

1 **TITLE: EXPLORING THE OVERLAP BETWEEN RHEUMATOID ARTHRITIS**

2 **SUSCEPTIBILITY LOCI AND LONG NON-CODING RNA ANNOTATIONS**

3 **RUNNING TITLE: RHEUMATOID ARTHRITIS RISK LOCI AND LONG NON-CODING RNA**

4 James Ding¹ (james.ding@manchester.ac.uk), Chenfu Shi¹, John Bowes^{1,2}, Stephen
5 Eyre^{1,2}, Gisela Orozco^{1,2}

6 ¹Centre for Genetics and Genomics Versus Arthritis. Division of Musculoskeletal and
7 Dermatological Sciences, School of Biological Sciences, Faculty of Biology, Medicine
8 and Health, The University of Manchester, UK.

9 ²NHR Manchester Biomedical Research Centre, Manchester University NHS
10 Foundation Trust, Manchester Academic Health Science Centre, Manchester, M13
11 9WL, UK.

12

13 **ABSTRACT**

14 Whilst susceptibility variants for many complex diseases, such as rheumatoid arthritis
15 (RA), have been well characterised, the mechanism by which risk is mediated is still
16 unclear for many loci. This is especially true for the majority of variants that do not
17 affect protein-coding regions. lncRNA represent a group of molecules that have been
18 shown to be enriched amongst variants associated with RA and other complex
19 diseases, compared to random variants. In order to establish to what degree direct
20 disruption of lncRNA may represent a potential mechanism for mediating RA
21 susceptibility, we chose to further explore this overlap. By testing the ability of
22 annotated features to improve a model of disease susceptibility, we were able to
23 demonstrate a local enrichment of enhancers from immune-relevant cell types
24 amongst RA susceptibility variants (\log_2 enrichment 3.40). This was not possible for
25 lncRNA annotations in general, however a small, but significant enrichment was
26 observed for immune-enriched lncRNA (\log_2 enrichment 0.867002). This enrichment
27 was no longer apparent when the model was conditioned on immune-relevant
28 enhancers (\log_2 enrichment -0.372734), suggesting that direct disruption of lncRNA
29 sequence, independent of enhancer disruption, does not represent a major
30 mechanism by which susceptibility to complex diseases is mediated. Furthermore, we
31 demonstrated that, in keeping with general lncRNA characteristics, immune-enriched
32 lncRNA are expressed at low levels that may not be amenable to functional
33 characterisation.

34 INTRODUCTION

35 In keeping with other complex diseases, rheumatoid arthritis (RA) susceptibility loci are
36 mainly non-coding, with relatively few variants having a potential impact upon the
37 coding sequence for a protein [1,2]. Enhancers have been identified as likely to
38 mediate disease susceptibility at many loci. Evidence to support this generalization
39 includes demonstrated effects of genome wide association study (GWAS) variants on
40 enhancers at individual loci [3], as well as an enrichment of RA GWAS variants amongst
41 enhancers from relevant cell types [4]. Alternative non-coding elements, such as long
42 non-coding RNA (lncRNA) may also play a role in mediating the increased risk
43 associated with non-coding variants.

44 lncRNA are a heterogeneous class of molecules that are defined based on a lack of
45 protein-coding potential and a minimum transcribed length of 200 nucleotides. Whilst
46 discrete subcategories exist, such as long intergenic non-coding RNAs, promoter
47 associated-lncRNAs or antisense lncRNA with some discriminatory characteristics,
48 including genomic context, overlapping chromatin marks, length and structure [5], it
49 can still be difficult to discriminate a genuine lncRNA annotation from a product of
50 spurious transcription. Many individual lncRNA have been functionally characterised,
51 with gene regulation featuring frequently amongst the wide variety of roles described.
52 One important subcategory was identified following observations of RNA polymerase II
53 recruitment and transcription at enhancers [6]. Often described as enhancer derived
54 RNA (eRNA), expression of these transcripts is highly correlated with enhancer activity

55 [7] and increasing evidence suggests that these ncRNA may contribute towards
56 enhancer function, although the precise mechanism is still unclear [8].

57 There is some evidence to suggest that GWAS susceptibility variants are enriched
58 amongst lncRNA [9,10]. However, using conventional methods of determining whether
59 annotations overlap more than can be expected by chance it is difficult to
60 appropriately account for confounding factors, such as chromosomal compartments or
61 chromatin accessibility. Using the *ab initio* MiTranscriptome assembly (58,648
62 lncRNA), which was generated using a large collection of RNA sequencing libraries,
63 GWAS SNPs were demonstrated to be enriched in lncRNA, compared to other SNPs
64 tested for in GWAS analyses [9]. This enrichment was also observed using either GWAS
65 SNPs or probabilistically identified causal SNPs (PICS) and lncRNA from the functional
66 annotation of the human genome (FANTOM) cap-analysis gene expression (CAGE)
67 associated transcriptome (CAT) assembly (27,919 lncRNA), generated using CAGE
68 datasets in combination with existing assemblies [10]. Using the tissue specific nature
69 of the FANTOM CAT annotation it was also possible to demonstrate that this
70 enrichment was markedly higher when testing specifically for an enrichment of
71 immune-relevant GWAS PICS in immune-expressed lncRNA transcripts.

72 Despite these studies, the relevance of lncRNA to the study of individual complex
73 diseases, such as RA, remains unclear. This is especially true given the overlap in
74 genomic locality and function between enhancers and lncRNA. We chose to investigate
75 the overlap between lncRNA annotations, enhancer annotations and GWAS SNPs
76 associated with RA susceptibility, with the aim of establishing to what degree the

77 direct disruption of lncRNA by RA-associated variants may contribute to the mediation
78 of disease risk. Central to our investigation is the use of the fgwas algorithm [11],
79 which tests the ability of individual annotations to improve a probabilistic model of
80 disease susceptibility, constructed using GWAS summary statistics. Using this method a
81 local enrichment is estimated, that takes into consideration the non-random
82 distribution of potentially confounding genomic features. In addition, it is possible to
83 model multiple traits and establish the degree to which they are independently
84 predictive.

85 **MATERIALS AND METHODS**

86 **ENRICHMENT TESTING**

87 Enrichment testing was performed using RA summary statistics [1] and fgwas v0.3.6
88 [11], with the “-cond” option called for conditional analyses. Chromatin state data was
89 obtained from the Roadmap Epigenomics project [12], with the expanded, 18-state
90 model used for all 98 corresponding epigenomes. The 18-states were combined to
91 form four exclusive annotations as follows: TSSs: Active TSS (1), Flanking TSS (2),
92 Flanking TSS upstream (3), Flanking TSS downstream (4), Bivalent/poised TSS (14).
93 Transcription: Strong transcription (5), Weak Transcription (6). Enhancers: Genic
94 enhancer 1 (7), Genic enhancer 2 (8), Active enhancer 1 (9), Active enhancer 2 (10),
95 Weak enhancer (11), Bivalent enhancer (15). Repressed chromatin: ZNF genes and
96 repeats (12), Heterochromatin (13), Repressed polycomb (16), Weak repressed
97 polycomb (17), Quiescent/low (18). The following lncRNA datasets were interrogated:
98 Lncipedia v5.2 [13], miTranscriptome v2 [9] and FANTOM CAT (robust) [10]. Immune-

99 relevant enhancers were defined as genomic regions annotated as enhancers in cell-
100 types defined as originating from “Blood and T cell” or “HSC and B cell” by the
101 Roadmap Epigenomics project. The definition of Immune-enriched lncRNA is based on
102 the underlying sample ontology and was wholly adopted from Hon et al. requiring:
103 detection in at least 50% of immune-relevant samples, 5 x higher expression in
104 immune-relevant samples than in other samples and $P < 0.05$ in a one-tailed Mann-
105 Whitney rank sum test [10].

106 **EXPRESSION PROFILING**

107 Raw RNA-seq reads were downloaded from the Roadmap Epigenomics project for
108 primary T-helper cells (SRA accession SRR644513 and SRR643766) [12]. Reads were
109 then filtered for quality, adapter content and polyA tails using fastp version 19.7, with
110 default settings and polyX tail trimming enabled. Transcripts were quantified using
111 Salmon version 13.1 [14] using suggested settings (“quant” mode and “-
112 validatemappings”) and using the reference index generated from the FANTOM CAT
113 robust database [10]. Transcripts quantifications (reported as Transcripts per million,
114 TPM) were then remapped to genes and summed for each gene. CAGE transcript
115 counts per million (CPM) were taken from FANTOM CAT for T-helper cells
116 (CL_0000084_T_cell), as published [10].

117 Statistical difference between the distributions of immune-enriched lncRNA and mRNA
118 abundance was established using a two-sided Welch’s t-test, with no assumption of
119 equal variance. This test assumes a normal distribution of the mean. Using the Mann-

120 Whitney u-test, which does not require this assumption, it was not possible to
121 estimate a p-value as it is too close to zero.

122 **RESULTS**

123 **ENHANCER ANNOTATIONS FROM IMMUNE-RELEVANT CELL TYPES ARE ENRICHED AMONGST RA**

124 **SUSCEPTIBILITY VARIANTS**

125 An enrichment of RA PICS has previously been demonstrated amongst *cis*-regulatory
126 elements that are active in T-helper cells and lymphoblastoid cells [4], using data from
127 the high density ImmunoChip custom SNP array. In order to establish confidence in the
128 ability of fgwas to identify similar enrichments we sought to validate this evidence of
129 enrichment using a more inclusive approach that incorporates the probability of
130 association for all SNPs tested in the most recent RA GWAS meta-analysis [1]. In order
131 to achieve this, chromatin state data taken from the Roadmap Epigenomics project
132 [12] was incorporated in a model of RA susceptibility. Using enrichment estimates
133 generated using the expanded 18-state model, it is possible to discern an enrichment
134 of certain chromatin states, such as genic enhancers, active enhancers and weak
135 enhancers in immune-relevant cell types, such as B cells, T cells and monocytes (panel
136 A in S1 Fig), however the confidence intervals (CIs) for estimates are broad for many
137 states, likely due to a reduced abundance of these annotations (panel B and C in S1
138 Fig). The size of CIs was improved by combining states to generate four more easily
139 interpretable annotations (enhancers, transcription start sites (TSSs), transcription,
140 and repressed chromatin; Fig 1). In keeping with our understanding of RA, the highest
141 level of enrichment was observed for enhancer annotations in immune-relevant cell

142 types, with regulatory T cells showing the highest enrichment (\log_2 enrichment 3.17,
143 95% CI 2.58; 3.75015). In immune-relevant cell types TSSs were also enriched, whilst
144 repressed chromatin was depleted (Fig 1).

145 **LNCRNA ANNOTATIONS SHOW NEGLIGIBLE ENRICHMENT AMONGST RA SUSCEPTIBILITY VARIANTS**

146 We applied fgwas to test for an enrichment of lncRNA amongst RA susceptibility
147 variants, using Incipedia, a large lncRNA database curated from a number of sources
148 [13], the MiTranscriptome assembly [9], and the FANTOM CAT assembly [10]. Using
149 fgwas, MiTranscriptome lncRNA genes show a level of depletion amongst RA
150 susceptibility variants comparable to repressed chromatin. MiTranscriptome lncRNA
151 exons and both genes and exons from either FANTOM CAT or Incipedia all show
152 negligible enrichment (Fig 2A).

153 **LNCRNA ANNOTATIONS WITH ENRICHED EXPRESSION IN IMMUNE-RELEVANT CELLS ARE SUBTLY** 154 **ENRICHED AMONGST RA SUSCEPTIBILITY VARIANTS**

155 Uniquely, FANTOM CAT transcripts are associated with tissue specific expression data.
156 As in the original FANTOM CAT publication, we took advantage of this additional
157 information, to test for an enrichment of lncRNA whose expression is enriched in
158 immune-relevant cell types amongst RA susceptibility loci. This approach
159 demonstrated a subtle enrichment of lncRNA genes (\log_2 enrichment 0.867, 95% CI
160 0.0554; 1.57) and similar level of enrichment for their exons, albeit with an increased
161 confidence interval (\log_2 enrichment 0.799, 95% CI 2.30; 1.94). FANTOM CAT mRNA
162 annotations, whose expression is enriched in immune-relevant cell types, were
163 included in order to provide a comparison. Genic mRNA annotations showed a similar

164 level of enrichment as lncRNA genes, with mRNA exons exhibiting slightly higher
165 enrichment (Fig 2B).

166 **THE SUBTLE ENRICHMENT OF IMMUNE-ENRICHED LNCRNA OBSERVED IS NOT INDEPENDENT OF**
167 **IMMUNE-RELEVANT ENHANCER ANNOTATIONS**

168 Given the strong enrichment of immune-relevant enhancers amongst RA susceptibility
169 loci and the established overlap between lncRNA annotations and enhancers, we were
170 interested to investigate the independence of these variables using fgwas. In this
171 conditional analysis, a residual enrichment of immune-enriched FANTOM CAT
172 annotations was tested after the enrichment of immune-relevant enhancers (\log_2
173 enrichment 3.40, 95% CI 2.54; 4.58) was accounted for (Fig 3Fig 3). Interestingly, both
174 lncRNA and mRNA annotations no longer show significant enrichment, indicating that
175 once enrichment of susceptibility variants in enhancers has been accounted for no
176 remaining enrichment of mRNA or lncRNA is apparent.

177 **IMMUNE-ENRICHED LNCRNA ARE EXPRESSED AT LOW LEVELS IN RA RELEVANT CELL TYPES**

178 lncRNA are generally considered to exhibit low expression levels and high tissue
179 specificity that can make them difficult to study using conventional methods. Given
180 their enrichment for RA susceptibility variants we were interested to establish whether
181 this description applied to FANTOM CAT immune-enriched lncRNA in an RA relevant
182 cell type. Using randomly primed Roadmap Epigenomics RNA-seq data from primary T-
183 helper cells, the distribution of expression levels for FANTOM CAT immune-enriched
184 lncRNA is significantly lower than that of FANTOM CAT immune-enriched mRNA ($p =$
185 4.46×10^{-34} , Fig 4A median lncRNA transcripts per million reads (TPM); 0.257, vs 124

186 for mRNA). The same is true using less-conventional expression profiling methods,
187 such as those employed by FANTOM CAT, which offer improved sensitivity for the
188 detection of transcripts of low abundance ($p = 1.61 \times 10^{-24}$, Fig 4B median lncRNA
189 counts per million reads (CPM); 1.43, vs 71.7 for mRNA). 90% of immune-enriched
190 lncRNA have abundance lower than 85.2% of immune-enriched mRNA in Roadmap
191 Epigenomics RNA-seq data or 75.6% of immune-enriched mRNA in FANTOM CAT CAGE
192 data.

193 **DISCUSSION**

194 By incorporating both cell-type specific enhancer and lncRNA annotations into a
195 probabilistic model of RA susceptibility it is possible to demonstrate their respective
196 levels of enrichment amongst RA susceptibility variants. Whilst previous studies have
197 demonstrated an enrichment of lncRNA compared with randomly shuffled annotations
198 these analyses fail to take into consideration the complex organisation of the genome
199 and are easily confounded by alternative features. In our analysis, which incorporated
200 lncRNA from various databases, it was only possible to demonstrate a subtle
201 enrichment of lncRNA whose expression was previously identified as being enriched in
202 relevant cell-types.

203 By conditioning a model of RA susceptibility on enhancer annotations from immune-
204 relevant cell types it is possible to test the independence of additional features. This
205 demonstrated that the subtle enrichment observed for immune-relevant lncRNA is
206 entirely explained by immune-relevant enhancer annotations. Interestingly, the same
207 is true of mRNA annotations, indicating that in both instances the primary influence of

208 susceptibility variants is in affecting non-coding regulatory elements, with any effect
209 on mRNA and lncRNA being secondary and/or indirect. This suggests that the majority
210 of genetic variance in disease susceptibility is mediated through disruption of
211 regulatory elements and not through direct disruption of transcript sequences. This
212 observation is in keeping with those made previously, relating to the minimal overlap
213 between the coding sequence of mRNA and RA susceptibility variants. It is also
214 consistent with the existence of well-characterised effects of RA susceptibility variants
215 on coding regions, such as for the HLA proteins [15], which convey a significant, but
216 not exhaustive proportion of risk.

217 Similarly, this analysis does not rule out the relevance of lncRNA at individual loci, in
218 fact it is worth noting that there is evidence to suggest that C5T1lncRNA may mediate
219 risk at a RA risk locus located at chromosomal position 9q33.2, with RA associated
220 variants falling within a C5T1lncRNA exonic region [16]. Our analysis, however,
221 precludes an effect of RA susceptibility variants on the transcribed sequence of lncRNA
222 independent of enhancer disruption. Sequence-specific functions of lncRNA are,
223 therefore, unlikely to mediate a significant proportion of the risk-modifying effect
224 associated with RA susceptibility variants. Furthermore, the distribution of FANTOM-
225 CAT immune-enriched lncRNA expression levels in primary T-helper cells highlight the
226 difficulties associated with studying lncRNA, whose expression is typically very low and
227 highly cell-type specific.

228 The analyses performed are specific to RA, however it is assumed that similar results
229 would be reached using GWAS data, enhancer annotations and lncRNA annotations

230 relevant to other diseases. This would, therefore, suggest that the enrichment of
231 lncRNA annotations amongst GWAS variants observed by others [9,10], may result
232 from a high degree of overlap between regulatory features such as enhancers and
233 lncRNA and other confounding features, such as active and inactive genomic
234 compartments. As a generalisation, when it comes to functional characterization of
235 variants associated with complex genetic disorders, sequence-specific lncRNA
236 functions is unlikely to represent an attractive area for study; lncRNA are not
237 independently enriched amongst such variants and are difficult to study, due to low
238 expression levels. Despite these results, dysregulation of lncRNA expression could still
239 play a role in RA and similar diseases, with disease associated variants affecting
240 regulatory elements, such as enhancers that control lncRNA expression.

241 fgwas represents a useful tool for studying the enrichment of different features
242 amongst susceptibility variants, especially when used in combination with Roadmap
243 Epigenomics data in order to identify cell-types and tissues that are relevant for
244 disease susceptibility. Whilst using this tool we observed that the number and size of
245 annotations that are tested have a strong influence, both on the confidence with
246 which any enrichment is estimated, as well as on the extent of that enrichment. It is
247 likely that this may explain some of the subtle differences observed between different
248 lncRNA databases.

249 Our analyses highlight the caveats associated with inferring functional relevance for a
250 given feature, based purely on the observation of enrichment over a genomic
251 background, as well as the care that must be taken when attempting to interpret such

252 enrichments. By deriving enhancer and lncRNA annotations from entirely different
253 sources we have tried to ensure that the demonstrated dependence is not self-
254 fulfilling, as it may have been if we defined immune-relevant lncRNA based on
255 underlying chromatin states.

256 In conclusion, using fgwas and Roadmap Epigenomics chromatin state data it is
257 possible to identify cell types and chromatin states of relevance to complex diseases,
258 such as RA. In the case of RA this is predominantly enhancers and transcription start
259 sites from immune-relevant cell types. It is also possible to test the association of
260 alternative features and establish their independence from chromatin states. Here, a
261 previously described enrichment of lncRNA amongst GWAS susceptibility loci was
262 explored for RA. Immune-enriched lncRNA from the FANTOM-CAT database were
263 found to be enriched amongst RA susceptibility loci, however, this enrichment was not
264 apparent when chromatin-state data was taken into account.

265 Our results suggest that regulatory elements, such as enhancers, are likely to mediate
266 the vast majority of variance in risk associated with RA and other complex diseases,
267 with no substantial independent contribution being made by direct disruption of
268 lncRNA sequences. Because of this, and the difficulties associated with detecting
269 transcripts of such low abundance, sequence-specific lncRNA function does not
270 represent the most attractive area for study with respect to RA susceptibility, except in
271 the case of in depth characterisation of individual loci.

272 **ACKNOWLEDGEMENTS**

273 The authors would like to acknowledge the assistance given by IT Services and the use
274 of the Computational Shared Facility at The University of Manchester.

275 REFERENCES

- 276 1. Okada Y, Wu D, Trynka G, Raj T, Terao C, Ikari K, et al. Genetics of rheumatoid
277 arthritis contributes to biology and drug discovery. *Nature*. 2013;506: 376–381.
278 doi:10.1038/nature12873
- 279 2. Eyre S, Bowes J, Diogo D, Lee A, Barton A, Martin P, et al. High-density genetic
280 mapping identifies new susceptibility loci for rheumatoid arthritis. *Nat Genet*.
281 2012;44: 1336–40. doi:10.1038/ng.2462
- 282 3. Simeonov DR, Gowen BG, Boontanrart M, Roth TL, Gagnon JD, Mumbach MR, et
283 al. Discovery of stimulation-responsive immune enhancers with CRISPR
284 activation. *Nature*. 2017;549: 111–115. doi:10.1038/nature23875
- 285 4. Farh KK-H, Marson A, Zhu J, Kleinewietfeld M, Housley WJ, Beik S, et al. Genetic
286 and epigenetic fine mapping of causal autoimmune disease variants. *Nature*.
287 2014;518: 337–43. doi:10.1038/nature13835
- 288 5. Jarroux J, Morillon A, Pinskaya M. History, discovery, and classification of
289 lncRNAs. *Advances in Experimental Medicine and Biology*. Springer, Singapore;
290 2017. pp. 1–46. doi:10.1007/978-981-10-5203-3_1
- 291 6. Kim T-K, Hemberg M, Gray JM, Costa AM, Bear DM, Wu J, et al. Widespread
292 transcription at neuronal activity-regulated enhancers. *Nature*. 2010;465: 182–
293 7. doi:10.1038/nature09033

- 294 7. Andersson R, Gebhard C, Miguel-Escalada I, Hoof I, Bornholdt J, Boyd M, et al.
295 An atlas of active enhancers across human cell types and tissues. *Nature*.
296 2014;507: 455–61. doi:10.1038/nature12787
- 297 8. Lam MTY, Li W, Rosenfeld MG, Glass CK. Enhancer RNAs and regulated
298 transcriptional programs. *Trends Biochem Sci*. 2014;39: 170–82.
299 doi:10.1016/j.tibs.2014.02.007
- 300 9. Iyer MK, Niknafs YS, Malik R, Singhal U, Sahu A, Hosono Y, et al. The landscape of
301 long noncoding RNAs in the human transcriptome. *Nat Genet*. 2015;47: 199–
302 208. doi:10.1038/ng.3192
- 303 10. Hon C-C, Ramilowski JA, Harshbarger J, Bertin N, Rackham OJL, Gough J, et al. An
304 atlas of human long non-coding RNAs with accurate 5' ends. *Nature*. 2017;543:
305 199–204. doi:10.1038/nature21374
- 306 11. Pickrell JK. Joint analysis of functional genomic data and genome-wide
307 association studies of 18 human traits. *Am J Hum Genet*. 2014;94: 559–73.
308 doi:10.1016/j.ajhg.2014.03.004
- 309 12. Kundaje A, Meuleman W, Ernst J, Bilenky M, Yen A, Heravi-Moussavi A, et al.
310 Integrative analysis of 111 reference human epigenomes. *Nature*. 2015;518:
311 317–330. doi:10.1038/nature14248
- 312 13. Volders P-J, Anckaert J, Verheggen K, Nuytens J, Martens L, Mestdagh P, et al.
313 LNCipedia 5: towards a reference set of human long non-coding RNAs. *Nucleic
314 Acids Res*. 2019;47: D135–D139. doi:10.1093/nar/gky1031

- 315 14. Patro R, Duggal G, Love MI, Irizarry RA, Kingsford C. Salmon provides fast and
316 bias-aware quantification of transcript expression. *Nat Methods*. 2017;14: 417–
317 419. doi:10.1038/nmeth.4197
- 318 15. Raychaudhuri S, Sandor C, Stahl EA, Freudenberg J, Lee H-S, Jia X, et al. Five
319 amino acids in three HLA proteins explain most of the association between MHC
320 and seropositive rheumatoid arthritis. *Nat Genet*. 2012;44: 291–6.
321 doi:10.1038/ng.1076
- 322 16. Messemaker TC, Frank-Bertoncelj M, Marques RB, Adriaans A, Bakker AM, Daha
323 N, et al. A novel long non-coding RNA in the rheumatoid arthritis risk locus
324 TRAF1-C5 influences C5 mRNA levels. *Genes Immun*. 2015;
325 doi:10.1038/gene.2015.54
- 326

327 **FIGURE CAPTIONS**

328 **FIG 1. ENRICHMENT OF CHROMATIN STATE GROUPS AMONGST RA SUSCEPTIBILITY**

329 **VARIANTS FOR 98 CELL TYPES**

330 Estimates for enrichment of combined chromatin state groupings are illustrated for all
331 98 cell types annotated within the Roadmap Epigenomics 18-state model. Cell-types
332 are ordered and coloured according to the clustering established by the Roadmap
333 Epigenomics project, with immune-relevant cell types coloured green. Estimates and
334 confidence intervals are clipped at axis limits, where applicable.

335 **FIG 2. ENRICHMENT OF LNCRNA ANNOTATIONS AMONGST RA SUSCEPTIBILITY**

336 **VARIANTS**

337 Estimates for the enrichment of genic (black circle) and exonic (grey diamond)
338 annotations from a variety of lncRNA containing databases, including 95% confidence
339 intervals (A). Separate estimates are included for annotations identified as exhibiting
340 enriched expression in immune-relevant cells (B).

341 **FIG 3. ENRICHMENT OF IMMUNE-ENRICHED LNCRNA AMONGST RA SUSCEPTIBILITY**

342 **VARIANTS AFTER CONDITIONING ON CHROMATIN STATE DATA**

343 The influence of immune-relevant enhancers (red circle) was fixed in a probabilistic
344 model of RA susceptibility to determine whether the subtle enrichment of FANTOM
345 CAT immune-enriched lncRNA or mRNA adds any additional predictive information and
346 is therefore independently enriched. Genic (black circle) and exonic (grey diamond)
347 annotations were both tested.

348 As may be expected given the magnitude of enrichments observed, after accounting
349 for the effect of immune-enriched FANTOM CAT annotations, the residual enrichment
350 of immune-relevant enhancers is not dramatically reduced (S2 Fig).

351 **FIG 4. DISTRIBUTION OF IMMUNE-ENRICHED LNCRNA AND MRNA EXPRESSION**
352 **LEVELS IN PRIMARY T-HELPER CELLS**

353 Staggered bars are used to illustrate the proportion of transcripts whose expression
354 falls in bins of 25 million transcripts, or counts, in Roadmap Epigenomics RNA-seq data
355 (A) and FANTOM CAT CAGE data (B), respectively.

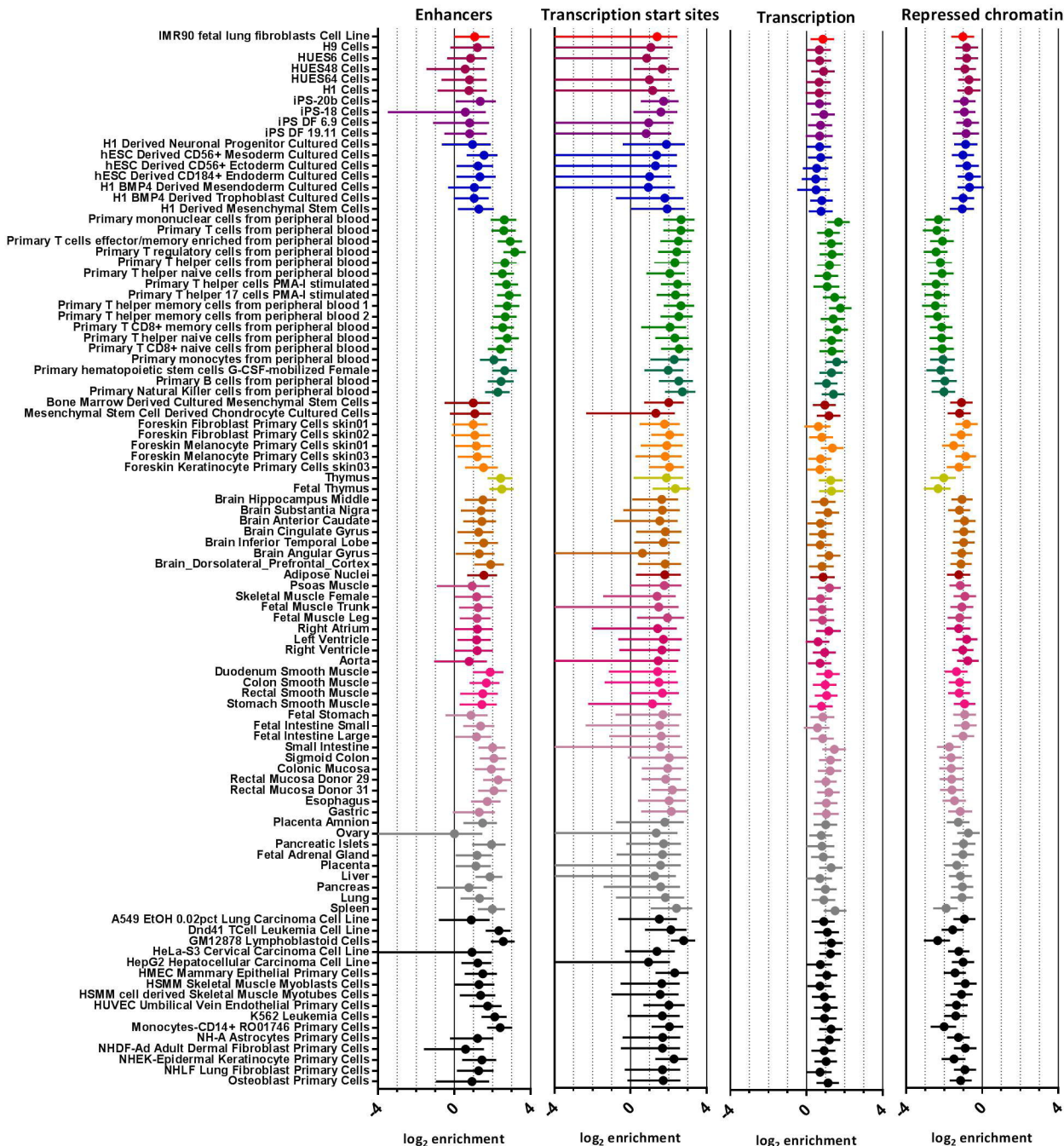
356 **SUPPORTING INFORMATION CAPTIONS**

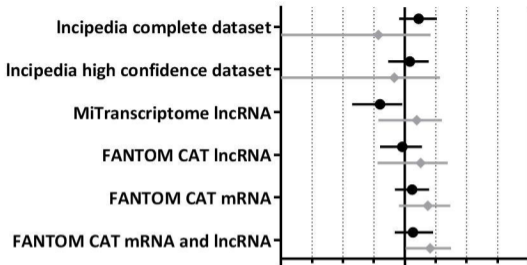
357 **S1 FIG. OF CHROMATIN STATE ANNOTATIONS AMONGST RA SUSCEPTIBILITY**
358 **VARIANTS**

359 Estimates for enrichment of individual states are illustrated for 98 cell types using the
360 Roadmap Epigenomics 18-state model (a). Similar states were grouped into four
361 groups for all immune-relevant primary cell-types (b), as individual states often gave
362 very broad 95% confidence intervals (c). Cell-types are ordered according to the
363 clustering established by the Roadmap Epigenomics project, with chromatin states
364 reordered according to their subsequent grouping. Estimates and confidence intervals
365 are clipped at axis limits, where applicable.

366 **S2 FIG. ENRICHMENT OF IMMUNE-RELEVANT ENHANCERS AMONGST RA**
367 **SUSCEPTIBILITY VARIANTS AFTER CONDITIONING ON IMMUNE-ENRICHED**
368 **TRANSCRIPTS**

369 The influence of FANTOM CAT immune-enriched lncRNA and mRNA was fixed in a
370 probabilistic model of RA susceptibility to confirm the independent enrichment of
371 immune-relevant enhancer chromatin states.



a**b**