# Chromatin three-dimensional interactions mediate genetic effects on gene expression — Source link ↗

Olivier Delaneau, Olivier Delaneau, Marianna Zazhytska, Christelle Borel ...+29 more authors

**Institutions:** Swiss Institute of Bioinformatics, University of Geneva, University of Lausanne, École Polytechnique Fédérale de Lausanne

**Topics:** Chromatin, Regulation of gene expression, Genetic variability, Gene and Genetic variation

Related papers:

- A 3D Map of the Human Genome at Kilobase Resolution Reveals Principles of Chromatin Looping

- Genetic effects on gene expression across human tissues.

- Topological domains in mammalian genomes identified by analysis of chromatin interactions

- Integrative analysis of 111 reference human epigenomes

- Comprehensive mapping of long-range interactions reveals folding principles of the human genome.

# Intra- and inter-chromosomal chromatin interactions mediate genetic effects on regulatory networks

O. Delaneau[123], M. Zazhytska[4], C. Borel[13], C. Howald[123], S. Kumar[56], H. Ongen[123], K. Popadin[47], D. Marbach[86], G. Ambrosini[56], D. Bielser[1], D. Hacker[9], L. Romano-Palumbo[1], P. Ribaux[1], M. Wiederkehr[4], E. Falconnet[1], P. Bucher[56], S. Bergmann[86], S. E. Antonarakis[13], A. Reymond[4*], E. T. Dermitzakis[123*]

[1] Department of Genetic Medicine and Development, University of Geneva, Geneva, Switzerland
[2] Swiss Institute of Bioinformatics (SIB), University of Geneva, Geneva, Switzerland
[3] Institute of Genetics and Genomics in Geneva, University of Geneva, Geneva, Switzerland
[4] Center for Integrative Genomics, Faculty of Biology and Medicine, University of Lausanne, Lausanne, Switzerland
[5] Swiss Institute for Experimental Cancer Research, Ecole Polytechnique Fédérale de Lausanne (EPFL), Lausanne, Switzerland
[6] Swiss Institute of Bioinformatics (SIB), CH-1015 Lausanne, Switzerland
[7] Immanuel Kant Baltic Federal University, Kaliningrad, Russia
[8] Department of Computational Biology, University of Lausanne, Lausanne, Switzerland
[9] Protein Expression Core Facility, School of Life Sciences, Ecole Polytechnique Fédérale de Lausanne (EPFL), Lausanne, Switzerland
* Corresponding authors (alexandre.reymond@unil.ch / emmanouil.dermitzakis@unige.ch )

**Summary:**

Genome-wide studies on the genetic basis of gene expression and the structural properties of chromatin have considerably advanced our understanding of the function of the human genome. However, it remains unclear how structure relates to function and, in this work, we aim at bridging both by assembling a dataset that combines the activity of regulatory elements (e.g. enhancers and promoters), expression of genes and genetic variations of 317 individuals and across two cell types. We show that the regulatory activity is structured within 12,583 Cis Regulatory Domains (CRDs) that are cell type specific and highly reflective of the local (i.e. Topologically Associating Domains) and global (i.e. A/B nuclear compartments) nuclear organization of the chromatin. These CRDs essentially delimit the sets of active regulatory elements involved in the transcription of most genes, thereby capturing complex regulatory networks in which the effects of regulatory variants are propagated and combined to finally mediate expression Quantitative Trait Loci. Overall, our analysis reveals the complexity and specificity of cis and trans regulatory networks and their perturbation by genetic variation.

1

**Introduction:**

A decade of unbiased genome-wide association studies revealed that most of the genetic determinants of complex traits likely reside in non-coding regions of the genome (Maurano et al., 2012; Nica et al., 2010; Nicolae et al., 2010) and are genetic variants modulating gene expression; usually called expression Quantitative Trait Loci or eQTLs (Pickrell et al., 2010). As a consequence, large genome-wide collections of eQTLs across multiple Human populations (Lappalainen et al., 2013), conditions (Quach et al., 2016) and tissues (GTEx Consortium, 2015) have been collected and offer now the resources to systematically identify the genes and tissues involved in complex traits (Ongen et al., 2016). We do not understand precisely the downstream effects of eQTLs on the regulatory machinery of the cell, yet recent developments in high-throughput assays that capture biochemical changes (e.g. ChIP-seq (Johnson et al., 2007)), open regions (e.g. ATAC-seq (Buenrostro et al., 2013)) and conformation (e.g. Hi-C (Lieberman-Aiden et al., 2009)) at the chromatin level provide the tools to achieve this task. This has already given key insights into the molecular mechanisms underlying eQTLs function and the consensus picture that now emerges involves the binding perturbation of transcription factors which result in variable activity of the corresponding regulatory elements (Spitz and Furlong, 2012). These effects then propagate to distal regulatory elements and genes through chromatin interactions (Grubert et al., 2015; Kasowski et al., 2013; Kilpinen et al., 2013; McVicker et al., 2013; Waszak et al., 2015) whose possible range is delimited by chromatin contact domains (usually called TADs for Topologically Associating Domains) which represent the structural units of the three-dimensional nuclear organization (Lieberman-Aiden et al., 2009; Pombo and Dillon, 2015; Rao et al., 2014).

To further characterize gene regulatory networks and their perturbation by genetic variants, we extended our previous work (Waszak et al., 2015) and assayed a combination of three well-studied histone modifications (H3K27ac, H3K4me1 and H3K4me3), which are known to reliably capture enhancer and promoter activities, in 317 lymphoblastoid cell lines (LCLs) and 78 primary fibroblasts derived from unrelated European individuals (**methods 1-2**). We complemented this dataset by profiling gene expression and genotyping millions of variants in each sample. By leveraging the population-scale design of this study, we explore the interplay between genetic variants, regulatory activity and gene expression, and show how this captures interactions between regulatory elements, reflects the three-dimensional

67    nuclear organization, varies across two cell types and propagates various forms of genetic

68    effects on gene expression.

69

70    **Results:**

71

72        *1. Inter individual variability of chromatin activity reveals Cis Regulatory Domains*

73    As an initial step to comprehend the inter-individual variability of chromatin activity, we first

74    integrated all chromatin assays into two matrices comprising 271,417 chromatin peaks

75    quantified in 317 LCLs and 78 fibroblast lines from unrelated European individuals,

76    respectively (**supplementary figures 1-2**) and then proceeded with chromatin Quantitative

77    Trait Locus mapping (cQTL). To properly deconvolve biological from technical variability and

78    therefore to maximize the power of cQTL discovery, we repeated the procedure multiple

79    times across variable sets of covariates (**supplementary figures 3-4**; **methods 4-6**). Using the

80    optimal configuration, we discovered 44,492 and 14,476 cQTLs for LCL and fibroblast,

81    respectively, which constitutes to our knowledge the largest available collection of cQTLs. As

82    expected, the variability of a substantial fraction of chromatin peaks results from nearby

83    genetic variants (pi1=30.4%; see **supplementary methods D6** for a description of pi1) often

84    located within open chromatin regions and exhibiting small allelic biases (i.e. same

85    proportions of negative and positive effect sizes for variants located within chromatin

86    peaks; **supplementary figures 5A-C**). Interestingly, the chromatin peaks under genetic

87    control tend to cluster together (**supplementary figure 5D**) possibly revealing the

88    coordination between nearby regulatory elements (Kilpinen et al., 2013; Waszak et al.,

89    2015).

90

91    To further characterize this coordinated behavior, we systematically measured in LCLs the

92    inter-individual correlation between nearby chromatin peaks (within 1Mb of each other)

93    and found a strong signal of correlation (pi1=18.7%; **methods 7**) that rapidly decays with

94    distance, slightly varies across ChIP-seq assay pairs, is driven by the most variable chromatin

95    peaks (**supplementary figures 6A-C**) and importantly forms well-delimited domains that we

96    call Cis Regulatory Domains (CRDs; **figure 1A-C**). We then looked at the fibroblast data to

97    infer the cell type specificity of the significant correlations and found that long range

98    correlations (e.g. > 100kb) are less shared between cell types than short range ones (e.g. <

3

99    100kb; **supplementary figure6D**); suggesting that they participate more to the regulatory

100   program responsible for cell type specificity.

101

102   We then searched for specific DNA properties that could underlie this correlation structure.

103   To this aim, we first used publicly available deep Hi-C data for LCLs (Rao et al., 2014) to

104   estimate the contact intensity between any pair of chromatin peaks (**methods 6**) and

105   compared the resulting contact intensity measures to the correlation values. Overall, we

106   find extremely good concordance between both measures (**figure 1D**). Then, we used a

107   collection of protein binding sites derived for LCL as part of the ENCODE project (Encode

108   Project Consortium, 2012) (**methods 3**) and found that the correlation decay is slower

109   between binding sites of structural proteins (CTCF, RAD21, SMC3 and ZNF143;

110   **supplementary figure6E**). Altogether, this shows that this correlation structure sufficiently

111   captures the local three-dimensional conformation of chromatin.

112

113   A genome-wide call set of CRDs effectively delimits the proximal regulatory elements that

114   likely interact together and therefore would help to decompose the genome into

115   independent cis regulatory units. To assemble them, we designed a calling algorithm based

116   on hierarchical clustering: chromatin peaks are iteratively grouped together into clusters

117   based on the correlation levels they exhibit and the resulting binary tree is then cut at the

118   level maximizing the total correlation mass captured (**methods 8**). This allowed us to

119   discover 12,583 and 10,442 CRDs for LCL and fibroblast, respectively. In LCL, the CRDs are

120   32kb long on average, cover ~8% of the genome and span point associations between two

121   distinct regulatory regions to complex association patterns that can encompass up to 80

122   distinct regulatory regions (**supplementary figures 7A-B**). By data subsampling, we assessed

123   the discovery potential of this approach and find that 150 samples are enough to saturate

124   the number of CRDs discovered while more samples are needed to better delimit the CRD

125   content and boundaries (**supplementary figures 7C-D**). In terms of chromatin peak

126   composition, we found that enhancers are the main component of CRDs (59.4% of

127   H3K4me1 peaks), that they are uniformly distributed along them while promoters tend to

128   locate at their boundaries (H3K4me3; **figure 1E**; **supplementary figure 7E**). This latter

129   observation suggests that most of the CRDs we detected reflect how enhancers are brought

130   in close proximity of promoters by chromatin looping.

4

131

### 2. CRDs delimit functionally active and variable TAD subdomains

133 Given the high concordance between Hi-C contacts and correlation data, we compared CRDs

134 to a collection of TADs derived from a deep LCL Hi-C experiment that constitutes the most

135 detailed definition of TADs so far (**methods 3**). CRDs are substantially smaller than TADs

136 (32kb versus 185kb; **figure 1F**) even though they can reach the same size in some cases (up

137 to 1Mb). When comparing the genomic regions defined by both CRDs and TADs, we find

138 that 90% of the TADs overlap at least one CRD, 57% at least two and that CRDs tend to be

139 fully encompassed within TAD boundaries (**supplementary figures 8A-B**). In addition, we

140 also observed a higher density of CTCF binding sites at CRD boundaries similar to what is

141 known for TADs (**supplementary figure 8C**) (Rao et al., 2014). All this constitutes

142 accumulating evidences that CRDs actually correspond to TAD subdomains.

143

144 We then aimed at explaining why some of these TAD sub-domains correspond to CRDs while

145 some others do not. To do so, we categorized the chromatin peaks into four categories

146 depending on their location relative to both TADs and CRDs. First, this showed that

147 chromatin peaks within CRDs are very likely to also fall within TADs (OR=1.75, fisher p-value

148 < 1e-300); confirming the similarity between chromatin activity and structure. When

149 focusing exclusively on TAD specific chromatin peaks and looking at their activity across all

150 individuals, we find that those that are also members of CRDs exhibit much higher mean

151 activity and variability than those that are not (**figure 1G**; **supplementary figure 8D**).

152 Together, this suggests that TADs mix together inactive (potential) and active (realized)

153 subdomains and we propose that active ones correspond to CRDs.

154

### 3. Inter-chromosomal correlation between CRDs reflects nuclear organization of DNA

156 By using fluorescence microscopy or Hi-C assays, multiple studies have characterized a high

157 level organization of the chromatin; within chromosome territories (Meaburn and Misteli,

158 2007), A/B nuclear compartments (Lieberman-Aiden et al., 2009) and sub-compartments

159 (Rao et al., 2014). Here, we used population variability of chromatin activity as perturbation

160 to uncover the interactions underlying the global chromatin architecture. Specifically, we

161 measured inter-individual correlations for all pairs of chromatin peaks located on distinct

162 chromosomes (n=~3.4e10 tests) and found a substantial signal of association (pi1=8.1%;

163    estimated by data sub-sampling). Given the difficulty to properly control for FDR in this

164    massive multiple-testing setting, we first examined the data by considering as significant any

165    pair of chromatin peaks with an association P-value below 1e-6. By doing so, we observed

166    that some chromosomes preferentially associate with some others (**figure 2A**); notably

167    chromosomes 16, 17, 19, 20 and 22 in line with the observation that these chromosomes

168    tend to locate at the center of the nucleus thereby increasing their chance to interact with

169    other chromosomes (Kaufmann et al., 2015). We also found that the per-chromosome

170    association frequency (i.e. the fraction of tests per-chromosome with a p-value < 1e-6)

171    correlates well with gene density (**figure 2E**) suggesting that gene rich chromosomes

172    participate more to inter-chromosomal chromatin interactions (Kalhor et al., 2011).

173    Consistent with this hypothesis, chromosomes 16, 17, 19, 20 and 22, are the richest in CpG

174    islands (Lander et al., 2001). By visual inspection of the correlation structure for some pairs

175    of chromosomes, we observed that strong correlations often concentrate at specific

176    genomic locations that frequently involve CRDs (**figure 2B-C**). To validate these 'hot spots' of

177    correlation, we again used deep Hi-C data for LCL and found good concordance between the

178    correlation structure and the inter-chromosomal contact intensities both locally

179    (**supplementary figure 9**) and globally (**figure 2D**). Altogether, this suggests a global

180    correlation structure between chromatin peaks above CRDs that reflects how chromatin is

181    packaged into the nucleus.

182

183    To move from this global pattern to individual associations between regulatory elements

184    while properly controlling for FDR, we next performed inter-chromosomal correlation

185    analysis between CRDs. This requires quantifying CRD activity on a per-individual basis, a

186    task that we performed using dimensionality reduction techniques to aggregate together

187    the quantifications of multiple peaks of a CRD into a single quantification vector (**methods**

188    **8**). As a result, we reduced the number of tests by three orders of magnitude (from

189    n=~3.4e10 tests to n=~7.5e7) and substantially increased the association signal (from

190    pi1=8.1% between chromatin peaks to pi1=16.9% between CRDs; **supplementary figure**

191    **10A**). Thus, CRDs are commonly associated across chromosomes as illustrated by the 25,315

192    significant CRD-CRD associations detected at 1% FDR. Interestingly, bringing all these

193    associations together into a network reveals the presence of hubs in which few CRDs play a

194    central role and connect together multiple CRDs from distinct chromosomes (**figure 2F**;

6

195   **supplementary figures 10B-C**). This led us to further refine our understanding of the global

196   correlation structure of chromatin activity by annotating each CRD as belonging to one of

197   the previously characterized nuclear compartments (A1, A2, B1, B2 or B3; **methods 3**) (Rao

198   et al., 2014). From this, we made three observations: (i) CRDs from the same chromosomal

199   compartment are more likely to be correlated, (ii) hubs usually bring together CRDs from the

200   same chromosomal compartment and (iii) associations between CRDs from the same

201   chromosomal compartment are more shared across cell types (**figure 2F**; **supplementary**

202   **figures 10D-E**). All this shows that inter-chromosomal associations between CRDs are

203   abundant in the genome, highly cell-type specific and form hubs with respect to how

204   chromatin segregates into A/B nuclear compartments.

205

206   ### 4.   *CRDs are fundamental functional units for gene regulation*

207   The population scale design of this study offers a unique opportunity of building a genome-

208   wide map interconnecting genes and CRDs at a resolution and a level of accuracy

209   inaccessible to simple approaches based on co-localization or genomic distance. Indeed, we

210   could directly assign genes to CRDs by assessing the significance of the correlation between

211   their respective quantifications (**methods 9**). In practice, this approach allowed us to unravel

212   a massive signal of associations: in LCL, the vast majority of the genes (pi1=82.4%) are

213   associated with at least one CRD in cis (+/- 1Mb) which converts into 12,204 CRD-gene

214   associations at 1% FDR (**figure 3A**; **supplementary figure 11A**). Often, these associations

215   show relatively high connectivity degree; in 23.3% of the cases a gene is associated with

216   multiple CRDs, while in 20.2% of the cases a CRD is associated with multiple genes (**figure**

217   **3B**). In terms of the linear position in the chromosome, we found that 46.8% of the

218   associations connect genes to CRDs that encompassed their TSS (Transcription Start Sites),

219   which also means that less than half of the associations are detectable through co-

220   localization. Overall, the TSS of the associated genes tends to locate close to the CRD

221   boundaries and usually exhibit a strong association signal when it falls within CRDs (**figure**

222   **3C**; **supplementary figure 11B**). In terms of histone mark composition, CRDs associated with

223   genes are enriched for active marks (e.g. H3K27ac, H3K4me1 and H3K9ac), depleted for

224   repressive (e.g. H3K27me3; **supplementary figures 11C-D**) and more often belong to the

225   active nuclear compartment A (**figure 3D**). Altogether, this suggests that these associations

226  bridge active CRDs to their target genes, thereby revealing local regulatory networks

227  involved in the regulation of a vast majority of the genes.

228

229  The observation that genes and CRDs frequently show high connectivity and do not overlap

230  led us to perform additional analyses. Notably, we first found that nearby genes linked to a

231  single CRD are much more correlated with each other than those that are linked to distinct

232  CRDs. Similarly, CRDs associated with the same gene are also highly correlated with each

233  other (**figures 3E, supplementary figure 12A**). Distant gene-CRD associations (separated by

234  from 10kb to 1Mb) occur within TAD boundaries more likely than what is expected by

235  chance and exhibit more cell type specificity than short range ones (**supplementary figures**

236  **12B-C**). There is therefore another layer of coordination between nearby CRDs that is cell

237  type specific and results in co-expression between nearby genes at the transcriptional level.

238  This also illustrates the limitations of using domain representations alone to capture the

239  hierarchical nature of regulatory interactions as additional layers of interactions are

240  systematically missed.

241

242  ### 5.  *CRDs are central mediators of genetic effects on gene expression*

243  The relatively large sample size of the LCL data set (n=317) offers the required statistical

244  power to analyze the genetic effects on CRDs . We therefore searched for Quantitative Trait

245  Loci (QTLs) for both genes and CRDs. Concerning CRDs, we did not only quantify their

246  activity as described before but also their structure by quantifying for each individual

247  whether its data contributed positively or negatively to the overall correlation structure

248  (**methods 8-9**). To summarize, we tested three molecular phenotypes for association with

249  nearby genetic variants (**methods 4**), thereby resulting in three different kinds of QTL

250  discoveries: at the gene expression level (eQTL), at the CRD activity level (aCRD-QTL) and at

251  the CRD structure level (sCRD-QTL). In total, we discovered 6,157 aCRD-QTLs (pi1=68.5%),

252  7,658 eQTLs (pi1=59.2%) and 110 sCRD-QTLs (pi1=11.3%; **figures 4A-B**; **supplementary**

253  **figure 13A**), which demonstrates that, similarly to genes, CRDs are under very strong genetic

254  control. In terms of genomic distribution, all three types of QTLs exhibit the same pattern:

255  they tend to become stronger and more frequent as we get closer to the genomic source of

256  the target phenotype (**supplementary figure 13B**). However, in terms of minor allele

257  frequency (MAF) spectrum, they show some differences: sCRD-QTLs involve almost

8

258   exclusively genetic variants with a MAF below 20% while both eQTLs and aCRD-QTLs span

259   the full MAF range (**figure 4C**); suggesting that the sCRD-QTLs are under intense selective

260   pressure. This is also supported by the observation that sCRD-QTLs involve more often a loss

261   of protein binding motifs compared to aCRD-QTLs and eQTLs (**supplementary methods D13**;

262   **supplementary figure13C**). To summarize, this suggests two kinds of genetic variations

263   affecting CRDs; the frequent case involving common genetic variants modulating the CRD

264   activity through perturbation of protein binding affinity and the more deleterious case

265   involving rare genetic variants that change the CRD structure by disrupting internal

266   interactions.

267

268   As a natural step forward, we then assessed the extent to which these molecular QTLs

269   overlap. In the context of our study, this task is greatly facilitated by the availability of a

270   dense genome-wide map interconnecting genes and CRDs, which, in practice, allows testing

271   for association QTLs of CRDs with relevant genes. Overall, we found extensive overlaps

272   between all types of molecular QTLs (**figure 4D**). More specifically, we first observed that

273   sCRD-QTLs are almost fully recapitulated by aCRD-QTLs (pi1=94.5%); indicating that

274   perturbations of the CRD structure also result in activity variations. We then found that

275   pi1=74.4% of the aCRD-QTLs affect gene expression, while pi1=74.7% of the eQTLs affect

276   CRD activity. Thus, most of the molecular QTLs affect both CRD activity and gene expression

277   and therefore likely tag the same causal variants. To identify the slight differences between

278   both QTL sets, we compared their respective chances of overlapping functional annotations

279   of the genomes (Encode Project Consortium, 2012; Whyte et al., 2013; Zarrei et al., 2015).

280   The aCRD-QTLs locate more frequently in open chromatin regions, protein binding sites and

281   enhancers, while eQTLs are enriched in promoters and transcribed regions (**supplementary

282   figure13D**). Both types of QTL exhibit therefore different trade-offs in capturing

283   distal/proximal genetic effects on gene expression.

284

285   We next examined the role of CRDs in mediating genetic effects on gene expression. First,

286   we tried to look at how multiple eQTLs with independent effects on a given gene actually

287   decompose into individual effects at the CRD level. We performed conditional analyses to

288   supplement our previous QTL discoveries with additional independent ones (**supplementary

289   methods D14**) and then compared the numbers of QTLs per gene and CRD by leveraging the

290    gene-CRD associations. We have observed that genes with multiple eQTLs associate more

291    often with (a) a single CRD having multiple aCRD-QTLs (OR=1.81) or (b) with more than 1

292    CRDs (OR=1.42) compared to genes with a single eQTL (**figure 4E**). This demonstrates that

293    the joint effect of multiple eQTLs on a gene results from molecular interactions at two

294    different levels: between regulatory elements within a single CRD or between distinct CRDs.

295    Then, we asked how genetic effects actually flow through genes and CRDs in order to

296    understand the causal relationships. A prerequisite of our analysis is to focus exclusively on

297    genetic variants that affect both molecular phenotypes simultaneously.  To detect these, we

298    applied a Principal Component Analysis on each of the 12,204 gene-CRD associations and

299    used the individual coordinates on PC1 as quantifications for QTL mapping. In total, we

300    discovered 8,706 QTLs for gene-CRD pairs (called eCRD-QTLs; CRD-expression QTLs) and

301    estimated that pi1=81.7% of these pairs mediate genetic effects. We then used Bayesian

302    Network modeling to unravel causal relationships, an approach that has already been used

303    for causal inference between genetic variants and molecular phenotypes (Gutierrez-Arcelus

304    et al., 2013) (**methods 10**). In our study, this provided very high call rates in most of the

305    cases (**supplementary figure 14A**); meaning that we were able to reliably assign a causal

306    scenario for a vast majority of the 8,706 gene-CRD-QTL triplets under consideration. Overall,

307    we find a clear signal that activity changes at CRDs are causal to changes in gene expression

308    when they originate from eCRD-QTLs distal to genes while they are reactive for proximal

309    eCRD-QTLs (**figure 4F**). Interestingly, this signal becomes even stronger when we only

310    focused on eCRD-QTLs directly located within CRDs (**supplementary figure 14B**). This led us

311    next to distinguish the functional elements that contribute to these two types of causal

312    mechanisms (i.e. causal or reactive). We found that eCRD-QTLs involving a causal role for

313    CRDs usually map in enhancers while those that involve a reactive role for CRDs map in

314    promoters (**supplementary figure 14C**). Altogether, this shows that regulatory interactions

315    within CRDs mediate the long range effects of enhancer related regulatory variations on

316    gene expression while they actually reflect these perturbations in the case of promoter

317    related regulatory variations.

318

319    The availability of four collections of molecular QTLs allows further functional comparisons.

320    First, we can determine the degree to which those QTLs are shared across cell types. To do

321    so, we quantified all the LCL features (CRDs, gene-CRD pairs and genes) in the 78 fibroblast

10

322  samples and computed the P-values of association in fibroblasts of all significant LCL

323  discoveries. We found that between pi1=27.3% and pi1=35.4% of the QTLs for CRDs are also

324  significant in fibroblasts, a value which goes up to pi1=48% for eQTLs (**supplementary figure**

325  **14D**). Thus, QTLs for CRDs are more cell type specific than eQTLs similar to the previous

326  observation that QTLs for CRDs capture more enhancer related effects. Second, we

327  examined if our QTL collections are enriched for known GWAS hits (MacArthur et al., 2017)

328  as it has already been shown for eQTLs (Maurano et al., 2012; Nica et al., 2010; Nicolae et

329  al., 2010). We performed an enrichment analysis (**supplementary methods 17**) and found

330  similar enrichment values across all types of QTLs, emphasizing their relevance in studying

331  diseases, with a slight advantage for sCRD-QTLs as a likely result of their predicted

332  deleterious effects (**supplementary figure 14E**).

333

334  We provide here two concrete examples how the dense genome-wide map interconnecting

335  genes, CRDs and regulatory variations that we progressively assembled throughout this

336  study could empower future association studies. As a first example, we used CRDs as a

337  genome annotation to study the cumulative effect of multiple genetic variants. Specifically,

338  we used whole genome and transcriptome sequencing data from the GEUVADIS study

339  (**methods 3**) (Lappalainen et al., 2013) to assess if the accumulation of rare variants within a

340  CRD is associated with the expression of nearby genes (i.e. rare-eQTL effects). To this aim, a

341  'genetic burden' was defined for each CRD by counting the number of minor alleles at rare

342  variants (MAF<5%) falling into its regulatory regions and was then tested for association

343  with the expression levels of nearby genes (**methods 11**). This amounts to a burden test at

344  the regulatory level and showed that a non-negligible fraction of genes (pi1=10.4%) are

345  affected by rare-eQTLs (**figures 5A-B**). At 5% FDR, this resulted in 33 CRDs whose baseline

346  activity is perturbed by the accumulation of rare mutations with a level of expression of the

347  downstream gene being decreased in 64% of the cases (**supplementary figure 15A**). This

348  demonstrates the potential of our approach to explore the cumulative effects of multiple

349  regulatory variations such as haplotypic effects. As a second example, we leveraged the

350  inter-chromosomal associations between CRDs to discover trans-eQTLs likely acting through

351  direct chromatin interactions and not diffuse factors (Bryois et al., 2014; Pierce et al., 2014).

352  We implemented this by combining three different layers of associations: (i) aCRD-QTLs, (ii)

353  CRD-CRD inter-chromosomal associations and (iii) gene-CRD associations. This defined

354 20,489 variant-gene pairs located on distinct chromosomes to be tested for association

355 (**methods 12**), a task that we performed in six distinct datasets compiling data from two

356 studies (the present study and EUROBATS; **methods 3**) and across five cell types (blood, fat,

357 LCL, skin and fibroblasts). Overall, we found association enrichments at particularly high

358 levels in the LCL datasets (**figure 5C**) potentially highlighting the cell type specificity of trans

359 effects. We then extracted the 27 hits significant at 5% FDR in the largest LCL data set

360 (EUROBATS LCL, n=765) and replicated 9 of them (p-value < 0.01 and same effect size

361 direction) in our LCL dataset (n=317; **figure 5D**; **supplementary figure15B**). Almost none of

362 these 27 hits are replicated in other cell types supporting the notion that they are real trans

363 effects and not due to technical biases such as read mapping artifacts. By assembling all the

364 relevant associations between CRDs, genes and trans-QTLs into networks, we find that all

365 associations actually relate to two distinct variants (rs2980236 and rs7758352); the former

366 variant affecting a total of 8 genes in trans, all of them connected by a CRD hub. Then, we

367 enforced the directionality of the edges in the networks to match our mechanistic

368 hypothesis that trans-eQTLs act through chromatin interactions and assessed whether this

369 assumption is supported by the data using Bayesian networks (**supplementary methods 20**).

370 Overall, the data provided extremely good support (**figure 5E**), thereby indicating that these

371 9 trans-eQTLs act through direct inter-chromosomal chromatin interactions.

372

373 **Discussion:**

374 In this study, we have analyzed genome-wide genotype, chromatin and gene expression

375 data to reveal first- and higher-order regulatory interactions in the genome by using genetic

376 variation as a perturbation. We illustrate the multiple benefits of integrating genome

377 variation of multiple molecular assays to dissect the effects of regulatory variation onto

378 gene expression at an unprecedented scale and resolution. We provide a very large set of

379 chromatin QTLs that can form the basis of many specific analyses at the gene or region level

380 in order to reveal genetic effects and interactions. We uncover the extensive interplay and

381 organization of chromatin around genes and regulatory elements that ultimately form Cis

382 Regulatory Domains (CRDs). These CRDs, while being embedded in Topologically Associating

383 Domains (TADs), have higher complexity and degree of interaction at the sub-TAD level,

384 capture local interactions with most genes and tend to be highly tissue specific. In addition,

385 they are key factors involved in the co-regulation between nearby genes. However, they can

386    also form a backbone for long-range direct interactions that result in regulation and co-
387    regulation of genes at larger distances and on distinct chromosomes, revealing a level of
388    organization that was previously suspected (Williams et al., 2010) but not documented
389    beyond a few rare examples. To summarize, we find that genome variability in the context
390    of regulatory activity and gene expression reflects the local and global nuclear organization
391    of chromatin, thereby revealing complex networks of chromatin interactions. These
392    regulatory networks are perturbed by genetic variations and serve as physical backbones to
393    propagate and combine their effects on gene expression which can result in distal eQTL,
394    independent eQTL, trans eQTL and rare eQTL effects. Overall, our work bridges discoveries
395    made by genome-wide studies focusing on the structural properties of the chromatin
396    together with those describing the genetic basis of chromatin activity and gene expression
397    variability.

398    By drawing a dense and high resolution genome-wide map of the interplay between genetic
399    variations, regulatory elements and gene expression across two cell types, we provide a
400    unique reference resource to better interpret disease associated variants, a proof-of-
401    principle that population-scale variability of molecular phenotypes is a powerful means to
402    unravel the regulatory networks underlying gene expression and that these regulatory
403    networks constitute informative priors to enhance future association studies. We also
404    introduced major methodological advances related to the processing of multi-layered
405    molecular phenotype data, notably for their integration through clustering and
406    dimensionality reduction techniques and for causal inference through Bayesian networks.
407    Future studies using such a population-scale approach should allow building on our work by
408    including larger samples to capture even weaker effects, more cell types to ultimately build
409    a cross-tissue catalog of regulatory networks, affected patients to characterize disease
410    specific regulatory signature and other molecular assays to either rationalize the
411    examination of active regulatory elements (e.g. ATAC-seq) or extend it to silencing effects
412    (e.g. ChIP-seq for H3K27me3).

413    The consideration of genetic effects in local and global scale, as we have performed in this
414    study, offers mechanistic insights that allow us to start assessing the true cellular impact of
415    genetic variants individually as altogether and their impact not only at the gene level but
416    also at the cellular level. This brings us closer to fully dissecting the properties and

417 interactions of genetic variants defining organismal phenotypes via intra- and inter-cellular

418 interactions.

419

420 **Author contributions**

421 E.T.D., A.R., S.E.A., S.B. and P.B. designed and supervised the study and contributed to data

422 interpretation. S.E.A. and C.B. provided the GenCord samples. D.H. cultured and processed

423 the lymphoblastoid cell lines. C.B., P.R and E.F. cultured and processed the fibroblasts. A.R.

424 designed the ChIP experiments. M.Z. and M.W. executed the ChIP experiments. L.R-P. and

425 D.B. prepared ChIP- and RNA-seq libraries and executed the sequencing. O.D. designed and

426 executed the primary data analysis. C.H., G.A., H.O., K.P. and S.K contributed to data

427 analysis. O.D. and E.T.D. performed the primary manuscript writing. A.R., S.E.A., S.B., D.M.,

428 C.B. and C.H. contributed to the manuscript writing.

429

430 **Data availability**

431 All the raw sequence data (i.e. BAM files) generated in this study are deposited in the Array

432 Express Archive (http://www.ebi.ac.uk/arrayexpress/ *[Available upon publication]*). All the

433 processed data (molecular phenotype matrixes, collections of molecular QTLs and CRDs) are

434 on a dedicated webpage (ftp://jungle.unige.ch/SGX *[Nicer website soon available]*). All the

435 code implementing most computational methods used in this study has been packaged,

436 documented and released on GitHub (https://github.com/odelaneau/clomics *[Better*

437 *documentation soon available]*).

438

439 **Acknowledgments**

445

446 **Methods**

447 Hereafter is given a summary of all the experimental and computation methods used

448 throughout this study. More details can be found in supplementary information and an

14

449   overview of how all computational methods are interconnected is given in **supplementary**
450   **figure 16**.

451

452   *1.  Sample collection (supplementary methods A)*
453   We collected samples from three different sources: (i) 46 LCLs that we generated as part of
454   a previous study (Waszak et al., 2015), (ii) 111 LCLs from the 1000 Genomes Project (1000
455   Genomes Project Consortium, 2015) and (iii) 160 LCLs and 78 fibroblasts from the GenCord
456   collection (Dimas et al., 2009). Altogether, this gave a total of 317 LCLs and 78 fibroblasts all
457   from individuals of European ancestry.

458

459   *2.  Experimental methods (supplementary methods B)*
460   We cultured all cells for which no data was already available (271 LCLs and 78 fibroblasts) to
461   perform both ChIP-seq and RNA-seq. Specifically, we performed chromatin immuno-
462   precipitation for three different histone modifications: H3K27ac, H3K4me1 and H3K4me3
463   and extracted total RNA. We sequenced all resulting libraries on a HiSeq2000 machine. In
464   addition, we also collected genotype data for all individuals by merging existing sequence
465   data from the 1000 Genomes project (n=145) with data we generated using Illumina Human
466   OMNI 2.5M SNP arrays (n=178).

467

468   *3.  External data (supplementary method C)*
469   To analyze the data we generated, we used multiple external data sets: (i) contact maps,
470   TAD definition and A/B nuclear compartments all derived for LCL from a deep Hi-C
471   experiment (Rao et al., 2014), (ii) multiple functional annotations of the genome (binding
472   sites for 50 transcription factors, 11 histone modifications, open chromatin regions and
473   genome segmentation) derived for LCLs in the frame of the ENCODE project (Encode Project
474   Consortium, 2012) and (iii) RNA-seq and genotype data for hundreds of European
475   individuals from the GEUVADIS and EUROBATS projects (Buil et al., 2015; Lappalainen et al.,
476   2013).

477

478   *4.  Genotype data preparation (supplementary methods D1)*
479   All genotype data coming from SNP array have been merged together and filtered using
480   standard procedures to remove low quality SNPs. The resulting genotype matrix was

481 imputed from the 1000 Genomes phase 3 reference panel, poorly imputed variants
482 removed and the rest combined with the genotype data coming from the 1000 Genomes
483 sequencing data (1000 Genomes Project Consortium, 2015). This finally gave us genotype
484 data for 323 individuals at 9,255,024 variants.

485

486 ***5. Molecular phenotype data preparation (supplementary methods D2-3)***
487 All sequence data were mapped onto the human genome (hg19) using either BWA (Li and
488 Durbin, 2009) for ChIP-seq data or GEM (Marco-Sola et al., 2012) for RNA-seq data. Gene
489 expression has been quantified using QTLtools (Delaneau et al., 2017) with GENCODE v19
490 (Harrow et al., 2012) as reference and filtered to retain genes expressed in more than 90%
491 of the samples. For the chromatin assays, we proceeded in two steps: we built a population
492 scale BAM file for each assay by sub-sampling 1 million reads from 50 LCLs and 50
493 fibroblasts, then we performed peak calling on these using Homer v4.7 (Heinz et al., 2010)
494 and quantified each sample in turn using the resulting peaks coordinates as reference
495 genome annotations. We finally residualized all these molecular phenotypes for multiple
496 covariates: sex, project of origin, genotyping platform, ancestry and technical variables
497 (done by maximizing QTL discovery; **methods 6**). Of note, we merged the outcome of all
498 chromatin assays together to obtain a total of four quantification matrices: (1) expression at
499 18,939 genes for 317 LCLs, (2) activity at 271,467 chromatin peaks for 317 LCLs, (3)
500 expression at 18,068 genes for 78 fibroblasts and (4) activity at 271,467 chromatin peaks for
501 78 fibroblasts.

502

503 ***6. Molecular QTL mapping (supplementary methods D4)***
504 We mapped molecular QTL for each molecular phenotype using (i) 1,000 permutations to
505 correct for the number of genetic variants being tested in cis (+/- 1 Mb from the genomic
506 feature boundaries) and (ii) the False Discovery Rate procedure implemented in the
507 R/qvalue package (Storey and Tibshirani, 2003) to correct for the number of phenotypes
508 being tested genome-wide. All these analyses were run using the QTLtools software package
509 (Delaneau et al., 2017).

510

511 ***7. Building correlation and contact maps (supplementary methods D5-7)***

512    We built correlation maps for the chromatin assays by systematically measuring inter-

513    individual correlation (i.e. Pearson correlation coefficient) between any pair of chromatin

514    peaks located on the same chromosome (intra-chromosomal) or on distinct chromosomes

515    (inter-chromosomal). In multiple analyses, we measured the enrichment of small correlation

516    P-values according to various parameters of interest by using the pi1 statistics of the

517    R/qvalue package (Storey and Tibshirani, 2003). This pi1 statistics essentially gives an

518    estimate of the number of statistical tests made under the alternative hypothesis of

519    association. In addition, we built peak-centered Hi-C contact maps by interpolating the

520    contact intensities between any pair of chromatin peaks from a deep Hi-C experiment made

521    for LCL (Rao et al., 2014). Overall, this means that between any two chromatin peaks, we got

522    two interaction measures: (a) the correlation between activity levels of the peaks and (b)

523    the contact intensity measured by Hi-C.

524

525    ### 8.  Calling and quantifying CRDs (supplementary methods D8-10)

526    We called Cis Regulatory Domains (CRDs) from correlation data in two steps: (i) we

527    performed hierarchical clustering on the chromatin data to get a binary tree for each

528    chromosome that regroups the chromatin peaks depending on the correlation they exhibit

529    and (ii) we cut the resulting binary tree at the level maximizing the total correlation mass

530    captured by using three empirical criteria (i.e. overall, edge and distal correlations) which

531    gave us clusters of highly correlated chromatin peaks (i.e. the actual CRDs). Once the CRDs

532    were called, we quantified them on a per sample basis at two levels: (i) we measured their

533    overall activity by simply taking the averaged peak quantification and (ii) we measured their

534    structure by assessing the contribution of each individual onto the mean correlation within

535    the CRD.

536

537    ### 9.  Mapping genes and QTLs for CRDs (supplementary methods D11-12)

538    Once the CRDs quantified, we searched for associations with nearby genes and genetic

539    variants using the permutation procedure implemented in the QTLtools software package

540    (Delaneau et al., 2017). This allowed us to efficiently correct for the fact that multiple

541    genetic variants, genes and CRDs are tested genome wide. Of note, we only searched for

542    associations within 1Mb of the CRD boundaries. For the genes, we used as genomic

543    coordinates the position of their TSSs and only considered protein coding genes or lincRNAs.

17

544    In some cases, we complemented this step by conditional analysis as implemented in the
545    forward-backward scan of QTLtools (Delaneau et al., 2017).

546

547        **10. Inferring causality (supplementary methods D15-16)**
548    To infer the causal relationships between QTLs, CRDs and genes, we proceeded in two steps.
549    First, we only focused on QTLs affecting both CRDs and genes. To get those, we quantified
550    each gene-CRD pair using PCA and mapped QTLs using the coordinates on PC1 as
551    quantifications. As a byproduct, this gave us QTL-CRD-gene triplets onto which we
552    performed Bayesian Network (BN) analysis to get the posterior probabilities for the three
553    possible causal models we considered in this study: (i) the causal model with QTL > CRD >
554    gene, (ii) the reactive model QTL > gene > CRD and (iii) the independent model with QTL >
555    CRD and QTL > gene. This has been implemented using the R/bnlearn package (Scutari,
556    2010).

557

558        **11. CRD based burden test (supplementary methods D18)**
559    We designed a burden test at the regulatory level by leveraging CRDs. This test relies on two
560    steps: (i) we measured the burden of rare mutations (minor allele at variant with a MAF <
561    5%) for each CRD and (ii) we then tested this burden with expression at relevant genes. By
562    relevant genes, we mean genes that have been previously found to be associated with the
563    CRD of interest. We performed this analysis on a dataset in which we have whole genome
564    and transcriptome data for 358 individuals (1000 Genomes Project Consortium, 2015;
565    Lappalainen et al., 2013).

566

567        **12. CRD driven trans-eQTL mapping (supplementary methods D19-20)**
568    We built a map of gene-variant pairs located on distinct chromosome to be tested for
569    association by leveraging the inter-chromosomal CRD-CRD associations we discovered in the
570    frame of this study. We assembled together (i) gene-CRD associations (cis links), (ii) CRD-
571    CRD associations (trans links) and (iii) QTL-CRD associations (cis links) which effectively
572    forms 20,489 gene-variant pairs to be tested in *trans*. Then, instead of performing the
573    billions of tests required by standard trans-eQTL analysis (millions of variants times
574    thousands of genes), we only focused on the 20,489 pairs this approach helped us to

575    assemble. We performed these tests in multiple data sets (e.g. EUROBATS (Buil et al., 2015))

576    and extracted the significant hits we obtained at 5% FDR.

577

578    **References**

579    1000 Genomes Project Consortium, D. (2015). A global reference for human genetic
580    variation. Nature *526*, 68-74.

581

582    Bryois, J., Buil, A., Evans, D.M., Kemp, J.P., Montgomery, S.B., Conrad, D.F., Ho, K.M., Ring,
583    S., Hurles, M., Deloukas, P.*, et al.* (2014). Cis and trans effects of human genomic variants on
584    gene expression. PLoS genetics *10*, e1004461.

585

586    Buenrostro, J.D., Giresi, P.G., Zaba, L.C., Chang, H.Y., and Greenleaf, W.J. (2013).
587    Transposition of native chromatin for fast and sensitive epigenomic profiling of open
588    chromatin, DNA-binding proteins and nucleosome position. Nature methods *10*, 1213-1218.

589

590    Buil, A., Brown, A.A., Lappalainen, T., Vinuela, A., Davies, M.N., Zheng, H.F., Richards, J.B.,
591    Glass, D., Small, K.S., Durbin, R.*, et al.* (2015). Gene-gene and gene-environment interactions
592    detected by transcriptome sequence analysis in twins. Nature genetics *47*, 88-91.

593

594    Delaneau, O., Ongen, H., Brown, A.A., Fort, A., Panousis, N.I., and Dermitzakis, E.T. (2017). A
595    complete tool set for molecular QTL discovery and analysis. Nature communications *8*,
596    15452.

597

598    Dimas, A.S., Deutsch, S., Stranger, B.E., Montgomery, S.B., Borel, C., Attar-Cohen, H., Ingle,
599    C., Beazley, C., Gutierrez Arcelus, M., Sekowska, M.*, et al.* (2009). Common regulatory
600    variation impacts gene expression in a cell type-dependent manner. Science *325*, 1246-
601    1250.

602

603    Encode Project Consortium (2012). An integrated encyclopedia of DNA elements in the
604    human genome. Nature *489*, 57-74.

605

606    Grubert, F., Zaugg, J.B., Kasowski, M., Ursu, O., Spacek, D.V., Martin, A.R., Greenside, P.,
607    Srivas, R., Phanstiel, D.H., Pekowska, A.*, et al.* (2015). Genetic Control of Chromatin States in
608    Humans Involves Local and Distal Chromosomal Interactions. Cell *162*, 1051-1065.

609

610    GTEx Consortium (2015). Human genomics. The Genotype-Tissue Expression (GTEx) pilot
611    analysis: multitissue gene regulation in humans. Science *348*, 648-660.

612

613    Gutierrez-Arcelus, M., Lappalainen, T., Montgomery, S.B., Buil, A., Ongen, H., Yurovsky, A.,
614    Bryois, J., Giger, T., Romano, L., Planchon, A.*, et al.* (2013). Passive and active DNA
615    methylation and the interplay with genetic variation in gene regulation. eLife *2*, e00523.

616

617    Harrow, J., Frankish, A., Gonzalez, J.M., Tapanari, E., Diekhans, M., Kokocinski, F., Aken, B.L.,
618    Barrell, D., Zadissa, A., Searle, S.*, et al.* (2012). GENCODE: the reference human genome
619    annotation for The ENCODE Project. Genome research *22*, 1760-1774.

620

621  Heinz, S., Benner, C., Spann, N., Bertolino, E., Lin, Y.C., Laslo, P., Cheng, J.X., Murre, C., Singh,
622  H., and Glass, C.K. (2010). Simple combinations of lineage-determining transcription factors
623  prime cis-regulatory elements required for macrophage and B cell identities. Molecular cell
624  *38*, 576-589.

625

626  Johnson, D.S., Mortazavi, A., Myers, R.M., and Wold, B. (2007). Genome-wide mapping of in
627  vivo protein-DNA interactions. Science *316*, 1497-1502.

628

629  Kalhor, R., Tjong, H., Jayathilaka, N., Alber, F., and Chen, L. (2011). Genome architectures
630  revealed by tethered chromosome conformation capture and population-based modeling.
631  Nature biotechnology *30*, 90-98.

632

633  Kasowski, M., Kyriazopoulou-Panagiotopoulou, S., Grubert, F., Zaugg, J.B., Kundaje, A., Liu,
634  Y., Boyle, A.P., Zhang, Q.C., Zakharia, F., Spacek, D.V.*, et al.* (2013). Extensive variation in
635  chromatin states across humans. Science *342*, 750-752.

636

637  Kaufmann, S., Fuchs, C., Gonik, M., Khrameeva, E.E., Mironov, A.A., and Frishman, D. (2015).
638  Inter-chromosomal contact networks provide insights into Mammalian chromatin
639  organization. PloS one *10*, e0126125.

640

641  Kilpinen, H., Waszak, S.M., Gschwind, A.R., Raghav, S.K., Witwicki, R.M., Orioli, A.,
642  Migliavacca, E., Wiederkehr, M., Gutierrez-Arcelus, M., Panousis, N.I.*, et al.* (2013).
643  Coordinated effects of sequence variation on DNA binding, chromatin structure, and
644  transcription. Science *342*, 744-747.

645

646  Lander, E.S., Linton, L.M., Birren, B., Nusbaum, C., Zody, M.C., Baldwin, J., Devon, K., Dewar,
647  K., Doyle, M., FitzHugh, W.*, et al.* (2001). Initial sequencing and analysis of the human
648  genome. Nature *409*, 860-921.

649

650  Lappalainen, T., Sammeth, M., Friedlander, M.R., t Hoen, P.A., Monlong, J., Rivas, M.A.,
651  Gonzalez-Porta, M., Kurbatova, N., Griebel, T., Ferreira, P.G.*, et al.* (2013). Transcriptome
652  and genome sequencing uncovers functional variation in humans. Nature *501*, 506-511.

653

654  Li, H., and Durbin, R. (2009). Fast and accurate short read alignment with Burrows-Wheeler
655  transform. Bioinformatics *25*, 1754-1760.

656

657  Lieberman-Aiden, E., van Berkum, N.L., Williams, L., Imakaev, M., Ragoczy, T., Telling, A.,
658  Amit, I., Lajoie, B.R., Sabo, P.J., Dorschner, M.O.*, et al.* (2009). Comprehensive mapping of
659  long-range interactions reveals folding principles of the human genome. Science *326*, 289-
660  293.

661

662  MacArthur, J., Bowler, E., Cerezo, M., Gil, L., Hall, P., Hastings, E., Junkins, H., McMahon, A.,
663  Milano, A., Morales, J.*, et al.* (2017). The new NHGRI-EBI Catalog of published genome-wide
664  association studies (GWAS Catalog). Nucleic acids research *45*, D896-D901.

665

666 Marco-Sola, S., Sammeth, M., Guigo, R., and Ribeca, P. (2012). The GEM mapper: fast,
667 accurate and versatile alignment by filtration. Nature methods *9*, 1185-1188.
668
669 Maurano, M.T., Humbert, R., Rynes, E., Thurman, R.E., Haugen, E., Wang, H., Reynolds, A.P.,
670 Sandstrom, R., Qu, H., Brody, J*., et al.* (2012). Systematic localization of common disease-
671 associated variation in regulatory DNA. Science *337*, 1190-1195.
672
673 McVicker, G., van de Geijn, B., Degner, J.F., Cain, C.E., Banovich, N.E., Raj, A., Lewellen, N.,
674 Myrthil, M., Gilad, Y., and Pritchard, J.K. (2013). Identification of genetic variants that affect
675 histone modifications in human cells. Science *342*, 747-749.
676
677 Meaburn, K.J., and Misteli, T. (2007). Cell biology: chromosome territories. Nature *445*, 379-
678 781.
679
680 Nica, A.C., Montgomery, S.B., Dimas, A.S., Stranger, B.E., Beazley, C., Barroso, I., and
681 Dermitzakis, E.T. (2010). Candidate causal regulatory effects by integration of expression
682 QTLs with complex trait genetic associations. PLoS genetics *6*, e1000895.
683
684 Nicolae, D.L., Gamazon, E., Zhang, W., Duan, S., Dolan, M.E., and Cox, N.J. (2010). Trait-
685 associated SNPs are more likely to be eQTLs: annotation to enhance discovery from GWAS.
686 PLoS genetics *6*, e1000888.
687
688 Ongen, H., Brown, A.A., Delaneau, O., Panousis, N., Nica, A.C., and Dermitzakis, E.T. (2016).
689 Estimating the causal tissues for complex traits and diseases. bioRxiv.
690
691 Pickrell, J.K., Marioni, J.C., Pai, A.A., Degner, J.F., Engelhardt, B.E., Nkadori, E., Veyrieras, J.B.,
692 Stephens, M., Gilad, Y., and Pritchard, J.K. (2010). Understanding mechanisms underlying
693 human gene expression variation with RNA sequencing. Nature *464*, 768-772.
694
695 Pierce, B.L., Tong, L., Chen, L.S., Rahaman, R., Argos, M., Jasmine, F., Roy, S., Paul-Brutus, R.,
696 Westra, H.J., Franke, L*., et al.* (2014). Mediation analysis demonstrates that trans-eQTLs are
697 often explained by cis-mediation: a genome-wide analysis among 1,800 South Asians. PLoS
698 genetics *10*, e1004818.
699
700 Pombo, A., and Dillon, N. (2015). Three-dimensional genome architecture: players and
701 mechanisms. Nature reviews Molecular cell biology *16*, 245-257.
702
703 Quach, H., Rotival, M., Pothlichet, J., Loh, Y.E., Dannemann, M., Zidane, N., Laval, G., Patin,
704 E., Harmant, C., Lopez, M*., et al.* (2016). Genetic Adaptation and Neandertal Admixture
705 Shaped the Immune System of Human Populations. Cell *167*, 643-656 e617.
706
707 Rao, S.S., Huntley, M.H., Durand, N.C., Stamenova, E.K., Bochkov, I.D., Robinson, J.T.,
708 Sanborn, A.L., Machol, I., Omer, A.D., Lander, E.S*., et al.* (2014). A 3D map of the human
709 genome at kilobase resolution reveals principles of chromatin looping. Cell *159*, 1665-1680.
710 Scutari, M. (2010). Learning Bayesian Networks with the bnlearn R Package Journal of
711 Statistical Software *35*, 1-22.
712

713 Spitz, F., and Furlong, E.E. (2012). Transcription factors: from enhancer binding to
714 developmental control. Nature reviews Genetics *13*, 613-626.

715

716 Storey, J.D., and Tibshirani, R. (2003). Statistical significance for genomewide studies.
717 Proceedings of the National Academy of Sciences *100*, 9440-9445.

718

719 Waszak, S.M., Delaneau, O., Gschwind, A.R., Kilpinen, H., Raghav, S.K., Witwicki, R.M., Orioli,
720 A., Wiederkehr, M., Panousis, N.I., Yurovsky, A.*, et al.* (2015). Population Variation and
721 Genetic Control of Modular Chromatin Architecture in Humans. Cell *162*, 1039-1050.

722

723 Whyte, W.A., Orlando, D.A., Hnisz, D., Abraham, B.J., Lin, C.Y., Kagey, M.H., Rahl, P.B., Lee,
724 T.I., and Young, R.A. (2013). Master transcription factors and mediator establish super-
725 enhancers at key cell identity genes. Cell *153*, 307-319.

726

727 Williams, A., Spilianakis, C.G., and Flavell, R.A. (2010). Interchromosomal association and
728 gene regulation in trans. Trends in genetics : TIG *26*, 188-197.

729

730 Zarrei, M., MacDonald, J.R., Merico, D., and Scherer, S.W. (2015). A copy number variation
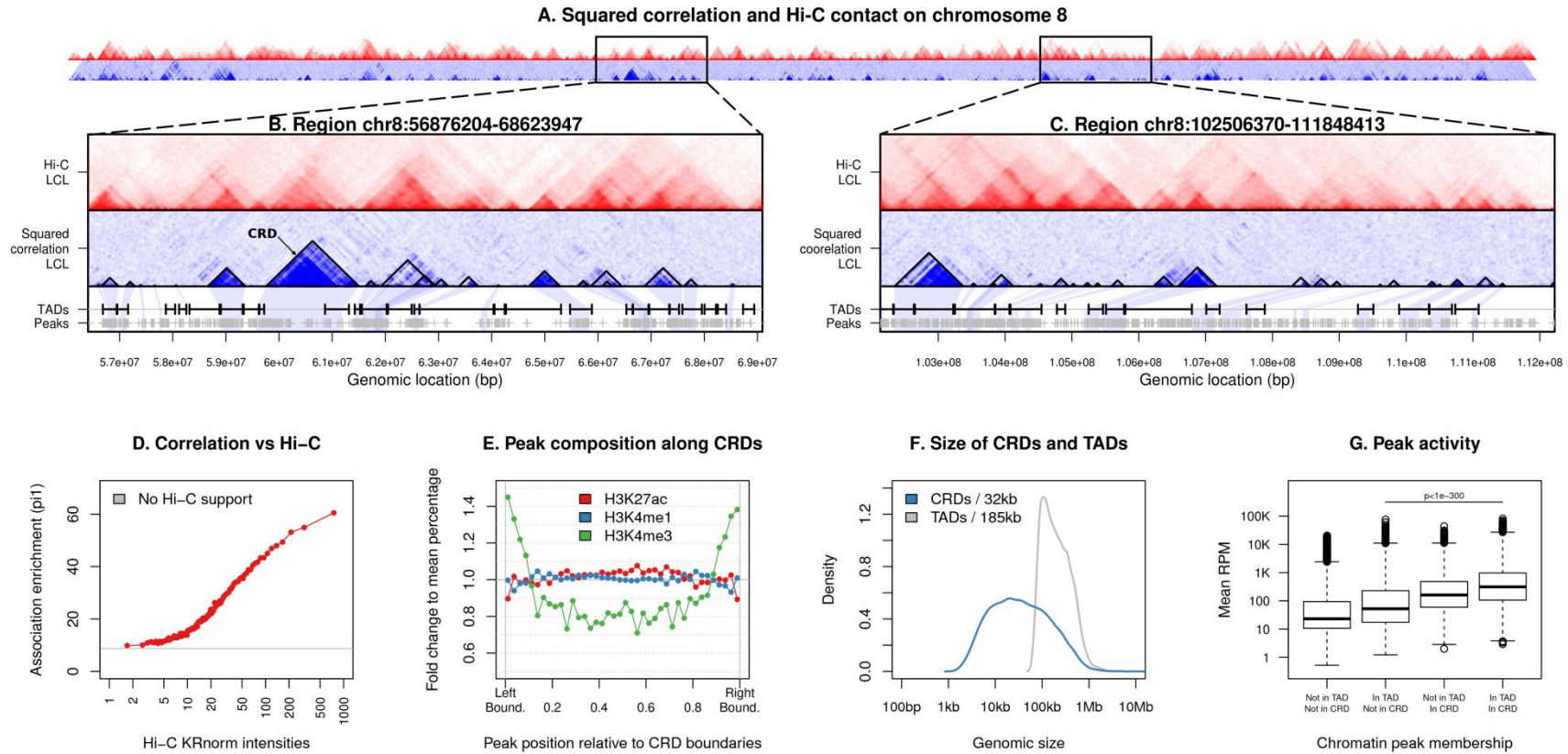731 map of the human genome. Nature reviews Genetics *16*, 172-183.

732

**Figure1: The *Cis* Regulatory Domains (CRDs).** The squared correlations and Hi-C contact intensities (scaled between 0 and 1) between nearby chromatin peaks on chromosome 8 are shown in panel A with shades of blues and reds, respectively. Two examples containing 1,000 peaks are also shown in panels B and C together with (i) the CRDs (black triangles), (ii) the TADs (black segments) and (iii) the chromatin peak locations (grey crosses). The concordance between Hi-C and the correlation is shown in panel D: the association enrichment in the full LCL data was measured within 100 bins of increasing Hi-C contact intensities. The chromatin peak density as we move along CRDs is shown in panel E. The x-axis represents the position scaled between 0 and 1 of the chromatin peaks relative to the left and right boundaries of CRDs. The panel F compares the lengths of CRDs and TADs. And finally, the panel G gives the chromatin peak activity depending on their TAD and CRD memberships and expressed as Reads Per Millions.
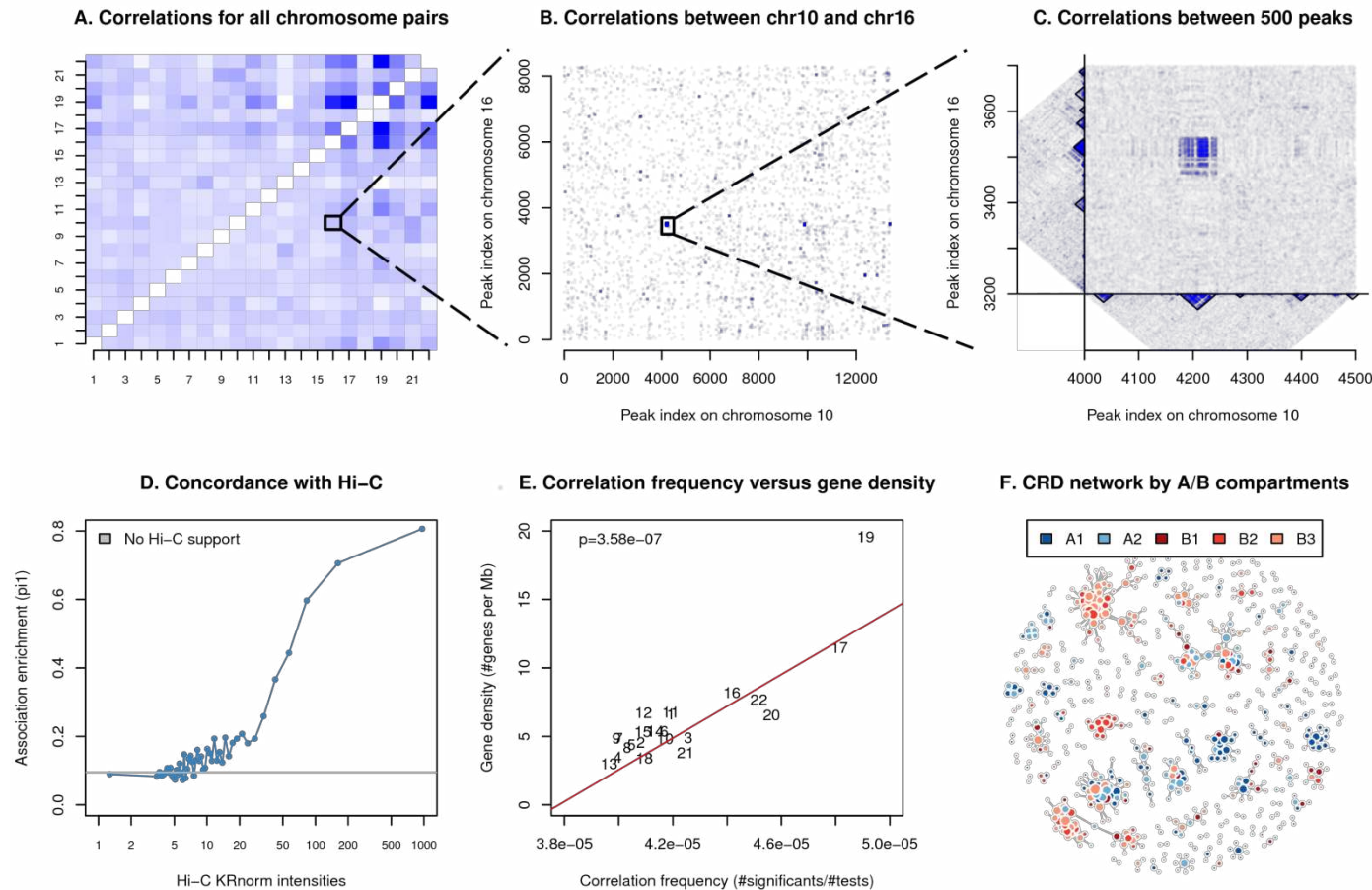
**Figure2: The inter-chromosomal interactions of chromatin.** The three first panels show the correlation structure between chromosomes at three different zoom levels: (A) the frequency of correlations with a P-value < 1e-6 for all pairs of chromosomes, (B) all pairs of peaks between chromosome 10 and 16 with a P-value < 1e-4 and (C) a 500x500 peaks region exhibiting high correlation; all this in the context of CRDs and local correlation between chromatin peaks. The panel D gives the association enrichments of inter-chromosomal correlations within 50 bins of increasing Hi-C contact intensities. The panel E gives the association frequency as a function of the gene density on a per-chromosome basis. The panel F shows a network representation of the 1,000 strongest inter-chromosomal associations between CRDs colored by nuclear compartment membership.
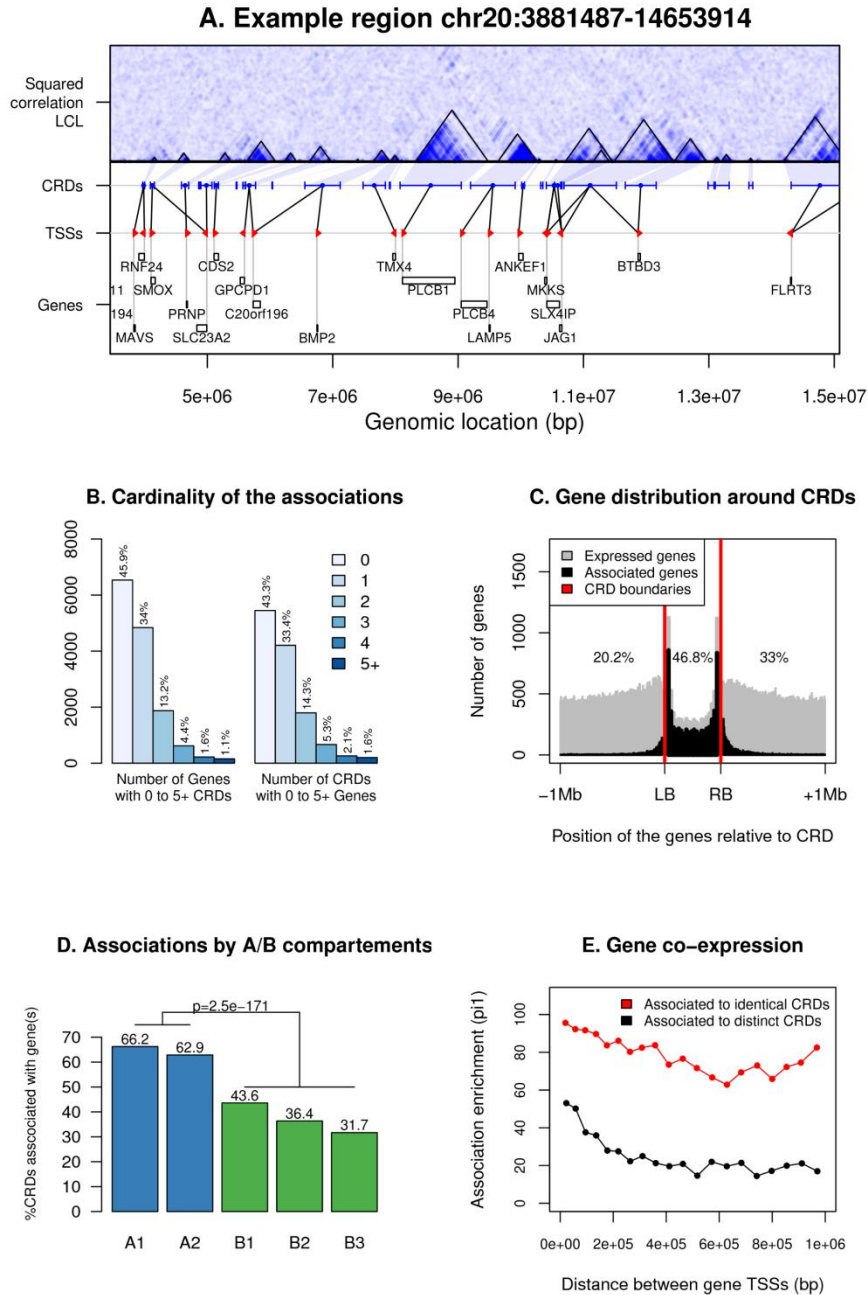
**A. Example region chr20:3881487-14653914**

**B. Cardinality of the associations**

**C. Gene distribution around CRDs**

**D. Associations by A/B compartements**

**E. Gene co-expression**

**Figure3: The regulatory function of CRDs.** The first panel shows the associations between CRDs (blue segments) and genes (gene TSSs with red arrows and gene body with black boxes) within an example region of chromosome 20. Then, the panel B shows the connectivity between genes and CRDs while the panel C the distributions of tested genes (in grey) and significantly associated genes (in black) (i) within CRDs (relative locations; LB and RB correspond to left and right CRD boundaries, respectively) and (ii) within the genomic regions surrounding CRDs (absolute positions +/- 1 Mb). Panel (D) shows the percentage of CRDs being associated with genes within A/B nuclear compartments (significance comparing A to B is shown on top). The panel E gives the association enrichment in between genes associated with identical (in red) or distinct (in black) CRDs as a function of the distance between them.
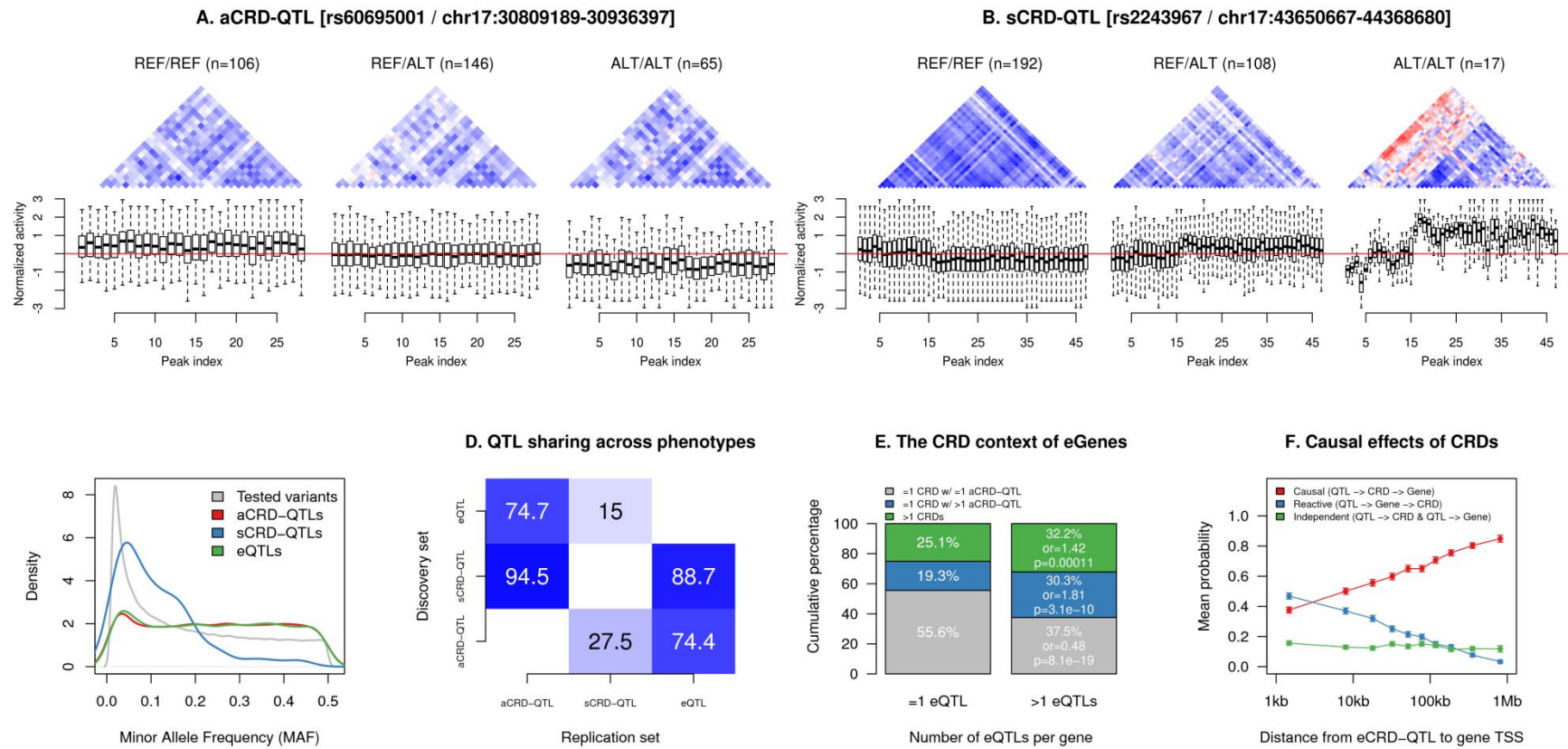
**Figure4: The genetic control of CRDs.** The panels A and B show an example of aCRD-QTL (CRD-activity QTL) and sCRD-QTL (CRD-structure QTL), respectively. In each case, the correlation structure (positive and negative correlations in blue and red, respectively) and the normalized activity distribution (bottom panels) for all chromatin peaks contained in the CRD are shown stratified by the QTL genotype (REF/REF, REF/ALT and ALT/ALT). The numbers of individuals carrying each genotype are shown on top of each panel. The panel C shows the minor allele frequency spectrum for all types of QTLs while the panel D gives the association enrichments of QTLs across all molecular phenotypes (i.e. overlaps or degree of sharing). The panel E shows the respective proportions of genes with one eQTL or more being associated with (i) one CRD having one aCRD-QTL (in grey), (ii) one CRD with two or more aCRD-QTLs (in blue) and (iii) two or more CRDs (in green). The panel F gives the mean probability of each causal model as a function of distance between the eCRD-QTLs and the target genes. Only data for eCRD-QTLs falling within their target CRD is shown here.
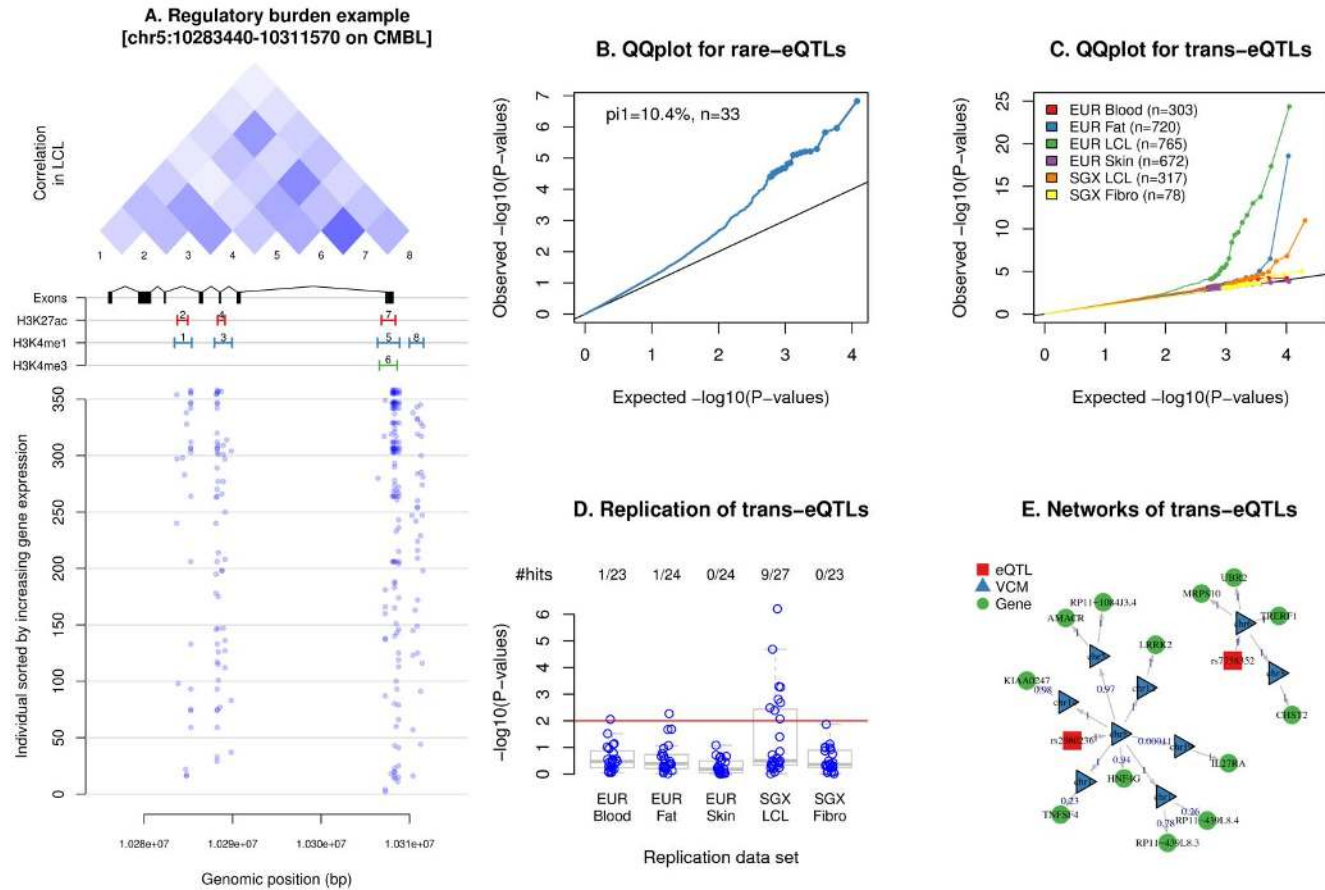
**Figure5:** **CRDs enhance association studies.** The left panel (A) shows an example of how CRD can be used to perform a burden test of rare variants on gene expression. From top to bottom are shown (i) the pairwise correlations between peaks within the CRD, (ii) the positions of the various genomic features (peaks and gene) and (iii) the positions of the minor alleles at variants overlapping peaks for each of the 358 tested individual (sorted by increasing gene expression; low expression at the bottom). The QQplots in panel B and C show the association enrichments obtained when mapping *rare*-eQTLs (B; burden test) and *trans*-eQTLs in 6 distinct data sets (C). The panel D gives the replication P-values of the 27 significant hits found to be significant in EUROBATS LCL data. The panel E gives a network representation of the 9 replicated trans-eQTLs together with the relevant *cis* and *trans* CRDs and genes. The edge directions are set up given causal hypotheses and their robustness assessed using Bayesian networks: each weight ranges between 0 and 1; the latter meaning a direction of effect highly supported by the data.