes at
207 the rs9901869 locus are yet to be functionally characterised [35]. Our data suggest the risk allele
208 rs9901869 [A] increases the expression of both *TBKBPI*, which plays an active role in the NF- κ B and
209 IFN- α signalling pathways ($\beta = 0.57$; $P = 2.47 \times 10^{-17}$) [42], and *ITGB3* - involved in the intestinal
210 immune pathway for IgA production ($\beta = 0.35$; $P = 1.53 \times 10^{-6}$) [43]. Similarly, we found the risk
211 allele rs9901869 [A] increases expression of both novel candidate-causal eGenes *NFATC2IP* ($\beta =$

212 0.25; $P = 6.18 \times 10^{-4}$) and *TUFM* ($\beta = 0.60$; $P = 2.82 \times 10^{-17}$). *TUFM* has been reported as the
213 causative gene at this locus for early onset inflammatory bowel disease [44], whereas *NFATC2IP*
214 (Nuclear Factor of Activated T-cells 2 Interacting Protein) has clear immunological roles in the
215 induction of IL-4 production and regulation of the TNF receptor family of proteins [45].
216 Our analysis has shed new light on the molecular genetics of AS and can be used in similar manner
217 for the functional dissection of the remaining 19 autoimmune diseases (www.insidegen.com).

218

219 **Functional genomic support for candidate-causal *cis*-eQTLs**

220 The resolution of RNA-Seq can be leveraged to map candidate-genes and isolate specific exons and
221 junctions perturbed by disease-associated variants. Functional genomic data can then be used to
222 support potential causal associations to deduce molecular mechanisms and epigenetically prioritize
223 causal variants.

224 The remaining AS associated variant rs1128905 is a candidate-causal *cis*-eQTL for both *CARD9* and
225 *SNAPC4* (Fig 4A). The candidate-gene at this AS locus is thought to be *CARD9* [36]. Our results also
226 draw attention to *SNAPC4* (Small Nuclear RNA-activating Complex Polypeptide 4). Using exon-level
227 RNA-Seq, the risk allele rs1128905 [C] decreased the expression of exons 18 ($\beta = -0.37$; $P = 3.65 \times$
228 10^{-07}) and 19 ($\beta = -0.27$; $P = 4.80 \times 10^{-04}$) of the canonical transcript of *SNAPC4* (Fig 4B).
229 Accordingly, the exon 18-19 boundary was also significantly decreased, captured by junction-level
230 quantification ($\beta = -0.26$; $P = 2.74 \times 10^{-04}$). As rs1128905 lies over 39kb away from the transcription
231 start site of *SNAPC4*, we used existing promoter capture Hi-C data in lymphoblastoid cell lines to
232 assess whether rs1128905 and associated SNPs may act distally upon *SNAPC4* to influence its
233 expression [32]. We found the bait region encompassing rs1128905 interacts with five targets with
234 great confidence CHiCAGO score > 12 (Fig 4C) [46]. Four of these are located within the *SNAPC4*
235 gene itself. Adding further evidence from histone marks from lymphoblastoid cell lines from the
236 RoadMap Epigenomics Project [31], we found two SNPs in near-perfect LD with rs1128905 ($r^2 >$
237 0.95; rs10870201 and rs10870202) were localised to the peaks of H3K4me3, H3K27ac, and H3K9ac
238 marks, and the region encompassed is predicted to be an active enhancer (S4 Fig). Our data therefore

239 suggest that associated SNPs rs10870201 and rs10870202 may perturb the enhancer-promoter
240 interaction with *SNAPC4* affecting expression. In fact, rs10870201 was the best *cis*-eQTL in the 1Mb
241 region for exons 18 and 19 of *SNAPC4*. Interestingly, although no autoimmune phenotype has been
242 documented with *SNAPC4*, an uncorrelated SNP rs10781500 (r^2 with rs1128905 < 0.5), associated
243 with Crohn's Disease, inflammatory bowel disease, and ulcerative colitis, has also been classified as a
244 candidate-causal *cis*-eQTL for *SNAPC4* but not *CARD9* in *ex vivo* human B-lymphocytes (the risk
245 allele is also correlated with reduced expression of *SNAPC4*) [47]. This effect holds true in our
246 analysis - rs10781500 is an eQTL for *SNAPC4* but not *CARD9*.

247 Our data point to candidate genes and molecular mechanisms but further functional characterization is
248 of course necessary to determine the true causative gene(s) at this locus.

249

250 **Detection of autoimmune associated *trans*-eQTLs using RNA-Seq**

251 We extended our RNA-Seq based eQTL investigation to include expression targets > 5Mb away from
252 each of the 846 lead autoimmune GWAS variants (S3 Table). Though we were relatively
253 underpowered for a *trans*-eQTL analysis, we were able to detect 26 *trans*-eQTLs at isoform-level,
254 eight at exon-level, six at gene-level, three at junction-level (Fig 5A). Many of the *trans*-eQTLs
255 detected were only associated with one eGene, and no *trans*-eQTLs were detected at intron-level. With
256 exon-level quantification however, we were able to identify an interesting effect of *trans*-eQTL
257 rs7726414 - associated with Systemic Lupus Erythematosus (SLE) [7]. We found rs7726414, was a
258 *trans*-eQTL for eight eGenes (Fig 5B). These comprise *SIPA1L2*, *PDPK1*, *IVNSIABP*, *HES2*, *JAZF1*,
259 *ULK4*, *RP11-51F16.8*, and *PPM1M*. We found the risk allele rs7726414 [T] was associated with
260 increased expression of each of these eight genes (Fig 5C). We highlight, *PDPK1*, a key regulator of
261 *IRF4* and inducer of apoptosis [48], and *JAZF1* which is genetically associated with many
262 autoimmune diseases including SLE itself [7]. The serine/threonine-protein kinase *ULK4* is also of
263 interest as its family member, *ULK3*, is also an SLE susceptibility gene [7]. Though we did not
264 classify rs7726414 as a candidate-causal *cis*-eQTL in our dataset, it has been documented as
265 candidate-causal in SLE using a larger eQTL cohort profiled in lymphoblastoid cell lines for eGenes

266 *TCF7* (Transcription Factor 7, T-Cell Specific, HMG-Box) and the ubiquitin ligase complex *SKP1*

267 [10].

268

269 Discussion

270 Elucidation of the functional consequence of non-coding genetic variation in human disease is a major
271 objective of medical genomics [49]. Integrative studies that map disease-associated eQTLs in relevant
272 cell types and physiological conditions are proving essential in progression towards this goal through
273 identification of causal SNPs, candidate-genes, and illumination of molecular mechanisms [50]. In
274 autoimmune disease, where there is considerable overlap of immunopathology, integrative eQTL
275 investigations have been able to connect discrete aetiological pathways, cell types, and epigenetic
276 modifications, to particular clinical manifestations [2,50–52]. Emerging evidence however suggests
277 that only a minority (~25%) of autoimmune associated SNPs share casual variants with *cis*-eQTLs in
278 primary immune cell-types [9].

279

280 Genetic variation can influence expression at every stage of the gene regulatory cascade - from
281 chromatin dynamics, to RNA folding, stability, and splicing, and protein translation [53]. As RNA-
282 Seq becomes the convention for genome-wide transcriptomics, be it for differential expression or
283 eQTL analysis, it is essential to maximise the ability to resolve and quantify discrete transcriptomic
284 features. It is now well documented that SNPs affecting these units of expression vary strikingly in
285 their genomic location and localisation to specific epigenetic marks [18]. The reasoning for our
286 investigation therefore was to delineate the limits of microarray and RNA-Seq based eQTL cohorts in
287 the functional annotation of autoimmune disease association signals. To map autoimmune disease
288 associated *cis*-eQTLs, we interrogated RNA-Seq expression data profiled at gene-, isoform, exon-,
289 junction-, and intron-level, and tested for a shared genetic effect at each significant association. We
290 found exon- and junction-level quantification led to the greatest frequency of candidate-causal *cis*-
291 eQTL and eGenes, and thoroughly outperformed gene-level analysis (Fig 2C). We argue however that
292 it is necessary to profile expression at all possible resolutions to diminish the likelihood of
293 overlooking potentially causal *cis*-eQTLs (Fig 3D). In fact, by combining our results across all
294 resolutions, we found 45% of autoimmune loci were candidate-causal *cis*-eQTLs for at least one

295 eGene. Our findings can be used as a resource to lookup causal eQTLs and candidate genes of
296 autoimmune disease (www.insidegen.com).

297

298 Gene-level expression estimates can generally be obtained in two ways – union-exon based
299 approaches [14,17] and transcript-based approaches [11,12]. In the former, all overlapping exons of
300 the same gene are merged into union exons, and intersecting exon and junction reads (including split-
301 reads) are counted to these pseudo-gene boundaries. Using this counting-based approach, it is also
302 possible to quantify meta-exons and junctions easily and with high confidence by preparing the
303 reference annotation appropriately [13,15,54]. Introns can be quantified in a similar manner by
304 inverting the reference annotation between exons and introns [18]. Conversely, transcript-based
305 approaches make use of statistical models and expectation maximization algorithms to distribute reads
306 among gene isoforms - resulting in isoform expression estimates [11,12]. These estimates can then be
307 summed to obtain the entire expression estimate of the gene. Greater biological insight is gained from
308 isoform-level analysis; however, disambiguation of specific transcripts is not trivial due to substantial
309 sequence commonality of exons and junctions. In fact, we found only 15% of autoimmune loci shared
310 a causal variant at transcript-level (Fig 3A). The different approaches used to estimate expression can
311 also lead to significant differences in the reported counts. Union-based approaches, whilst
312 computationally less expensive, can underestimate expression levels relative to transcript-based, and
313 this difference becomes more pronounced when the number of isoforms of a gene increases, and when
314 expression is primarily derived from shorter isoforms [20]. The GEUVADIS study implemented a
315 transcript-based approach to obtain whole-gene expression estimates. A gold standard of eQTL
316 mapping using RNA-Seq is essential therefore for comparative analysis across datasets.

317

318 Our findings support recent evidence that suggests exon-level based strategies are more sensitive and
319 specific than conventional gene-level approaches [21]. Subtle isoform variation and expression of less
320 abundant isoforms are likely to be masked by gene-level analysis. Exon-level allows for detection of
321 moderate but systematic changes in gene expression that are not captured at gene-level, and also,
322 gene-level summary counts can be shifted in the direction of extreme exon outliers [21]. It is therefore

323 important to note that a positive exon-level eQTL association does not necessarily mean a differential
324 exon-usage or splicing mechanism is involved; rather a systematic expression effect across the whole
325 gene may exist that is only captured by the increased sensitivity. Additionally, by combining exon-
326 level with other RNA-Seq quantification types, inferences can be made on the particular isoforms and
327 functional domains affected by the eQTL which can later aid biological interpretation and targeted
328 follow-up investigations [10].

329

330 We found intron-level quantification also generated more candidate-causal *cis*-eQTLs than gene-
331 level. As the library was synthesised from poly-A selection, these associations are unlikely due to
332 differences in pre-mRNA abundance. Rather, they are likely derived from either true retained introns
333 in the mature RNA or from coding exons that are not documented in the reference annotation used.
334 We observed multiple instances where a candidate-causal *cis*-eQTL at intron-level was detected, yet a
335 previous investigation had detected an exonic effect using a different reference annotation. For
336 example, an intronic-effect was detected for SLE candidate eGenes *IKZF2* and *WDFY4* in this
337 analysis (which used the GENCODE v12 basic reference annotation). Using the comprehensive
338 reference annotation of GENCODE v12, we found these effects were in fact driven by transcribed
339 exons located within the intronic region of the basic annotation – and were validated *in vitro* by qPCR
340 [10]. The choice of reference annotation therefore has a profound effect on expression estimates [55];
341 and so again, a gold standard is necessary prevent misinterpretation and increase consistency of eQTL
342 associations.

343

344 Lastly, we show how our findings can be leveraged to comprehensively dissect GWAS results of
345 autoimmune diseases. We found 14 of the 31 SNPs associated with Ankylosing Spondylitis (AS)
346 were candidate-causal *cis*-eQTLs for at least on eGene (Fig 4). The majority of these eQTLs
347 influenced the expression of multiple eGenes which had direct relevance to biological pathways
348 associated with autoimmunity. In fact, the majority of the candidate genes detected (for example
349 *RUSC1*, *TBKBPI*, *NFATC2IP*, *TNFRSF1A*, and *PDE4A*) support the involvement of TNF- α and NF-
350 κ B in the pathology of AS [35]. We finally show at the *CARD9-SNAPC4* locus, how existing

351 functional genomic data from chromatin interaction and epigenetic modification experiments can
352 strengthen evidence of the eQTL associations detected by RNA-Seq and allow for functional
353 prioritization of causal variants (Fig 4C). We also highlight the benefit of exon-level analysis to also
354 detect disease associated *trans*-eQTLs (Fig 5).

355

356 Taken together, we have provided a deeper mechanistic understanding of the genetic regulation of
357 gene expression in autoimmune disease by profiling the transcriptome at multiple resolutions using
358 RNA-Seq. Similar analyses in new and existing datasets using relevant cell types and context-specific
359 conditions will undoubtedly increase our understanding of how associated variants alter cell
360 physiology and ultimately contribute to disease risk.

361 **Materials and Methods**

362

363 **Autoimmune disease associated SNPs**

364 SNPs were taken from the ImmunoBase resource (www.immunobase.org). It comprises summary
365 case-control association statistics from twenty diseases: twelve originally targeted by the
366 ImmunoChip consortium (Ankylosing Spondylitis, Autoimmune Thyroid Disease, Celiac Disease,
367 Crohn's Disease, Juvenile Idiopathic Arthritis, Multiple Sclerosis, Primary Biliary Cirrhosis, Psoriasis,
368 Rheumatoid Arthritis, Systemic Lupus Erythematosus, Type 1 Diabetes, Ulcerative Colitis), and eight
369 others (Alopecia Areata, Inflammatory Bowel Disease, IgE and Allergic Sensitization, Narcolepsy,
370 Primary Sclerosing Cholangitis, Sjogren Syndrome, Systemic Scleroderma, Vitiligo). For eQTL
371 analysis, we took the lead SNPs for each disease - defined as a genome-wide significant SNP with the
372 lowest reported *P*-value (S1 Table). X-chromosome associations and SNPs with minor allele
373 frequency < 5% were omitted from analysis, leaving 846 SNPs. A total of 262 SNPs were pruned
374 using the '*--indep-pairwise*' function of PLINK 1.9 with a window size of 100kb and an r^2 threshold
375 of 0.8, leaving 584 independent loci.

376

377 **RNA-Seq gene expression data**

378 Normalised RNA-Seq expression data of 373 lymphoblastoid cell lines from four European sub-
379 populations (CEU, GBR, FIN, TSI) of the 1000Genomes Project (Geuvadis) [18] were obtained from
380 EBI ArrayExpress (E-GEUV-1; full methods can be found in <http://geuvadiswiki.crg.es/>). In
381 summary, transcripts, splice-junctions, and introns were quantified using Flux Capacitor against the
382 GENCODE v12 basic reference annotation [16]. Reads belonging to single transcripts were predicted
383 by deconvolution per observations of paired-reads mapping across all exonic segments of a locus.
384 Gene-level expression was calculated as the sum of all transcripts per gene. Annotated splice
385 junctions were quantified using split read information, counting the number of reads supporting a
386 given junction. Intronic regions that are not retained in any mature annotated transcript, and reported
387 mapped reads in different bins across the intron to distinguish reads stemming from retained introns

388 from those produced by not yet annotated exons. Meta-exons were quantified by merging all
389 overlapping exonic portions of a gene into non-redundant units and counting reads within these bins
390 [15]. Reads were excluded when the read pairs map to two different genes. Quantifications were
391 corrected for sequencing depth and gene length (RPKM). Only expression elements quantified in > 50
392 % of individuals were kept and Probabilistic Estimation of Expression Residuals (PEER) was used to
393 remove technical variation [25].

394

395 ***Cis* and *trans*-eQTL analysis**

396 An overview of the integration pipeline is depicted in Fig 1. Genotypes were obtained from EBI
397 ArrayExpress (E-GEUV-1). The 41 individuals genotyped on the Omni 2.5M SNP array were
398 previously imputed to the Phase 1 v3 release as described [18]. PCA of genotype data was performed
399 using the Bioconductor package SNPRelate (S2 Fig) [26]. Only bi-allelic SNPs with MAF > 0.05,
400 imputation call-rates ≥ 0.8 , and HWE $P < 1 \times 10^{-04}$ were used. All eQTL association testing was
401 performed with a linear-regression model in R. Normalized expression residuals (PEER factor
402 normalized RPKM) for each quantification type were transformed to standard normal and the first
403 three principle components used as covariates in the eQTL model as well as the binary imputation
404 status. *Cis* and *trans*-eQTL mapping was performed for genes within +/-1Mb of the lead SNP and for
405 genes > 5Mb from the lead SNP respectively. Adjustment for multiple testing of eQTL results per
406 quantification type (corrected total genes, isoforms, exons, junctions, and introns) was undertaken
407 using an FDR of 0.05 for *cis* and 0.01 for *trans* analysis (MHC associations were excluded in *trans*).

408

409 **Analysis of shared causal variant**

410 The Regulatory Trait Concordance (RTC) method was used to assess the likelihood of a shared causal
411 variant between the GWAS SNP and the *cis*-eQTL signal [27]. SNPs were firstly classified according
412 to their position in relation to recombination hotspots (based on genome-wide estimates of hotspot
413 intervals) [28]. For each significant *cis*-eQTL association, the residuals from the linear-regression of
414 the best *cis*-eQTL (lowest association *P*-value within the hotspot interval) was extracted against the

415 expression quantification for the expression unit in hand. Regression was then performed using all
416 SNPs within the defined hotspot interval against these residuals. The RTC score was then calculated
417 as $(N_{SNPs} - Rank_{GWAS\ SNP} / N_{SNPs})$. Where N_{SNPs} is the total number of SNPs in the recombination hotspot
418 interval, and $Rank_{GWAS\ SNP}$ is the rank of the GWAS SNP association P -value against all other SNPs in
419 the interval from the linear-association against the residuals of the best *cis*-eQTL. Disease associated
420 SNPs with statistically significant association to gene expression ($q < 0.05$) and an RTC score > 0.95
421 were classified as ‘candidate-causal eQTLs’. Genes whose expression is modulated by the eQTL were
422 defined as ‘candidate-causal eGenes’.

423

424 **Data visualisation and resources**

425 R version 3.3.1 was used to create heatmaps, box-plots (ggplot2), and circularized chromosome
426 diagrams (circlize). Genes were plotted in UCSC Genome Browser [29] and IGV [30]. Roadmap
427 epigenetic data were downloaded from the web resource [31], and chromatin interaction data were
428 taken from the CHiCP web resource [32].

429

430 **Acknowledgements**

431 We thank Dr David L Morris for helpful discussions throughout this work. We also thank Philip

432 Tombleson for his assistance with data uploading.

433 The GEUVADIS 1000 Genomes RNA-Seq data was downloaded from the EBI ArrayExpress Portal

434 (accession E-GEUV-1).

435

436 **References**

- 437 1. Fever FM. NIH Progress in Autoimmune Diseases Research. in National Institute of Health
438 Publication. 2005; 17–7576.
- 439 2. Parkes M, Cortes A, van Heel DA, Brown MA. Genetic insights into common pathways and
440 complex relationships among immune-mediated diseases. *Nat Rev Genet*. Nature Publishing
441 Group; 2013;14: 661–73. doi:10.1038/nrg3502
- 442 3. Hindorf LA, Sethupathy P, Junkins HA, Ramos EM, Mehta JP, Collins FS, et al. Potential
443 etiologic and functional implications of genome-wide association loci for human diseases and
444 traits. *Proc Natl Acad Sci U S A*. 2009;106: 9362–9367. doi:10.1073/pnas.0903103106
- 445 4. Westra H-J, Franke L. From genome to function by studying eQTLs. *Biochim Biophys Acta*.
446 Elsevier B.V.; 2014;1842: 1896–1902. doi:10.1016/j.bbadis.2014.04.024
- 447 5. Klionsky DJ. Crohn’s disease, autophagy, and the Paneth cell. *N Engl J Med*. 2009;360: 1785–
448 1786. doi:10.1056/NEJMcibr0810347
- 449 6. Hu X, Kim H, Raj T, Brennan PJ, Trynka G, Teslovich N, et al. Regulation of Gene
450 Expression in Autoimmune Disease Loci and the Genetic Basis of Proliferation in CD4+
451 Effector Memory T Cells. *PLoS Genet*. 2014;10. doi:10.1371/journal.pgen.1004404
- 452 7. Bentham J, Morris DL, Cunninghame Graham DS, Pinder CL, Tombleson P, Behrens TW, et
453 al. Genetic association analyses implicate aberrant regulation of innate and adaptive immunity
454 genes in the pathogenesis of systemic lupus erythematosus. *Nat Genet*. Nature Publishing
455 Group; 2015;47: 1457–1464. doi:10.1038/ng.3434
- 456 8. Fairfax BP, Knight JC. Genetics of gene expression in immunity to infection. *Curr Opin*
457 *Immunol*. Elsevier Ltd; 2014;30: 63–71. doi:10.1016/j.coi.2014.07.001
- 458 9. Chun S, Casparino A, Patsopoulos NA, Croteau-chonka DC, Raby BA, Jager PL De, et al.
459 Limited statistical evidence for shared genetic effects of eQTLs and autoimmune-disease-
460 associated loci in three major immune-cell types. *NatGenet*. 2017; doi:10.1038/ng.3795
- 461 10. Odhams CA, Cortini A, Chen L, Roberts AL, Viñuela A, Buil A, et al. Mapping eQTLs with
462 RNA-Seq Reveals Novel Susceptibility Genes, Non-Coding RNAs, and Alternative-Splicing

- 463 Events in Systemic Lupus Erythematosus. *Hum Mol Genet.* 2017;26: 1003–1017.
464 doi:10.1093/hmg/ddw417
- 465 11. Trapnell C, Roberts A, Goff L, Pertea G, Kim D, Kelley DR, et al. Differential gene and
466 transcript expression analysis of RNA-seq experiments with TopHat and Cufflinks. *Nat*
467 *Protoc.* 2012;7: 562–78. doi:10.1038/nprot.2012.016
- 468 12. Li B, Dewey CN. RSEM: accurate transcript quantification from RNA-Seq data with or
469 without a reference genome. *BMC Bioinformatics.* 2011;12: 323. doi:10.1186/1471-2105-12-
470 323
- 471 13. Schuierer S, Roma G. The exon quantification pipeline (EQP): a comprehensive approach to
472 the quantification of gene, exon and junction expression from RNA-seq data. *Nucleic Acids*
473 *Res.* 2016; gkw538. doi:10.1093/nar/gkw538
- 474 14. Anders S, Pyl PT, Huber W. HTSeq-A Python framework to work with high-throughput
475 sequencing data. *Bioinformatics.* 2015;31: 166–169. doi:10.1093/bioinformatics/btu638
- 476 15. Anders S, Reyes A, Huber W. Detecting differential usage of exons from RNA-seq-
477 npre20126837-2.pdf. *Genome Res.* 2012;12: 1088–9051. doi:10.1101/gr.133744.111
- 478 16. Montgomery SB, Sammeth M, Gutierrez-Arcelus M, Lach RP, Ingle C, Nisbett J, et al.
479 Transcriptome genetics using second generation sequencing in a Caucasian population.
480 *Nature.* Nature Publishing Group; 2010;464: 773–777. doi:10.1038/nature08903
- 481 17. Liao Y, Smyth GK, Shi W. FeatureCounts: An efficient general purpose program for assigning
482 sequence reads to genomic features. *Bioinformatics.* 2014;30: 923–930.
483 doi:10.1093/bioinformatics/btt656
- 484 18. Lappalainen T, Sammeth M, Friedländer MR, 't Hoen P a C, Monlong J, Rivas M a, et al.
485 Transcriptome and genome sequencing uncovers functional variation in humans. *Nature.*
486 2013;501: 506–11. doi:10.1038/nature12531
- 487 19. Battle A, Mostafavi S, Zhu X, Potash JB, Weissman MM, McCormick C, et al. Characterizing
488 the genetic basis of transcriptome diversity through RNA-sequencing of 922 individuals.
489 *Genome Res.* 2014;24: 14–24. doi:10.1101/gr.155192.113
- 490 20. Zhao S, Xi L, Zhang B. Union exon based approach for RNA-seq gene quantification: To be

- 491 or not to be? PLoS One. 2015;10: e0141910. doi:10.1371/journal.pone.0141910
- 492 21. Laiho A, Elo LL. A note on an exon-based strategy to identify differentially expressed genes
493 in RNA-seq experiments. PLoS One. 2014;9: 1–12. doi:10.1371/journal.pone.0115964
- 494 22. Gaidatzis D, Burger L, Florescu M, Stadler MB. Analysis of intronic and exonic reads in
495 RNA-seq data characterizes transcriptional and post-transcriptional regulation. Nat Biotech.
496 Nature Publishing Group; 2015;33: 722–729. doi:10.1038/nbt.3269
- 497 23. Mortazavi A, Williams BA, McCue K, Schaeffer L, Wold B. Mapping and quantifying
498 mammalian transcriptomes by RNA-Seq. Nat Methods. 2008;5: 621–628.
499 doi:10.1038/nmeth.1226
- 500 24. Trynka G, Westra HJ, Slowikowski K, Hu X, Xu H, Stranger BE, et al. Disentangling the
501 Effects of Colocalizing Genomic Annotations to Functionally Prioritize Non-coding Variants
502 within Complex-Trait Loci. Am J Hum Genet. The Authors; 2015;97: 139–152.
503 doi:10.1016/j.ajhg.2015.05.016
- 504 25. Stegle O, Parts L, Durbin R, Winn J. A bayesian framework to account for complex non-
505 genetic factors in gene expression levels greatly increases power in eQTL studies. PLoS
506 Comput Biol. 2010;6: 1–11. doi:10.1371/journal.pcbi.1000770
- 507 26. Zheng X, Levine D, Shen J, Gogarten SM, Laurie C, Weir BS. A high-performance computing
508 toolset for relatedness and principal component analysis of SNP data. Bioinformatics. 2012;28:
509 3326–3328. doi:10.1093/bioinformatics/bts606
- 510 27. Nica AC, Montgomery SB, Dimas AS, Stranger BE, Beazley C, Barroso I, et al. Candidate
511 causal regulatory effects by integration of expression QTLs with complex trait genetic
512 associations. PLoS Genet. 2010;6: e1000895. doi:10.1371/journal.pgen.1000895
- 513 28. McVean GA. The fine-scale structure of recombination rate variation in the human genome.
514 Science (80-). 2004;304: 581. Available: <http://dx.doi.org/10.1126/science.1092500>
- 515 29. Kent WJ, Sugnet CW, Furey TS, Roskin KM. The Human Genome Browser at UCSC W. J
516 Med Chem. 2002;19: 1228–31. doi:10.1101/gr.229102.
- 517 30. Thorvaldsdottir H, Robinson JT, Mesirov JP. Integrative Genomics Viewer (IGV): high-
518 performance genomics data visualization and exploration. Brief Bioinform. 2013;14: 178–192.

- 519 doi:10.1093/bib/bbs017
- 520 31. Kundaje A, Meuleman W, Ernst J, Bilenky M, Yen A, Heravi-Moussavi A, et al. Integrative
521 analysis of 111 reference human epigenomes. *Nature*. 2015;518: 317–330.
522 doi:10.1038/nature14248
- 523 32. Schofield EC, Carver T, Achuthan P, Freire-Pritchett P, Spivakov M, Todd JA, et al. CHiCP:
524 A web-based tool for the integrative and interactive visualization of promoter capture Hi-C
525 datasets. *Bioinformatics*. 2016;32: 2511–2513. doi:10.1093/bioinformatics/btw173
- 526 33. Marston B. B cells in the pathogenesis and treatment of rheumatoid arthritis. *Curr Opin*
527 *Rheumatol*. 2011;22: 307–315. doi:10.1097/BOR.0b013e3283369cb8.B
- 528 34. Dörner T, Giesecke C, Lipsky PE. Mechanisms of B cell autoimmunity in SLE. *Arthritis Res*
529 *Ther*. 2011;13: 243. doi:10.1186/ar3433
- 530 35. Cortes A, Hadler J, Pointon JP, Robinson PC, Karaderi T, Leo P, et al. Identification of
531 multiple risk variants for ankylosing spondylitis through high-density genotyping of immune-
532 related loci. *Nat Genet*. 2013;45: 730–8. doi:10.1038/ng.2667
- 533 36. Reveille JD, Sims A-M, Danoy P, Evans DM, Leo P, Pointon JJ, et al. Genome-wide
534 association study of ankylosing spondylitis identifies non-MHC susceptibility loci. *Nat Genet*.
535 2010;42: 123–7. doi:10.1038/ng.513
- 536 37. O’Rielly DD, Uddin M, Rahman P. Ankylosing spondylitis: beyond genome-wide association
537 studies. *Curr Opin Rheumatol*. 2016;28: 337–45. doi:10.1097/BOR.0000000000000297
- 538 38. Napolitano G, Mirra S, Monfregola J, Lavorgna A, Leonardi A, Ursini MV. NESCA: A new
539 NEMO/IKKgamma and TRAF6 interacting protein. *J Cell Physiol*. 2009;220: 410–417.
540 doi:10.1002/jcp.21782
- 541 39. Nair RP, Duffin KC, Helms C, Ding J, Stuart PE, Goldgar D, et al. Genome-wide scan reveals
542 association of psoriasis with IL-23 and NF-kappaB pathways. *Nat Genet*. 2009;41: 199–204.
543 doi:10.1038/ng.311
- 544 40. Cheng H, Li Y, Zuo X-B, Tang H-Y, Tang X-F, Gao J-P, et al. Identification of a Missense
545 Variant in LNPEP that Confers Psoriasis Risk. *J Invest Dermatol*. Elsevier Masson SAS;
546 2013;134: 1–22. doi:10.1038/jid.2013.317

- 547 41. Biroc SL, Gay S, Hummel K, Magill C, Palmer JT, Spenc DR, et al. Cysteine protease activity
548 is up-regulated in inflamed ankle joints of rats with adjuvant-induced arthritis and decreases
549 with in vivo administration of a vinyl sulfone cysteine protease inhibitor. *Arthritis Rheum.*
550 2001;44: 703–711. doi:10.1002/1529-0131(200103)44:3<703::AID-ANR120>3.0.CO;2-2
- 551 42. Goncalves A, Bürckstümmer T, Dixit E, Scheicher R, Górna MW, Karayel E, et al. Functional
552 dissection of the TBK1 molecular network. *PLoS One.* 2011;6.
553 doi:10.1371/journal.pone.0023971
- 554 43. Reynier F, Pachot A, Paye M, Xu Q, Turrel-Davin F, Petit F, et al. Specific gene expression
555 signature associated with development of autoimmune type-I diabetes using whole-blood
556 microarray analysis. *Genes Immun.* 2010;11: 269–278. doi:10.1038/gene.2009.112
- 557 44. Narahara M, Higasa K, Nakamura S, Tabara Y, Kawaguchi T, Ishii M, et al. Large-scale East-
558 Asian eQTL mapping reveals novel candidate genes for LD mapping and the genomic
559 landscape of transcriptional effects of sequence variants. *PLoS One.* 2014;9.
560 doi:10.1371/journal.pone.0100924
- 561 45. Hashiguchi K, Ozaki M, Kuraoka I, Saitoh H. Establishment of a human cell line stably
562 overexpressing mouse Nip45 and characterization of Nip45 subcellular localization. *Biochem*
563 *Biophys Res Commun.* Elsevier Inc.; 2013;430: 72–77. doi:10.1016/j.bbrc.2012.11.020
- 564 46. Cairns J, Freire-Pritchett P, Wingett SW, Dimond A, Plagnol V, Zerbino D, et al. CHiCAGO:
565 Robust Detection of DNA Looping Interactions in Capture Hi-C data. *Genome Biol. Genome*
566 *Biology;* 2016;17: 28068. doi:10.1101/028068
- 567 47. Fairfax BP, Makino S, Radhakrishnan J, Plant K, Leslie S, Dilthey A, et al. Genetics of gene
568 expression in primary immune cells identifies cell type–specific master regulators and roles of
569 HLA alleles. *Nat Genet.* Nature Publishing Group; 2012;44: 502–510. doi:10.1038/ng.2205
- 570 48. Chinen Y, Kuroda J, Shimura Y, Nagoshi H, Kiyota M, Yamamoto-Sugitani M, et al.
571 Phosphoinositide protein kinase PDPK1 is a crucial cell signaling mediator in multiple
572 myeloma. *Cancer Res.* 2014;74: 7418–7429. doi:10.1158/0008-5472.CAN-14-1420
- 573 49. Lappalainen T. Functional genomics bridges the gap between quantitative genetics and
574 molecular biology. *Genome Res.* 2015;25: 1427–1431. doi:10.1101/gr.190983.115.

- 575 50. Albert FW, Kruglyak L. The role of regulatory variation in complex traits and disease. *Nat*
576 *Rev Genet.* Nature Publishing Group; 2015;16: 197–212. doi:10.1038/nrg3891
- 577 51. Farh KK, Marson A, Zhu J, Kleinewietfeld M, Housley WJ, Beik S, et al. Genetic and
578 epigenetic fine mapping of causal autoimmune disease variants. *Nature.* Nature Publishing
579 Group; 2015;518: 337–343. doi:10.1038/nature13835
- 580 52. Trynka G, Sandor C, Han B, Xu H, Stranger BE, Liu XS, et al. Chromatin marks identify
581 critical cell types for fine mapping complex trait variants. *Nat Genet.* Nature Publishing
582 Group; 2013;45: 124–30. doi:10.1038/ng.2504
- 583 53. Li YI, Geijn B Van De, Raj A, Knowles DA, Petti AA, Golan D, et al. RNA splicing is a
584 primary link between genetic variation and disease. *Science.* 2016;352: 600–4.
585 doi:10.1126/science.aad9417
- 586 54. Ongen H, Dermitzakis ET. Alternative Splicing QTLs in European and African Populations.
587 *Am J Hum Genet.* The Authors; 2015;97: 567–575. doi:10.1016/j.ajhg.2015.09.004
- 588 55. Zhao S, Zhang B. A comprehensive evaluation of ensembl, RefSeq, and UCSC annotations in
589 the context of RNA-seq read mapping and gene quantification. *BMC Genomics.* 2015;16: 97.
590 doi:10.1186/s12864-015-1308-8
- 591
- 592

593 **Figure captions**

594

595 **Fig 1. *Cis*-eQTL analysis pipeline to detect candidate-causal eQTLs of autoimmune disease.**

596 The 846 autoimmune disease associated SNPs per disease are documented in S1 Table and were LD
597 pruned to 584 independent loci (see Methods). Genotypes of 1000Genomes individuals were quality
598 controlled and subset to regions of recombination hotspots. If the lead GWAS SNP was found
599 between a recombination hotspot, then all SNPs were between the recombination hotspot intervals
600 were used in the Regulatory Trait Concordance (RTC) analysis. If the lead GWAS SNP was found
601 within a recombination hotspot itself, then all SNPs before or after the summit (including the between
602 summit SNPs) were used in the RTC (upper-interval and lower-interval hotspot respectively).
603 Normalized RNA-Seq expression data at gene-, isoform-, exon-, junction-, and intron-level were
604 obtained for the 1000Genomes individuals of the GEUVAIDS cohort in lymphoblastoid cell lines.
605 Disease associated SNPs with statistically significant association with gene expression ($q < 0.05$) and
606 an RTC score > 0.95 were classified as ‘candidate-causal eQTLs’. Genes whose expression is
607 modulated by the eQTL were defined as ‘candidate-causal eGenes’.

608

609 **Fig 2. Frequency of candidate-causal *cis*-eQTLs detected across each RNA-Seq quantification** 610 **type.**

611 Expression targets correspond to the quantification type under consideration (i.e. number of isoforms
612 captured at isoform-level, the number of exons-captured at exon-level). (A) Number of expression
613 targets and corresponding genes (by referencing the expression target back to the gene it belongs to)
614 interrogated in *cis* ($\pm 1\text{Mb}$) to the 846 autoimmune SNPs. (B) Number of significant *cis*-eQTL
615 associations that pass an FDR q -value threshold of 0.05, comprising the number of expression targets,
616 eQTLs, and eGenes. Only unique associations are reported (for example if two independent eQTLs
617 act on the same eGene, the eGene is only counted once). (C) Number of candidate-causal *cis*-eQTLs
618 ($q < 0.05$ and $\text{RTC} \geq 0.95$). (D) Percentage of statistically significant *cis*-eQTLs ($q < 0.05$) that are
619 candidate-causal ($q < 0.05$ and $\text{RTC} \geq 0.95$) to show the dropout of *cis*-eQTL associations that do not

620 appear to share the same causal variant as disease. (E) Percentage of the 846 autoimmune associated
621 SNPs that are candidate-causal *cis*-eQTLs, and the percentage of the 8,927 genes in *cis* that are
622 candidate-causal eGenes.

623

624 **Fig 3. Breakdown of autoimmune candidate-causal *cis*-eQTLs per RNA-Seq quantification type.**

625 (A) Percentage and number of candidate-causal *cis*-eQTLs detected per RNA-Seq quantification type,
626 following LD pruning of associated SNPs to 584 independent susceptibility loci. (B) Total candidate-
627 causal *cis*-eQTLs per disease across all five levels of RNA-Seq quantification (full results in found at
628 www.insidegen.com), using the 20 diseases of the ImmunoBase resource. In orange are disease-
629 associated SNPs that show no shared association with expression across any quantification type. In
630 blue are the disease-associated SNPs that are also candidate-causal *cis*-eQTLs. 47% of SNPs across
631 all diseases were candidate-causal *cis*-eQTLs on average. (C) Candidate-causal *cis*-eQTLs per disease
632 broken down by quantification type. (D) Candidate-causal *cis*-eQTLs detected per quantification type.
633 Percentage of candidate-causal *cis*-eQTLs captured are shown as a percentage of the 362 total.

634

635 **Fig 4. Functional annotation of Ankylosing Spondylitis risk loci using RNA-Seq.**

636 (A) Heatmap of the 14 candidate-causal *cis*-eQTLs and 27 eGenes of Ankylosing Spondylitis detected
637 across the five RNA-Seq quantification types (full results found at www.insidegen.com). Heat is
638 relative to *P*-value of association. To normalize across quantification types, relative significance of
639 each association per column was calculated as the $-\log_2(P/P_{max})$; where *P_{max}* is the most significant
640 association per quantification type. If a candidate-causal association is detected at any level of
641 quantification, it is shown and marked with an *. Associations not marked with an * are not
642 candidate-causal. (B) Isolation of the effect of Ankylosing Spondylitis associated SNP and candidate-
643 causal *cis*-eQTL, rs1128905, on the expression of *SNAPC4*. The risk allele is rs1128905 [C]. At exon-
644 level, rs1128905 is a candidate-causal *cis*-eQTL for exons 18 and 19 of *SNAPC4*, and at junction-
645 level for the exon 18-19 junction. (C) Incorporating existing functional genomic data from promoter
646 capture Hi-C data in lymphoblastoid cell lines [32]. The bait region encompassing rs1128905 interacts

647 with five targets with great confidence CHiCAGO score > 12 including four interactions with
648 *SNAPC4* which lies ~39kb away.

649

650 **Fig 5. Autoimmune associated *trans*-eQTLs detected using RNA-Seq.**

651 (A) Number of *trans*-eQTLs and *trans*-eGenes ($q < 0.01$) detected across all five RNA-Seq
652 quantification types (S3 Table). (B) Genome-wide depiction of *trans*-eQTL rs7726414, which is
653 associated with eight genes in *trans* detected using exon-level RNA-Seq. (C) Box-plots of rs7726414
654 on the eight *trans*-genes, the risk allele is rs7726414 [T].

655

656 **Supporting information**

657 **S1 Table.** Disease associated SNPs from the ImmunoBase resource taken forward for eQTL analysis.

658 Only autosomal SNPs with MAF > 5% included (noted with 'y'). SNPs are classified as being

659 associated with the 20 immune-related diseases with 'y' (associated) or 'n' (not associated).

660 Abbreviations for diseases are taken from ImmunoBase (www.immunobase.org).

661

662 **S2 Table.** Candidate-causal *cis*-eQTLs detected across all five RNA-Seq quantification types. See tab

663 2 for a definition of all column headers. Only candidate-causal ($q < 0.05$ and $RTC \geq 0.95$) are shown.

664 Data are also accessible through the web-portal (www.insidegen.com).

665

666 **S3 Table.** *Trans*-eQTLs detected across all five RNA-Seq quantification types.

667

668 **S1 Fig.** Correlation of SNP-Gene association pairs across RNA-Seq quantification types. The bottom

669 panel shows the correlation coefficient of the effect sizes (beta) of candidate-causal *cis*-eQTL

670 associations across each RNA-Seq quantification type. The top panel shows the same data but is

671 adjusted so to force the same direction of effect.

672

673 **S2 Fig.** Processing of genotype data and principle component analysis. Genotype data in VCF format

674 of 1000Genomes individuals were downloaded from E-GEUV1 (ArrayExpress). Insertion-deletion

675 sites were removed, and bi-allelic SNPs kept only. SNPs with $HWE < 0.0001$ were removed and the

676 VCF converted to 0,1,2 format using PLINK. Principle component analysis was performed on

677 genotype data using the R package SNPRelate on chromosome 20. The first 3 components were

678 included in the eQTL regression model as well as the binary imputation status (see methods).

679

680 **S3 Fig.** Replication of *cis*-eQTL effect of rs30187 on *ERAP1* and rs2910686 on *ERAP2* expression-

681 levels using gene-level quantification.

682

683 **S4 Fig.** Functional prioritization of the rs1128905 (*CARD9-SNAPC4*) locus. rs1128905, associated
684 with Ankylosing Spondylitis, was found to be a candidate-causal *cis*-eQTL for *CARD9* and *SNAPC4*.
685 The rs1128905 locus interacts with the *SNAPC4* locus (~39kb away) via a chromatin interaction in
686 lymphoblastoid cell lines (Fig 4C). To find potential causal SNPs, we took all associated SNPs with in
687 strong LD $r^2 > 0.8$ with rs1128905, and looked for colocalisation with genome-wide epigenetic marks
688 from histone and DNase from the ENCODE project in lymphoblastoid cell lines (GM12878). Two
689 SNPs, rs10870201 and rs10870202, were found to be in the peak summit of enhancer/promoter
690 histone marks H3K4me3, H3K27ac, and H3K9ac, as well as in a region of DNase hypersensitivity.
691 We hypothesize variation at these SNPs causes loss of interaction with *SNAPC4*, which reduces the
692 expression of *SNAPC4* (Fig 4B).

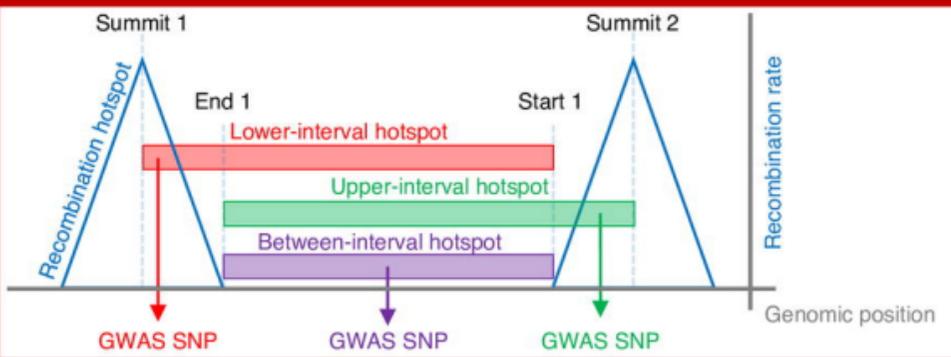
1. Autoimmune disease-associated SNPs

846 lead-SNP associations (584 independent loci), passing genome-wide significance curated by the ImmunoBase Cohort.

2. Genotype

1000 Genomes
373 Europeans (CEU, FIN, GBR, TSI)
SNPS: biallelic, HWE $p > 0.0001$, MAF ≥ 0.05

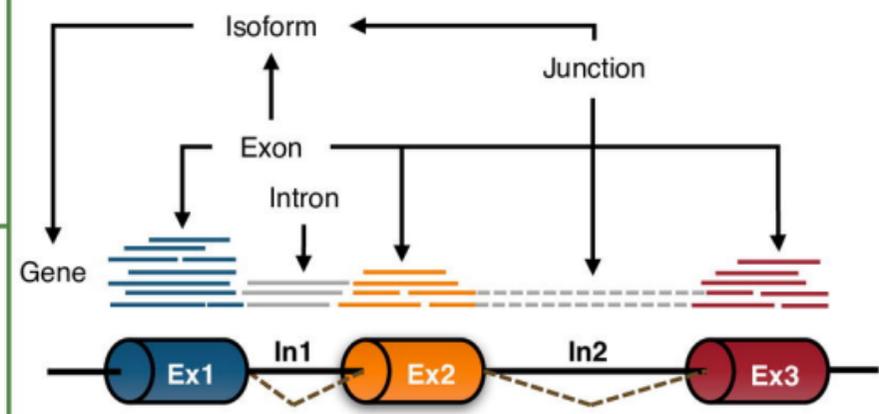
Subset and retrieve all SNPs within hotspot intervals around each GWAS SNP.
846 independent SNPs in **451** unique hotspots, corresponding to a total of **137,509** SNPs.



3. Expression

GEUVADIS Consortium
RNA-Seq of 373 Europeans
Lymphoblastoid cell lines
QC passed (expressed in $> 50\%$ of individuals)
PEER Factor normalized and SN transformed RPKM
Subset to *cis*-intervals around each GWAS SNP (± 1 Mb)

Quantification of expression targets in *cis* (± 1 Mb)
Genes (7,969 genes)
Isoforms (28,220 isoforms, 7,636 genes)
Exons (54,043 exons, 5,621 genes)
Junctions (49,909 junctions, 4,366 genes)
Introns (35,662 introns, 6,443 genes)

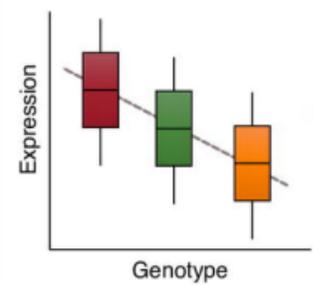


4. *Cis*-eQTL association analysis

Per quantification type, perform linear regression of all hotspot SNPs against all *cis* (± 1 Mb) normalized expression targets. Include four genotype principle components in regression model.

5. Test for shared causal variant

1) Identify best hotspot *cis*-eQTL by linear regression for each unit:



2) Perform linear regression of all hotspot *cis*-eQTLs against residuals of best then rank:

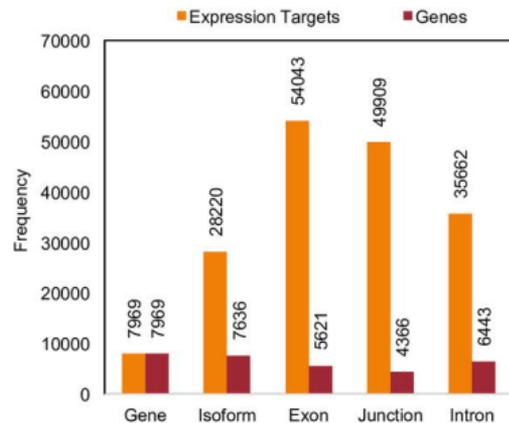
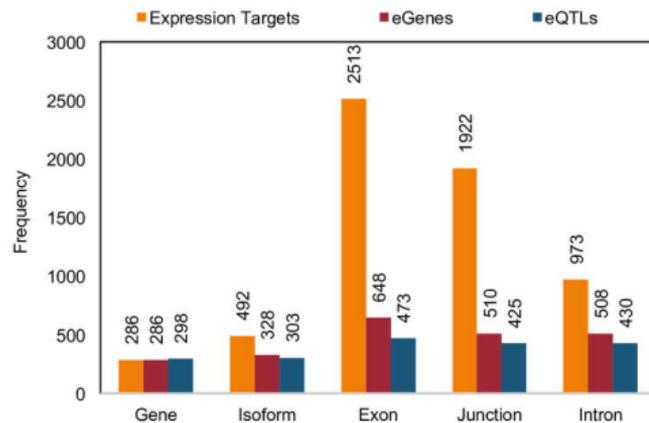
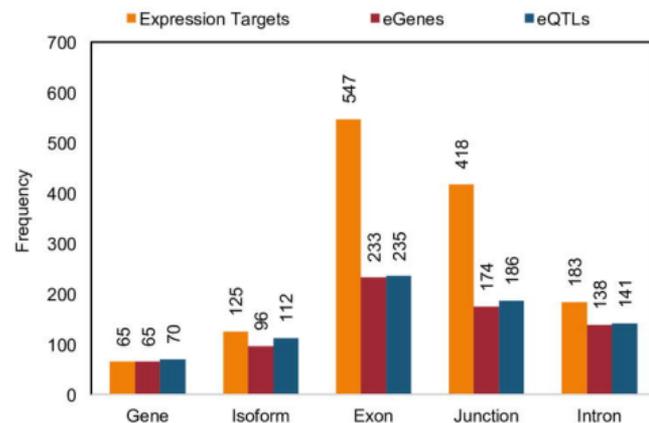
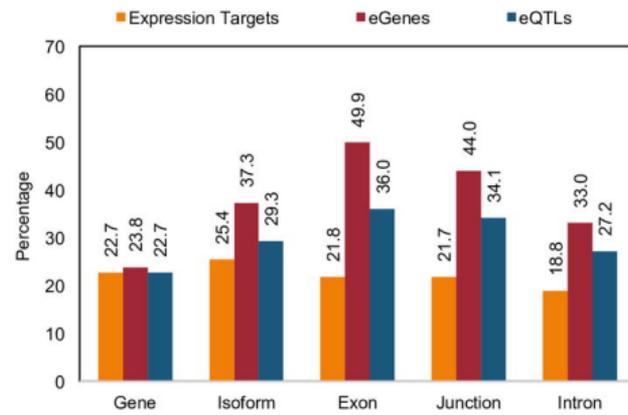
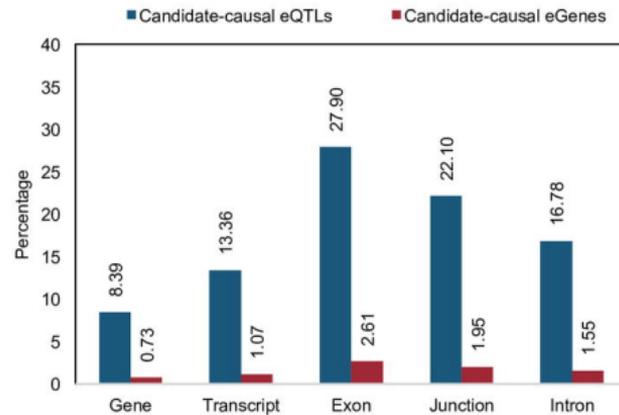
SNP	P-val	Rank
Best	Least sig.	1
GWAS
SNP _i
SNP _j	Most sig.	N _{SNPs}

3) Calculate RTC:

$$RTC = \frac{N_{SNPs} - Rank_{GWAS\ SNP}}{N_{SNPs}}$$

6. Define candidate-causal eQTLs and eGenes

FDR adjustment per quantification type.
q < 0.05
RTC ≥ 0.95

A Number of expression targets & corresponding genes interrogated
(in *cis* to the 846 autoimmune associated SNPs; +/-1Mb)**B** Number of statistically significant *cis*-eQTL associations
($q < 0.05$)**C** Number of candidate-causal *cis*-eQTL associations
($q < 0.05$; $RTC > 0.95$)**D** Percentage of statistically significant *cis*-eQTLs that are candidate-causal**E** Percentage of total SNPs & genes interrogated that are candidate-causal
(846 SNPs, 8,927 genes total)

Candidate-causal *cis*-eQTLs per independent loci (584 total)

■ Candidate-causal *cis*-eQTL ■ Not a candidate-causal *cis*-eQTL

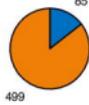
All RNA-Seq resolutions



Gene-level



Isoform-level



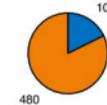
Exon-level



Junction-level

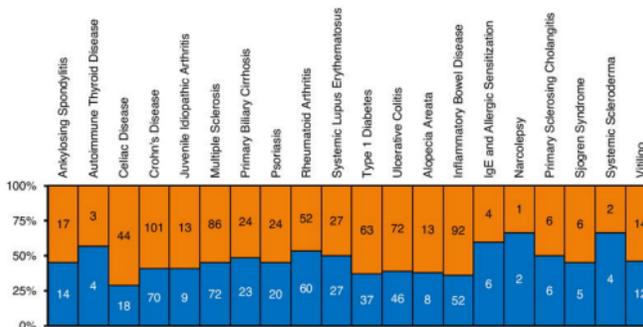


Intron-level



Total candidate-causal *cis*-eQTLs per disease across all levels of RNA-Seq

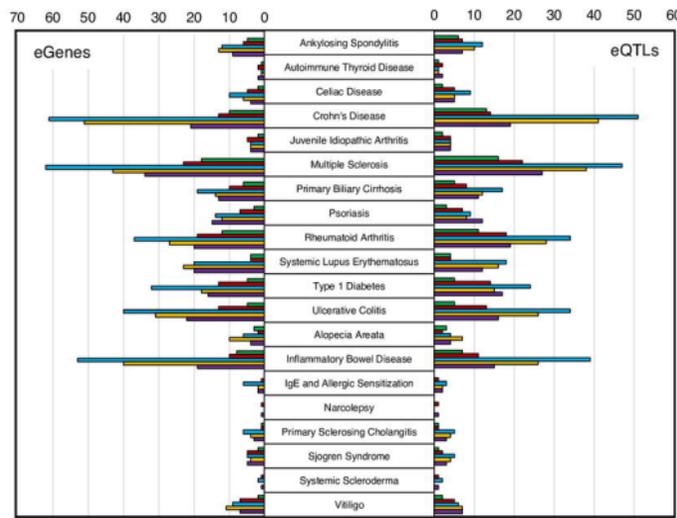
■ Candidate-causal *cis*-eQTL ■ Not a candidate-causal *cis*-eQTL



C

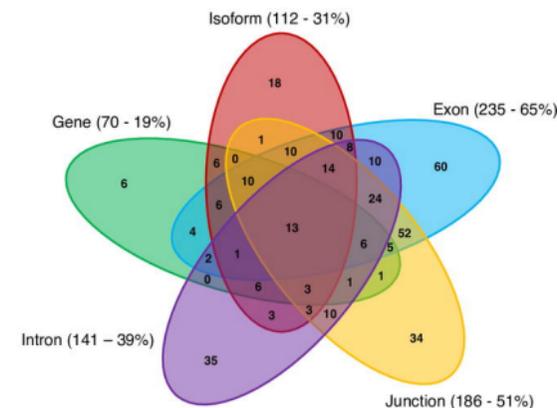
Candidate-causal *cis*-eQTLs and eGenes per disease

■ Gene ■ Isoform ■ Exon ■ Junction ■ Intron

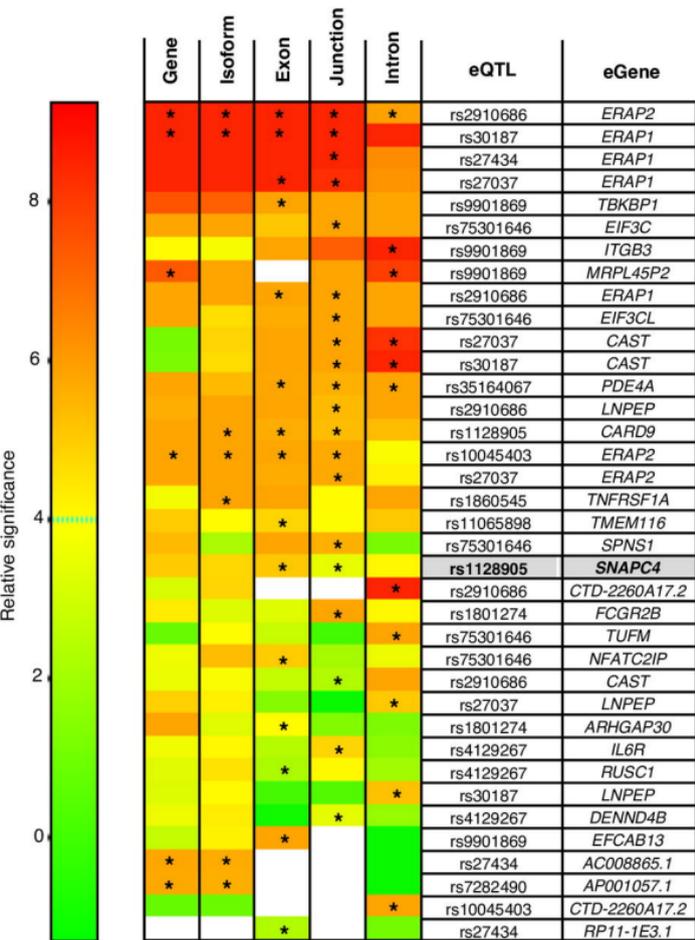


D

Candidate-causal *cis*-eQTLs per quantification type

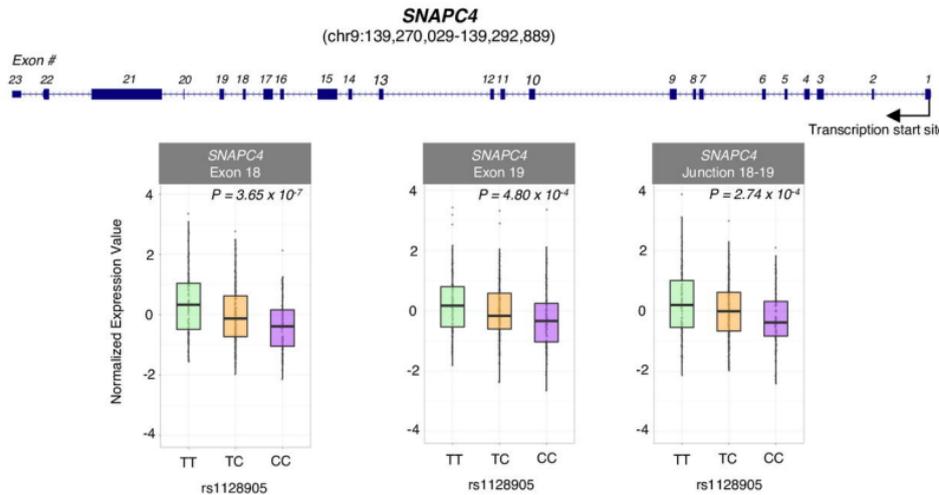


Candidate-causal *cis*-eQTLs and eGenes of Ankylosing Spondylitis
14 candidate-causal *cis*-eQTLs, 27 eGenes total



B

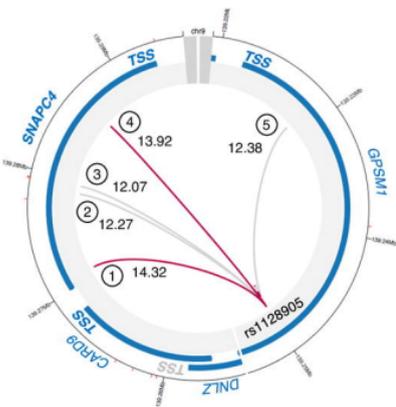
rs1128905 is a candidate-causal *cis*-eQTL of *SNAPC4* at exon-level and junction-level



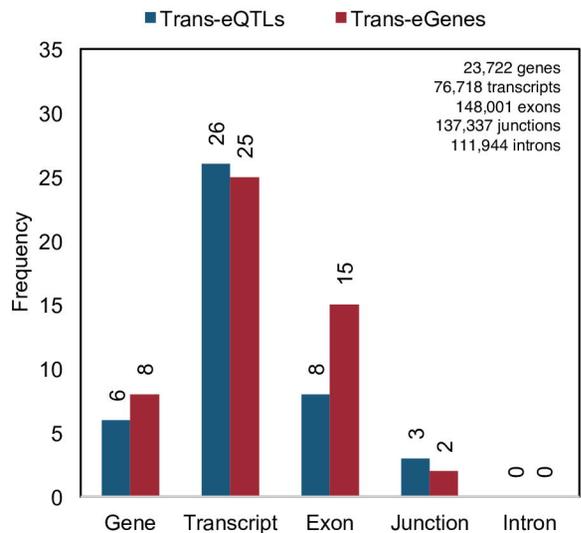
C

Chromatin interaction of rs1128905 with *SNAPC4* in LCLs

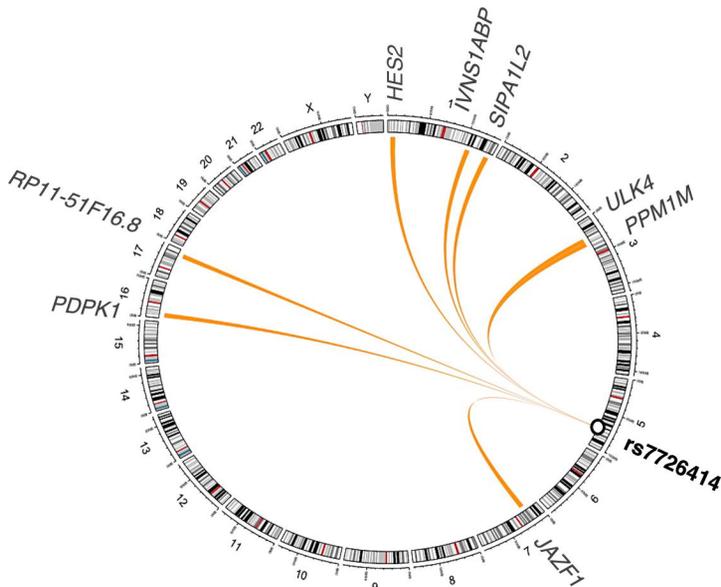
(chr9:139,219,258-139,295,668)



A Number of trans-eQTLs and trans-eGenes
($q < 0.01$)



B SLE associated rs7726414 is a *trans*-eQTL for eight eGenes
Exon-level ($q < 0.01$)



C

Trans-eGenes of rs7726414 detected at exon-level

