

**A peer-reviewed version of this preprint was published in PeerJ on 28 September 2017.**

[View the peer-reviewed version](https://peerj.com/articles/3742) (peerj.com/articles/3742), which is the preferred citable publication unless you specifically need to cite this preprint.

Dzida T, Iqbal M, Charapitsa I, Reid G, Stunnenberg H, Matarese F, Grote K, Honkela A, Rattray M. 2017. Predicting stimulation-dependent enhancer-promoter interactions from CHIP-Seq time course data. PeerJ 5:e3742  
<https://doi.org/10.7717/peerj.3742>

# Predicting context specific enhancer-promoter interactions from CHIP-Seq time course data

**Tomasz Dzida** <sup>Corresp., 1</sup>, **Mudassar Iqbal** <sup>1</sup>, **Iryna Charapitsa** <sup>2</sup>, **George Reid** <sup>2</sup>, **Henk Stunnenberg** <sup>3</sup>, **Filomena Matarese** <sup>3</sup>, **Antti Honkela** <sup>4</sup>, **Magnus Rattray** <sup>Corresp., 1</sup>

<sup>1</sup> Faculty of Biology, Medicine and Health, University of Manchester, Manchester, United Kingdom

<sup>2</sup> Chemical Biology Core Facility, European Molecular Biology Laboratory, Heidelberg, Germany

<sup>3</sup> Department of Molecular Life Sciences, Radboud University Nijmegen, Nijmegen, Netherlands

<sup>4</sup> Helsinki Institute for InformationTechnology (HIIT), Department of Computer Science, University of Helsinki, Helsinki, Finland

Corresponding Authors: Tomasz Dzida, Magnus Rattray

Email address: tomasz.dzida@gmail.com, magnus.rattray@manchester.ac.uk

We have developed a machine learning approach to predict context specific enhancer-promoter interactions using evidence from changes in genomic protein occupancy over time. The occupancy of estrogen receptor alpha (ER $\alpha$ ), RNA polymerase (Pol II) and histone marks H2AZ and H3K4me3 were measured over time using ChIP-Seq experiments in MCF7 cells stimulated with estrogen. A Bayesian classifier was developed which uses the correlation of temporal binding patterns at enhancers and promoters and genomic proximity as features to predict interactions. This method was trained using experimentally determined interactions from the same system and was shown to achieve much higher precision than predictions based on the genomic proximity of nearest ER $\alpha$  binding. We use the method to identify a genome-wide confident set of ER $\alpha$  target genes and their regulatory enhancers genome-wide. Validation with publicly available GRO-Seq data demonstrates that our predicted targets are much more likely to show early nascent transcription than predictions based on genomic ER $\alpha$  binding proximity alone.

# Predicting context specific enhancer-promoter interactions from ChIP-Seq time course data

Tomasz Dzida<sup>1</sup>, Mudassar Iqbal<sup>1</sup>, Iryna Charapitsa<sup>2</sup>, George Reid<sup>2</sup>, Henk Stunnenberg<sup>3</sup>, Filomena Matarese<sup>3</sup>, Antti Honkela<sup>4</sup>, and Magnus Rattray<sup>1</sup>

<sup>1</sup>Faculty of Biology, Medicine and Health, University of Manchester, Manchester, UK

<sup>2</sup>Chemical Biology Core Facility, European Molecular Biology Laboratory, Heidelberg, Germany

<sup>3</sup>Department of Molecular Life Sciences, Radboud University Nijmegen, Nijmegen, The Netherlands

<sup>4</sup>Helsinki Institute for Information Technology (HIIT), Department of Computer Science, University of Helsinki, Helsinki, Finland

Corresponding author:

Tomasz Dzida, Magnus Rattray

Email address: Tomasz.Dzida@gmail.com, Magnus.Rattray@manchester.ac.uk

## ABSTRACT

We have developed a machine learning approach to predict context specific enhancer-promoter interactions using evidence from changes in genomic protein occupancy over time. The occupancy of estrogen receptor alpha (ER $\alpha$ ), RNA polymerase (Pol II) and histone marks H2AZ and H3K4me3 were measured over time using ChIP-Seq experiments in MCF7 cells stimulated with estrogen. A Bayesian classifier was developed which uses the correlation of temporal binding patterns at enhancers and promoters and genomic proximity as features to predict interactions. This method was trained using experimentally determined interactions from the same system and was shown to achieve much higher precision than predictions based on the genomic proximity of nearest ER $\alpha$  binding. We use the method to identify a genome-wide confident set of ER $\alpha$  target genes and their regulatory enhancers genome-wide. Validation with publicly available GRO-Seq data demonstrates that our predicted targets are much more likely to show early nascent transcription than predictions based on genomic ER $\alpha$  binding proximity alone.

## INTRODUCTION

Gene expression is dependent upon the binding of transcription factor (TF) proteins to genomic regions which regulate transcriptional initiation (Nagarajan et al., 2014). In eukaryotic cells, these regulatory genomic regions are referred to as promoters and enhancers. The transcriptional competence of DNA in eukaryotes is determined by its organization in chromatin. Chromatin structure is dynamically regulated at multiple levels, including ATP-dependent chromatin remodelling and histone modifications (Bernstein et al., 2005; Bannister and Kouzarides, 2011; Zhu et al., 2013; Stasevich et al., 2014). Enhancers can act upstream or downstream of their target gene promoters and are often distal, separated by large inter-genic regions (Schoenfelder et al., 2010; Sanyal et al., 2012; Shen et al., 2012). Enhancer-promoter interactions require protein-mediated physical contact through formation of chromatin loops (Tolhuis et al., 2002). Although most contacts are intra-chromosomal, there are some interactions between loci from different chromosomes (Fullwood et al., 2009; Li et al., 2010, 2012). Interactions can also exist as part of large multi-gene and multi-enhancer complexes (Fullwood et al., 2009; Li et al., 2012).

Recent progress in experimental techniques such as ChIA-PET, 3C and its derivatives 4C, 5C, and Hi-C (Fullwood et al., 2009; Dekker et al., 2002; Hagège et al., 2007; Zhao et al., 2006; Dostie et al., 2006; Simonis et al., 2007; van Steensel and Dekker, 2010; Nagano et al., 2013; Jin et al., 2013) have mapped large numbers of chromatin interactions, including enhancer-promoter interactions. However, these methods are technically challenging and genome-wide methods, such as HiC, typically lack the

46 resolution required to identify individual interacting enhancer elements. Some methods are also thought to  
47 produce a high false negative rate (in case of ChIA-PET, 5C; Li et al., 2012; He et al., 2014) or cannot be  
48 applied on a genome-wide scale (3C,4C; Simonis et al., 2007). Capture-HiC methods have recently been  
49 developed (Mifsud et al., 2015; Javierre et al., 2016) to improve genomic resolution through focussing on  
50 predetermined genomic regions, e.g. promoters, and show promise but are not yet widely used. Data from  
51 these technologies can also be noisy and subject to various sources of bias which can be problematic to  
52 correct (van Steensel and Dekker, 2010). In addition, the physical contact between two chromatin regions  
53 does not determine a functional interaction (Shlyueva et al., 2014) with stimulus-dependant behaviour of  
54 chromatin looping adding a further layer of complexity (Drissen et al., 2004; Vakoc et al., 2005). For  
55 these reasons, complementary approaches to infer enhancer-promoter interactions by exploiting readily  
56 available sources of genomic data, such as ChIP-Seq and RNA-Seq data, are of interest.

57 ChIP-seq experiments enable the discovery of the genomic location of transcriptionally relevant  
58 proteins such as TFs, RNA polymerase and modified histones. Multiple ChIP-Seq datasets can be  
59 combined with data from other relevant genomics assays to identify active promoters and enhancers  
60 using genomic segmentation algorithms (Zhu et al., 2013; Ernst et al., 2011). Others have also used  
61 ChIP-seq and RNA-seq datasets to infer enhancer-promoter interactions. For example, Ernst et al. (2011)  
62 used histone mark data from multiple cell-types to identify active enhancers and promoters from which  
63 enhancer-associated data was correlated with expression data from genes within 125kbp to identify likely  
64 interactions. Thurman et al. (2012) used DNase I hypersensitivity (DHS) data from multiple cell-types to  
65 correlate and link distal DNase hypersensitivity sites (within 500kbp) to those within putative gene targets.  
66 Similarly, Andersson et al. (2014) predicted enhancer-promoter links by correlating CAGE enhancer RNA  
67 to CAGE promoter RNA.

68 Approaches for discovering cell-type specific interactions include PreSTIGE (Corradin et al., 2014),  
69 RIPPLE (Roy et al., 2015), and the method developed by Marstrand and Storey (2014). PreSTIGE uses a  
70 method based on the Shannon entropy to identify cell-type specific interactions between enhancers and  
71 genes using H3K4me1 and RNA-seq data respectively. The regions are linked within promoter-centric  
72 domains, bounded on each side by the minimal distance of 100kbp up to the first CTCF binding site from a  
73 TSS. RIPPLE uses ENCODE data from four cell-lines each with 11 ChIP-Seq datasets (RNA-seq, CTCF,  
74 RAD21, DNase1, TBP and histone marks) to train a random forest classifier which predicts enhancer-gene  
75 interactions within 1MB distance. The features used are two joint binary vectors of presence/absence of  
76 dataset signal peak over a promoter and enhancer, correlation of entries of the vectors, as well as gene  
77 expression of the promoter controlled gene. Marstrand and Storey (2014) developed a method to aggregate  
78 RNA-seq data over genes and DHS data over  $\pm 200$ kb regions surrounding them for twenty different cell  
79 lines. The method searches through each gene and cell-line for unexpected DHS/RNA-seq ratios and  
80 once found, scans across the gene vicinities in search of causal, local DHS variabilities. Lastly, a method  
81 proposed by He et al. (2014) uses a random forest classifier to find enhancer-gene interactions. The  
82 method uses three features: evolutionary conservation, correlation of enhancer scores derived from histone  
83 marks from RNA-seq data, and an average of correlations between TF ChIP-Seq and gene expression  
84 across 12 cell-types. A distance constraint is also imposed to aid inference.

85 The majority of the above methods require data from multiple cell-types and therefore do not allow  
86 discovery of interactions given data from one cell-type. Most existing methods also assume a stringent  
87 distance constraint and are therefore unable to discover distal links beyond this constraint. Finally, these  
88 methods do not take into account evidence from time course data.

89 We show how ChIP-Seq time course data that reports TF and RNA polymerase occupancy at multiple  
90 time points after cellular stimulation can be used to predict enhancer-promoter interactions within  
91 chromosomes. We have developed a Bayesian classifier that combines evidence from the correlation  
92 of ChIP-Seq time course data at enhancers and across gene bodies with the genomic separation of  
93 interacting elements as features. We apply our method to time course data from MCF7 breast cancer cells  
94 after stimulation with estradiol and we benchmark performance against publically available ChIA-PET  
95 data from this system. We show that our method performs much better than association by proximity,  
96 identifying many more interactions than predictions based on proximity alone. Estrogen Receptor (ER- $\alpha$ )  
97 and RNA polymerase (Pol II) ChIP-Seq time course data are shown to be highly informative for predicting  
98 interactions. We also stratify our predicted interactions to those that lie within Topologically Associating  
99 Domains (TADS; Dixon et al., 2012) and those that span TADs, showing that our classifier can make  
100 useful predictions in both categories. Finally, we use our predictions to provide a highly confident list of

101 directly ER-regulated target genes in this system and validate it against a GRO-seq dataset. Our predicted  
102 targets are much more likely to show early nascent transcription than predictions based on genomic ER- $\alpha$   
103 binding proximity alone and predicted targets are involved in many biological processes associated with  
104 breast cancer. Our model thus offers biologically meaningful insight into the early transcriptional response  
105 to ER- $\alpha$ .

## 106 MATERIALS AND METHODS

### 107 Data Preparation

108 The aim of our experiment was to uncover the early response to estradiol (E2) in MCF7 breast cancer cells.  
109 Our previous studies included only the Pol-II and RNA-Seq time course data from these experiments  
110 (wa Maina et al., 2014; Honkela et al., 2015) and here we include additional ChIP-Seq datasets. The first  
111 step was to create a reference sample in a ligand free environment. For that, the cells were placed into  
112 estradiol free media for 3 days, which reduced the binding between ER- $\alpha$  and E2. The cells were then  
113 ready to be re-exposed to E2. Following the introduction of E2, the resultant changes were tracked by  
114 multiple ChIP-seq experiments. The experiments were performed at 0, 5, 10, 20, 40, 80, 160, 320, 640 and  
115 1280 minutes after the stimulation. Each ChIP-seq experiment was carried out with a different antibody to  
116 measure genome-wide changes in genomic occupancy of their specific protein targets. Specifically, the  
117 studied protein factors and histone modifications were: ER- $\alpha$ , H3K4me3, and H2AZ (data available from  
118 GEO: accession GSM2467201). Other previously published data from the same set of experiments are  
119 available for Pol-II ChIP-Seq and RNA-Seq (GEO accession GSE62789 and GSE44800; wa Maina et al.,  
120 2014; Honkela et al., 2015).

121 *Preparation of MCF-7 cells:* The MCF-7 human breast cancer cell line originates from a 69-year old  
122 Caucasian woman and is estrogen receptor (ER) positive, progesterone positive (PR) and HER2 negative.  
123 Here MCF-7 cells (a clonal isolate obtained from the ATCC (catalogue number HTB-22) kindly provided  
124 by Prof. Edison Liu, Jackson Laboratories, Maine, USA) were grown in 15cm plates to 80% confluency.  
125 Plates were then washed 2 times with PBS and overlaid with 20 ml of phenol-red free high glucose  
126 DMEM (Gibco) containing 2% charcoal stripped FCS (Sigma). After 24 hours of incubation, the cells  
127 were again washed with PBS and fresh media containing 2% charcoal stripped FCS was added. This  
128 process was repeated over a three day period to generate cells devoid of estrogen. The time course (5, 10,  
129 20, 40, 80, 160, 320, 640 and 1280 minutes) was initiated by replacing media with prewarmed media  
130 containing 10 nM E2. In addition, an untreated sample was included in the experiment as a zero time  
131 point.

132 *ChIP-seq protocols and methods:* Cells were fixed for 10 minutes at room temperature by the addition  
133 of formaldehyde to a final concentration of 1%, after which glycine was added to a concentration of 100  
134 mM. Cells were then washed twice with PBS and collected into 2 ml of lysis buffer (150 mM NaCl, 20  
135 mM Tris pH 8.0, 2 mM EDTA, 1% triton X-100, protease inhibitor [complete EDTA free, Roche, 04 693  
136 132 001], 100 mM PMSF). The lysate was sonicated for 3  $\times$  30 seconds using a Branson ultrasonicator  
137 equipped with a microtip on a power setting of 3 and a duty cycle of 90%. Samples were cooled on  
138 ice between rounds of sonication. Alternatively, a Bioruptor sonicator was used (power high, 15 mins  
139 total, 30 s on 30 s off; total volume of sample –1 ml) to fragment chromatin. In either case, the resulting  
140 sonicate was centrifuged at 4000xg for 5 minutes, an aliquot of 10% retained for input and the remaining  
141 material transferred to a fresh tube. Four mg of anti-ERAantibody (HC-20, rabbit polyclonal, Santa Cruz,  
142 sc-543), 2 mg of anti-RNA Polymerase II antibody (AC-055-100, monoclonal, Diagenode, 001), 3 mg  
143 of anti-H3K4me3 antibody (pAb-MEHAHS-024, rabbit polyclonal, Diagenode, HC-0010) and 2 mg  
144 anti-Histone H2A.Z (acetyl K4+K7+K11) antibody (ab18262, sheep polyclonal, Abcam, 659355) were  
145 added to the samples, which were then incubated overnight at 40C with rotation. Chromatin antibody  
146 complexes were isolated, either by addition of 10 ml of protein G labeled magnetic beads (Millipore  
147 Pureproteome protein G magnetic beads, LSKMAGG10) prewashed in lysis buffer or with 20 ml protein  
148 A/G beads (Santa Cruz). Afterwards, the complexes obtained with protein G magnetic beads were washed  
149 three times with lysis buffer, then reverse crosslinked in 0.5 ml 5 M guanidine hydrochloride, 20 mM  
150 Hepes, 30% isopropanol, 10 mM EDTA for a minimum of 4 hours at 650C. Recovered DNA was then  
151 purified using a Qiaquick spin column and eluted in 50 ml of 10 mM Tris pH 8.0. Where protein A/G  
152 beads were used, the complexes were washed sequentially with three different buffers at 40C: two times  
153 with solution of composition 0.1% SDS, 0.1% DOC, 1% Triton, 150 mM NaCl, 1 mM EDTA, 0.5 mM  
154 EGTA, 20 mM HEPES pH 7.6, once with the solution as before but with 500 mM NaCl, once with

155 solution of composition 0.25 M LiCl, 0.5% DOC, 0.5% NP-40, 1 mM EDTA, 0.5 mM EGTA, 20 mM  
156 HEPES pH 7.6 and two times with 1 mM EDTA, 0.5 mM EGTA, 20 mM HEPES pH 7.6. A control  
157 library was generated by sequencing input DNA (non-ChIP genomic DNA). Immunopurified chromatin  
158 was eluted with 200 ml of elution buffer (1% SDS, 0.1 M NaHCO<sub>3</sub>), incubated at 650C for 4 h in the  
159 presence of 200 mM NaCl, isolated using a Qiaquick spin column and eluted in 50 ml of 10 mM Tris  
160 pH 8.0. Libraries were prepared for Illumina sequencing according to the manufacturer's protocols  
161 (Illumina). Briefly, DNA fragments were subject to sequential end repair and adaptor ligation. DNA  
162 fragments were subsequently size selected (approx. 300 base pair [bp]). The adaptor-modified DNA  
163 fragments were amplified by limited PCR (14 cycles). Quality control and concentration measurements  
164 were made by analysis of the PCR products by electrophoresis (Experion, BioRad) and by fluorometric  
165 dye binding using a Qubit fluorometer with the Quant-iT dsDNA HS Assay Kit (Invitrogen, Q32851)  
166 respectively. Cluster generation and sequencing-by-synthesis (36 bp) was performed using the Illumina  
167 Genome Analyzer IIx (GAIIx) according to standard protocols of the manufacturer (Illumina).

### 168 **Alignment to a reference human genome**

169 Raw reads from the experiments were mapped onto the human reference genome (NCBI\_build37) using  
170 the Genomatix Mining Station (version 3.5.2) to enable further analysis. The sequencing depth, i.e. the  
171 total number of sequenced reads, was very similar for each dataset, however, on average only 81%, 76%,  
172 67%, 61%, 64% of ER- $\alpha$ , Pol-II (rep 1), Pol-II (rep 2), H3K4me3, and H2AZ ChIP-seq reads were  
173 mapped uniquely to the genome. The non-uniquely mapped reads were discarded from further analysis.  
174 Using the statistical criterion provided by MACS, we established that our sequencing depth allows for no  
175 duplicates of reads, thus we discarded any duplicated reads as they are most likely an artefact in ChIP-Seq.

### 176 **ER- $\alpha$ Binding Locations**

177 The MACS package (v2.0, p-value:  $1e-7$ , no control, estimation of  $\lambda_{local}$  off) Zhang et al. (2008) was used  
178 for peak-calling and applied to each of the 0, 5, 10, 20, ..., 320 min time course datasets to estimate ER- $\alpha$   
179 binding locations. The last two time points (640 and 1280 mins) were not included as the number of ER- $\alpha$   
180 mapped reads was found to be very low at these times compared to earlier times. Persistent co-occurring  
181 ER- $\alpha$  binding locations (i.e occurring at least twice across two time points after  $t = 0$ ) were merged by a  
182 union operation (similar to the mergeBED method from BEDTools (Quinlan and Hall, 2010)), otherwise  
183 they were discarded. The method is illustrated in Figure S1. Since our analysis is aimed at intergenic  
184 ER- $\alpha$ -bound enhancers, we ignored the consensus peaks which overlapped with either gene bodies or  
185 upstream 300bp-long regions by which the genes were extended to account for a promoter region.

### 186 **Time-Series Construction**

187 We calculated the mapped read counts for each individual time point ChIP-seq dataset over the consensus  
188 ER- $\alpha$  binding sites to create time series over enhancer regions for each of our antibodies. To normalise  
189 the counts, we divided each read count over the total number of uniquely mapped and non-duplicated  
190 reads across all time points and multiplied the resultant values by the total number of mapped reads in  
191 the  $t = 0$  min dataset. We concatenated the normalised counts to produce time series for each ChIP-seq  
192 dataset. We refer to each enhancer time series as  $\mathbf{X}_{j,n}$ , where  $j \in J$  (number of intergenic enhancers) and  
193  $n \in N$  (number of time course ChIP-seq datasets). We repeated the process for the gene regions to create  
194 the analogous time series over gene regions, extending the genes by 300bp upstream from their canonical  
195 TSS. We refer to each time series over gene as  $\mathbf{Y}_{k,n}$  where  $k \in K$  (number of genes). We filtered out genes  
196 and intergenic enhancers from consideration if the total number of mapped reads across any time series  
197 was less than 30.

### 198 **Clustering**

199 To help visualise the occupancy dynamics of Pol II and ER- $\alpha$  at enhancers and genes we clustered the  
200 data with the R-implementation of Affinity Propagation (AP) (Frey and Dueck, 2007). AP is a clustering  
201 method based on belief propagation and works iteratively by passing messages between data points  
202 until exemplars (cluster centres) automatically emerge. A preference parameter  $p$  has an effect on the  
203 final number of clusters. The R implementation of AP can search through values of  $p$  to achieve an  
204 approximately pre-specified number of clusters. The method is similar to k-means but can achieve much  
205 better optimisation of the k-means objective function than the standard EM algorithm.

206 To reduce the effect of noise, for Pol II we clustered only the pairs of the time series for which the  
 207 Pearson correlation coefficient was at least 0.2 between replicates and the total number of mapped reads  
 208 was at least 30. For ER- $\alpha$ , due to lack of replicates, we only clustered the time series with more than 100  
 209 reads in total across all times. Prior to the clustering we standardized each time series to z-scores to bring  
 210 all time series onto the same scale. We obtained 20 and 22 clusters for Pol II time series over enhancer and  
 211 genes, respectively. Similarly we obtained 21 and 21 clusters for ER- $\alpha$  time series over enhancer and  
 212 genes. We also jointly clustered time series of PolII and ER- $\alpha$ . The results of the clustering can be seen  
 213 in Figure S2.

## 214 **Enhancer-centric model**

215 Suppose that an enhancer  $j = 1, \dots, J$  regulates a gene  $k = 1, \dots, K$  at a number of time points, and  
 216 that their contact is mediated by a protein. We can expect that the time course data of ChIP-seq data at  
 217 an enhancer  $j$  i.e.  $\mathbf{X}_j = (x_{j,1}, \dots, x_{j,D})$  and gene  $k$  i.e.  $\mathbf{Y}_k = (y_{k,1}, \dots, y_{k,D})$  would on average be more  
 218 correlated for interacting pairs than their non-interacting counterparts. Here, we intend to learn the  
 219 underlying distribution of correlations of the two classes of pairs for four complementary datasets and on  
 220 their basis jointly classify a new unobserved instance. In addition, we combine the time course derived  
 221 attributes with the corresponding distribution of genomic separation for interacting and non-interacting  
 222 elements.

## 223 **Definition of the model**

224 Our model is defined in terms of two  $K$ -dimensional random variables  $\mathbf{I}_j = I_{j,1}, \dots, I_{j,K}$  and  $\mathbf{D}_j =$   
 225  $D_{j,1}, \dots, D_{j,K}$ . The first variable  $\mathbf{I}_j$  encodes a structure of simultaneous contacts of a given enhancer  $j$   
 226 with its surrounding  $K$  putative target genes. It has  $K$  binary entries  $I_{j,k}$  indicating whether  $(E_j, G_k)$  forms  
 227 an interacting ( $I_{j,k} = 1$ ) or non-interacting pair ( $I_{j,k} = 0$ ). The variable  $\mathbf{D}_j$  is a  $K \times N$ -dimensional matrix  
 228 of observed attributes with each row ( $D_{j,k}$ ) consisting of  $N$  values of pair-wise comparisons between  
 229 time series of an enhancer  $j$  and a gene  $k$ , and their genomic location. The first set of comparisons  
 230 rely on Pearson correlation and involves calculating its value  $c_{j,k,n}$  for each pair  $(E_j, G_k)$ , i.e. its time  
 231 series  $(\mathbf{X}_{j,n}, \mathbf{Y}_{k,n})$ , and for each dataset  $n \in N$ , where  $N$  is a number of time course ChIP-seq datasets.  
 232 Additionally, the data vector also contains the Euclidean distance  $d_{j,k}$  calculated between the genomic  
 233 coordinates of the canonical TSS of a gene  $k$  to the centre of an enhancer  $j$ .

The joint likelihood of the model can be written as:

$$P(\mathbf{D}_j, \mathbf{I}_j) = P(\mathbf{D}_j | \mathbf{I}_j) P(\mathbf{I}_j). \quad (1)$$

234 The model provides a probability of observing a particular  $\mathbf{D}_j$  under a given structure  $\mathbf{I}_j$ . Due to its  
 235 regulatory role, an enhancer is unlikely to regulate a high number of genes, thus we can expect that the  
 236 true  $P(\mathbf{I}_j)$ , which in the Bayesian treatment is a prior distribution over the structures, would be sparse.  
 237 Moreover, we could expect that  $D_{j,k}$  and  $D_{j,k'}$  of any two interacting pairs  $k, k'$  would be interlinked, as  
 238 correlations between gene-enhancer pairs are not independent variables. These dependencies would be  
 239 reflected in a true form of the likelihood  $P(\mathbf{D}_j | \mathbf{I}_j)$ . Lastly, we could also expect that the  $N + 1$  attributes  
 240 i.e. correlations  $c_{j,k,n}$  and distance  $d_{j,k}$  of a pair  $j, k$  of the vector  $D_{j,k}$  would also be correlated.

## 241 **Simplifying the likelihood and Naive Bayes**

242 The modelling of all dependencies however is difficult given the relative sparsity of our training data.  
 243 We therefore restrict the form of the joint distribution and construct an approximate joint probability  
 244 of enhancer-gene contacts. Pairwise correlations provide a valid likelihood if we restrict our model to  
 245 consider one gene per enhancer.

### 246 **a) The joint distribution factorises**

We assume that the likelihood  $P(\mathbf{D}_j | \mathbf{I}_j)$  can be factorised and written in the form:

$$P(\mathbf{D}_j | \mathbf{I}_j) = \prod_{\{k: I_{j,k}=1\}} P(D_{j,k} | I_{j,k} = 1) \prod_{\{k: I_{j,k}=0\}} P(D_{j,k} | I_{j,k} = 0) \quad (2)$$

247 where  $\mathbf{I}_j = I_{j,1}, \dots, I_{j,K}$  and  $\mathbf{D}_j = D_{j,1}, \dots, D_{j,K}$ . Hence the distribution of each  $D_{j,k}$  is conditionally  
 248 independent of other allocations and conditional only on the indicator variable  $I_{j,k}$ .

249 **b) An enhancer regulates a single gene**

We assume further, that an enhancer  $j$  can interact with only one gene  $k$ . We restrict the event space of  $P(\mathbf{D}_j, \mathbf{I}_j)$  to its subspace  $P(\mathbf{D}_j, \mathbf{I}_{j,k}^{(1)})$ , where  $\mathbf{I}_{j,k}^{(1)} = 0, \dots, 1, \dots, 0$ . From (2) the events are given by:

$$P(\mathbf{D}_j | \mathbf{I}_{j,k}^{(1)} = 0, \dots, 1, \dots, 0) = P(\mathbf{D}_{j,k} | I_{j,k} = 1) \prod_{\{l:l \neq k\}} P(\mathbf{D}_{j,l} | I_{j,l} = 0). \quad (3)$$

The prior distribution  $P(\mathbf{I}_j)$  follows a multivariate Bernoulli distribution, and thus the restriction is equivalent to setting the probabilities of all the structures  $\mathbf{I}_j$  with non-singular number of contacts i.e.  $\mathbf{I}_j^{(2)}, \mathbf{I}_j^{(3)}, \dots, \mathbf{I}_j^{(K)}$  to zero. For the remaining  $\mathbf{I}_{j,k}^{(1)}$  we assume that the prior is uniform across these sparse vectors, i.e.

$$P(\mathbf{I}_{j,k}^{(1)} = 0, \dots, 1, \dots, 0) = 1/K, \quad (4)$$

250 so that each  $\mathbf{I}_{j,k}^{(1)}$  is equally likely *a priori*.

251 **c) The distribution of attributes is independent**

Assuming that the attributes are conditionally independent, the likelihood component  $P(\mathbf{D}_{j,k} | I_{j,k})$  becomes:

$$P(\mathbf{D}_{j,k} | I_{j,k}) = P(d_{j,k}, c_{j,k,1}, \dots, c_{j,k,N} | I_{j,k}) = P(d_{j,k} | I_{j,k}) \prod_{n \in N} P(c_{j,k,n} | I_{j,k}) \quad (5)$$

252 where  $d_{j,k}$  is a distance from the centre of an enhancer  $j$  to the TSS of a gene  $k$ , whereas  $c_{j,k,n}$  is a  
253 correlation between the time series of the  $n^{th}$  time course dataset between an enhancer  $j$  and gene  $k$ .

Combining the assumption of the factorisable likelihood (2) with the conditional independence of attributes (5) yields,

$$P(\mathbf{D}_j | \mathbf{I}_j) = \prod_{k=1}^K P(\mathbf{D}_{j,k} | I_{j,k}) = \prod_{k=1}^K \left[ P(d_{j,k} | I_{j,k}) \prod_{n \in N} P(c_{j,k,n} | I_{j,k}) \right]. \quad (6)$$

Restricting the event space to single enhancer-gene events (3) results in,

$$P(\mathbf{D}_j | \mathbf{I}_{j,k}^{(1)}) = \left[ P(d_{j,k} | I_{j,k} = 1) \prod_{n \in N} P(c_{j,k,n} | I_{j,k} = 1) \right] \prod_{\{l:l \neq k\}} \left[ P(d_{j,l} | I_{j,l} = 0) \prod_{n \in N} P(c_{j,l,n} | I_{j,l} = 0) \right]. \quad (7)$$

254 The assumption of conditional independence of features in (5) and the fact that each vector  $\mathbf{I}_{j,k}^{(1)}$  is a 1-of-K  
255 (i.e one-to-one relation) representation of  $K$  class indicators makes this algorithm a special case of Naive  
256 Bayes (NB) model.

257 **Posterior**

The posterior distribution under the model is:

$$P(\mathbf{I}_{j,k}^{(1)} | \mathbf{D}_j) = \frac{P(\mathbf{D}_j | \mathbf{I}_{j,k}^{(1)}) P(\mathbf{I}_{j,k}^{(1)})}{\sum_{k=1}^K P(\mathbf{D}_j | \mathbf{I}_{j,k}^{(1)}) P(\mathbf{I}_{j,k}^{(1)})}. \quad (8)$$

258 The posterior distribution can be used to find the probability of each structure  $\mathbf{I}_{j,k}^{(1)}$  given the pair-wise  
259 comparisons in  $\mathbf{D}_j$ , i.e. the values of the data-specific correlations and distance for each pair  $(E_j, G_k)$  and  
260 all complementary pairs  $(E_j, G_{\{l:l \neq k\}})$ . The posterior probabilities can be used to infer the most likely  
261 target of an enhancer  $j$  out of  $K$  genes.

262 **Positive set of interactions and background negatives**

263 We overlap the distal enhancers and promoter-extended-genes with the combined set of ChIA-PET  
264 predicted links using both ER- $\alpha$  and Pol II antibodies from ENCODE/GIS-Ruan (Li et al., 2012)[GEO  
265 accession numbers GSM970209 and GSM970212]. The overall design and processing of the datasets is  
266 described under GEO accession number GSE39495. The sources contain the high-confidence binding



267 sites and protein-mediated chromatin interactions with 3 and 4 replicates for ChIA-PET with antibodies  
 268 for ER- $\alpha$  and Pol II respectively. Overlapping the enhancers and genes with the concatenated set of  
 269 empirically confirmed interactions revealed a total of 2733 enhancer-promoter links, and shows that 2087  
 270 of our distal enhancers interact with at least one promoter.

271 To define the negative set, we restricted ourselves to all enhancer-gene pairs involving known interact-  
 272 ing enhancers coming from the positive set and all the remaining non-targeted genes. Enhancers without  
 273 any confirmed interactions from ChiA-PET data were not used for training as we have no information  
 274 about their target genes.

### 275 **Data features and their distributions**

276 The method uses five features of two types, i.e. four correlations and one distance. To obtain the first  
 277 four we correlated ChIP-seq time series at enhancers with those at promoter-extended genes, for each  
 278 dataset, for all enhancer-gene pairs in the positive and negative set (as defined above). For Pol II we  
 279 used the average correlation across the two replicates. For the distance feature we used the  $\log_{10}$  of  
 280 genomic distance between the centre of the enhancer and the canonical TSS of an extended gene. We used  
 281 the training set to estimate the distributions  $P(c_{j,k,n}|I_{j,k})$  and  $P(d_{j,k}|I_{j,k})$  using kernel density estimation  
 282 (KDE) with a Gaussian kernel. To ensure that the bandwidths of positive distributions are biologically  
 283 meaningful and robust, we used cross-validation. As part of the approach, we sequentially removed all  
 284 features of each chromosome from their total set across all chromosomes and at each time calculated  
 285 the log-likelihood of KDE for the reduced set of features. We then used the value of the bandwidth with  
 286 the highest log-likelihood over left-out data. In contrast, due to a large number of negative examples  
 287 and computational cost associated with KDE, employing the same approach for negatives was infeasible.  
 288 Their size, however, also entails less requirement for optimised fitting, and thus to select the bandwidth  
 289 we resorted to the Scott's rule (Scott, 2015).

### 290 **Model Validation**

291 We trained the classifier on the odd chromosomes and estimated the training error. Similarly, we tested  
 292 the method on the even chromosomes and obtained the test error. Since the test data is not used to build  
 293 the classifier (i.e. fit the feature densities), its predictions on the test data can be considered unbiased.  
 294 We measured the performance in two ways. Firstly, we evaluated and plotted precisions against the True  
 295 Positive Rate (TPR or recall) of 10%, 20%, and 30% for various combinations of features. Secondly, we  
 296 used an alternative MAP measure. Under our model each enhancer possesses a maximum a posteriori  
 297 (MAP) gene which is our best guess of enhancer's target. The MAP measure is the percentage of times  
 298 the MAP inferred target gene is confirmed by the positive set of interactions in the ChIA-PET data.

### 299 **Performance within and outside TADs**

300 We stratified our predicted interactions at 10%, 20%, and 30% thresholds into those that lie within  
 301 domains and those that crossed domain boundaries. Each TPR threshold maps to a subsets of negative  
 302 and positive links, and therefore each subset was partitioned into inter- and intra- domain interactions. We  
 303 then tested precisions for each of the subsets. For details of TAD preparation refer to the Supplementary  
 304 Material (suppl: Domains conserved between mESC, mouse Cortex, hESC and IMR90 converted from  
 305 hg18 to hg19 using <http://www.ncbi.nlm.nih.gov/genome/tools/remap>)

### 306 **Prediction of target genes**

We used our model to infer gene targets with strong evidence of being regulated by at least one enhancer.  
 The probability of gene  $k$  having at least one active regulatory link from an enhancer under our model is  
 defined,

$$P(\text{card}(\{j \in J : I_{j,k} = 1\}) > 0) = 1 - \prod_{\{j \in J : I_{j,k} = 1\}} (1 - P(I_{j,k}^{(1)} | \mathbf{D}_j)) \quad (9)$$

307 where the product above is equal to the probability that no enhancers regulate the gene.

308 Hah et al. (2011) carried out GRO-Seq experiments (GEO accession number GSM678536) to detect  
 309 whether Pol II molecules are engaged in transcription at the start of the experiment. The experiments  
 310 were performed with the same cell-line and stimulation as ours and were used to determine the early  
 311 transcriptional response of genes following E2 treatment. Using these data and the regulation probability

312 scores defined in Eqn. (9), we assessed how many of our predicted distally regulated genes were differ-  
313 entially expressed at early time points. Using the EdgeR processed GRO-seq data we filtered the GRO-seq  
314 determined DE genes at 10, 40, 160 min after E2 stimulation with q-value (multiple hypotheses testing  
315 adjusted p-values from EdgeR) of less than 0.05, 0.01, 0.001. For each q-value, we combined the DE  
316 genes from each of the time points into a single list.

## 317 **RESULTS AND DISCUSSION**

318 We demonstrate our method using ChIP-Seq time course data collected from the MCF7 breast cancer  
319 cell-line stimulated by estrogen. After stimulation, the ER- $\alpha$  TF associates with numerous enhancers to  
320 regulate transcription of target genes. ER- $\alpha$ , encoded by the ESR1 gene, is a particularly well studied  
321 example of a nuclear receptor due to its role in breast cancer development. Its genome-wide binding  
322 pattern under stimulation with estrogen has been established through ChIP-seq experiments (Liu and  
323 Cheung, 2014; Magnani and Lupien, 2014; Ross-Innes et al., 2012). Here, the genome-wide occupancy of  
324 ER- $\alpha$  along with RNA polymerase (Pol II) and two histone marks (H3K4me3 and H2AZ) associated with  
325 transcriptional competence, were measured via ChIP-seq at eight consecutive time-points after exposure  
326 of cells in estrogen free media to estradiol. ChIA-PET data are also available in this system and were  
327 used to evaluate our method's performance (Fullwood et al., 2009; Li et al., 2010, 2012).

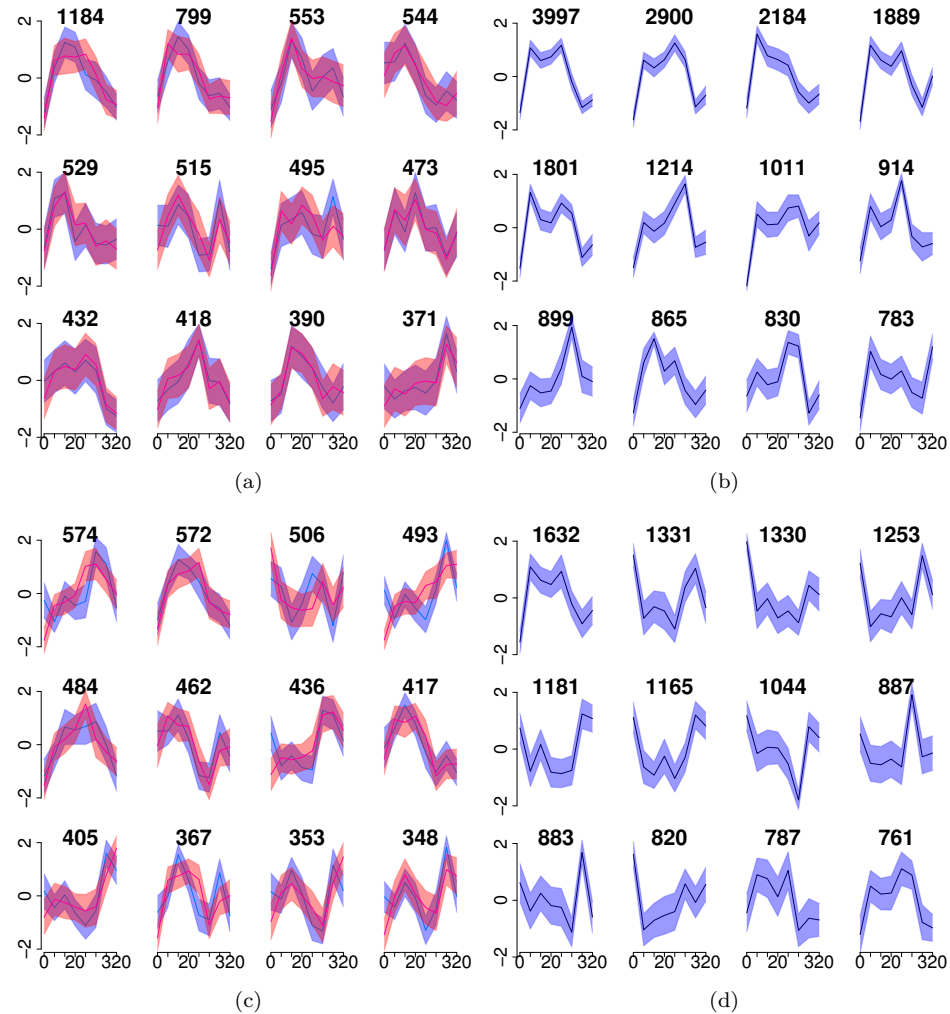
### 328 **ER- $\alpha$ bound enhancers overlap experimentally determined promoter interaction regions**

329 To locate binding events formed after stimulation with estradiol, we determined a set of genomic loci  
330 associated with ER- $\alpha$  in at least two time points. Among these 47921 regions, 21336 overlapped with  
331 a known gene or within a 300bp region upstream from its TSS (promoter-extended gene region) while  
332 26585 were distant from genes (distal enhancers).

333 Next, we determined how many of our distal ER- $\alpha$ -bound enhancers are known to form links with  
334 promoter-extended genes. Overlapping regions with interactions derived from two public ChIA-PET  
335 datasets that used the same ER- $\alpha$  and Pol II antibodies revealed a total of 2733 enhancer-promoter links.  
336 These interactions were used as a positive set for the purpose of developing our classifier. Missing  
337 interactions involving the same enhancers and other promoters in the same chromosome were used as  
338 the negative set. When training and testing the classifier, we did not include enhancers that did not have  
339 any interactions according to the ChIA-PET data. These enhancers are most likely not detected by the  
340 ChIA-PET method due to its limited sensitivity and their inclusion would introduce many false negatives  
341 into our training and testing data. However, we apply the classifier to all enhancers when making target  
342 gene predictions.

### 343 **ChIP-seq time series data**

344 We calculated the number of mapped reads for each of our ChIP-seq datasets over promoter-extended-gene  
345 bodies and over our consensus ER- $\alpha$  binding sites to create time series data for genes and enhancers (see  
346 Materials and Methods). We clustered the ER- $\alpha$  and Pol II data to help visualise the occupancy dynamics  
347 at enhancers and genes. As shown in Fig. 1, the clusters show substantial differences in occupancy  
348 dynamics across both genes and enhancers. This is expected for Pol II which shows a broad range of  
349 response profiles in this system (Honkela et al., 2015). Additionally, some differences in ER- $\alpha$  profiles  
350 were also detected, suggesting that occupancy is not solely determined by the nuclear concentration of  
351 ER- $\alpha$ .



**Figure 1.** ChIP-seq time course data show a variety of dynamic profiles which are exploited by our classifier. (a, c) show profiles of the first (blue) and the second (magenta) replicate of Pol-II for enhancers and genes, respectively. (c, d) show profiles of ER- $\alpha$  for enhancers and genes, respectively. X-axis shows time, Y-axis shows  $\pm$  one standard deviation of z-scores in each cluster. The headers show the number of time series in each cluster.

352 **Time series correlation and distance-based features are informative about enhancer-**  
 353 **promoter interactions**

354 We calculated the Pearson correlation coefficient between enhancer and gene time series data for every  
 355 enhancer-promoter pair in the positive and negative set. Figure 2 shows the distribution of correlations for  
 356 each dataset in our training data (odd chromosomes). The distribution for positive interactions differs  
 357 substantially from the background for all four datasets, with interacting regions more highly correlated  
 358 on average. This difference is most pronounced for ER- $\alpha$  and Pol II (Fig. 2a and Fig. 2b) while there  
 359 is a much smaller difference for the histone marks H2AZ and H3K4me3 (Fig. 2c and Fig. 2d). We  
 360 also compare the distribution of genomic separation for interacting and non-interacting promoters and  
 361 enhancers in Fig. 2e. Although a highly informative feature, there is a substantial overlap in the positive  
 362 and background distance densities due to a large separation of many ER- $\alpha$  bound enhancers from their  
 363 target promoters; therefore, distance alone is insufficient for accurate prediction of interactions. We note  
 364 that our ChIA-PET data does not contain very short ChIA-PET links. Links of a size shorter than 4.5kB  
 365 are usually considered to be the result of self-ligations and are filtered out Li et al. (2010). In Figure S3 we  
 366 plotted the corresponding histograms using data from all chromosomes. We observe that the distribution

367 does not change with the addition of data from even chromosomes.

### 368 **Naive Bayes classifier performance**

369 We developed a Naive Bayes classifier which integrates several discriminative features to estimate the  
370 probability of interactions between enhancer and putative target genes. Fig. 3 shows predicted interactions  
371 with only a small number confirmed by ChIA-PET (green). Interactions are shown using different shading  
372 for classification probabilities above 0.72, 0.54, 0.49 thresholds corresponding to 0.2, 0.25, 0.3 FDR levels  
373 (posterior probabilities with the highest TPR which are associated with the selected FDRs (1-precision))  
374 estimated using the training data (combination of features: Pol II, ER, distance).

375 We evaluated classifier performance using precision-recall (PR) curves (Fig. 4a and Fig. 4b). The  
376 classifier was trained on data from odd chromosomes and the results were used to establish which  
377 combination of features is most informative. Data from even chromosomes was then used as an unbiased  
378 test set to establish the performance of the selected model and to estimate decision cut-off levels. However,  
379 we do not observe significant over-fitting, probably due to the small number of features used by the  
380 classifier. Comparison of different combinations of correlations and distance features, including distance-  
381 alone and correlation-alone variants, shows that data from ER- $\alpha$  can be combined with distance to  
382 greatly enhance predictive performance (results for all possible feature combinations are shown in the  
383 Supplementary Material) while data from Pol II provides a smaller improvement in performance. The  
384 H2AZ and H3K4me3 time course data were found to not be particularly informative, consistent with  
385 Fig. 2 which shows these histone marks to have a less pronounced difference in distribution for positive  
386 and negative links. Table 1 shows that using the probability cut-offs to infer links across 23 chromosomes  
387 our model (combination of features: PolII, ER, distance) consistently outperforms the distance-alone  
388 model in terms of the number of uncovered true links. We show that at FDR equal to 0.20 our model  
389 infers 26.7 times more interactions than predictions based on proximity alone (see Table 1). In addition to  
390 considering precision-recall curves, we also tested how often using maximum a posteriori probabilities  
391 (MAP) to link all enhancers (in the training and test data) to their most probable promoters would result in  
392 correct assignments according to the ChIA-PET data (right-most column of plots in Fig. 4a and Fig. 4b).  
393 The mean performance in the MAP case is reduced and the added value of the ChIP-Seq data relative  
394 to the proximity information is also reduced. This is because for many enhancers the ChIP-Seq data  
395 signal is relatively weak and therefore focussing on the enhancer-promoter pairs with higher classification  
396 probabilities (as in the PR curves approach) produces better quality prediction on average than when we  
397 make predictions for all enhancers.

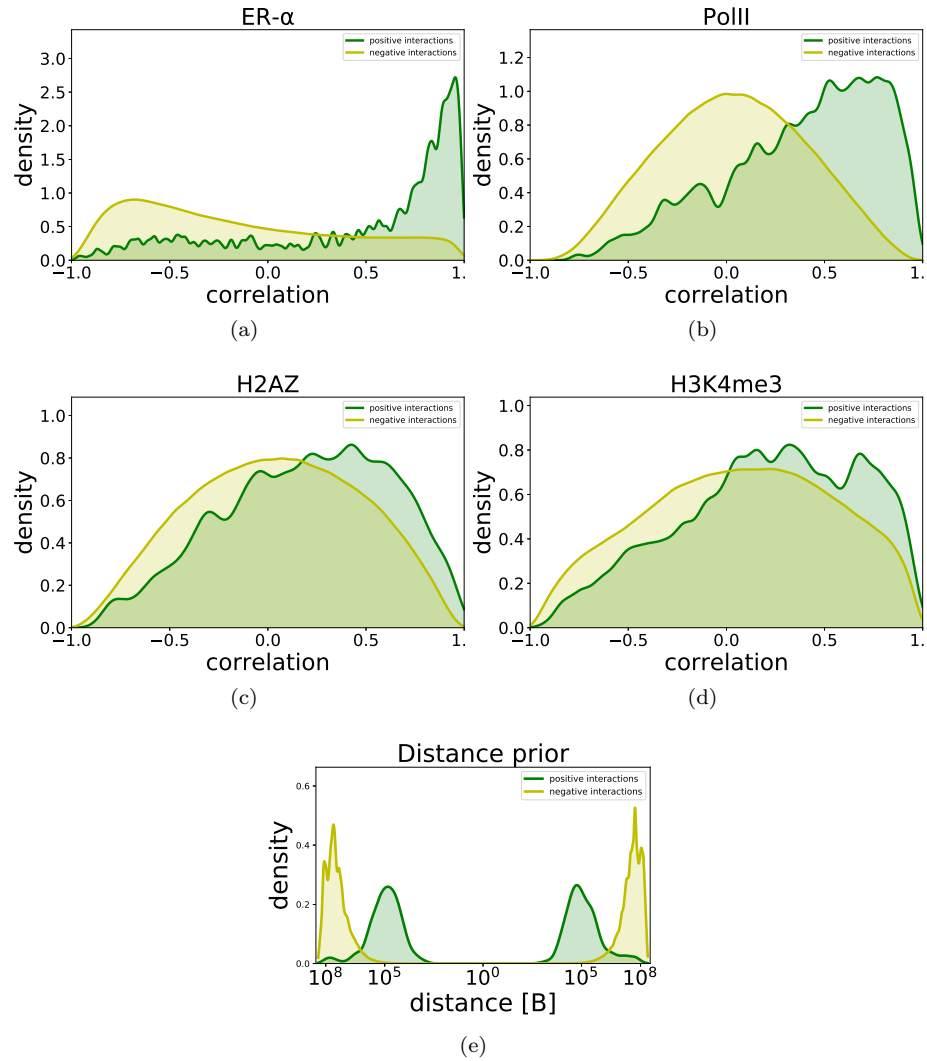
### 398 **Inter-domain and Intra-domain predictions**

399 Most enhancer-promoter interactions are thought to occur within the same Topologically Associating  
400 Domain (TAD) and we were interested in whether our method can discover interactions across TAD  
401 boundaries. In order to assess the performance of the model on discovery of intra-domain interactions and  
402 the ones involving elements from two different domains, we stratified our predicted interactions into those  
403 two groups, and recomputed precision-recall and MAP performance (Fig. 4c/d-4e/f).

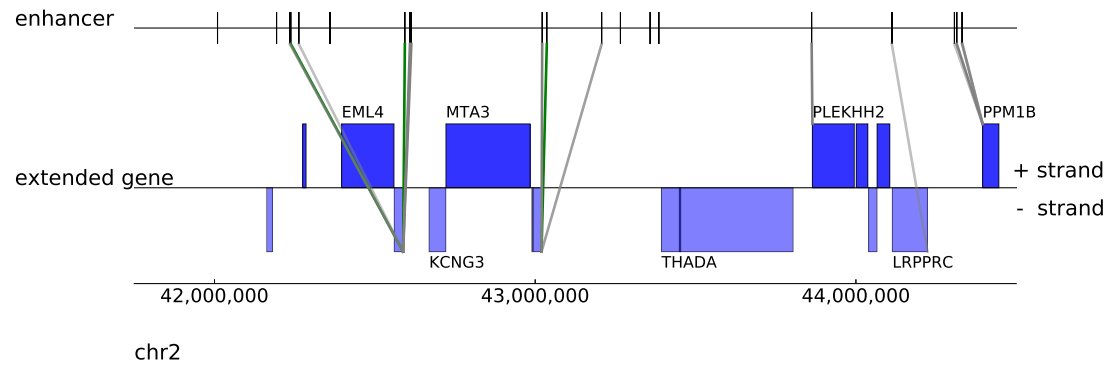
404 The majority (79%) of enhancer-promoter interactions lie within domains. The PR curves in Fig. 4d  
405 and Fig. 4e show that the ER- $\alpha$  and distance features provide the greatest contribution to performance.  
406 The Pol-II feature is also informative but does not add much to performance when combined with the  
407 ER- $\alpha$  data. Interestingly, within domains the “data-alone” model possesses much higher predictive  
408 power than in the chromosome-wide model. By excluding the possibility of long-range interactions

FDR	data/distance	distance	ratio
0.4	14217	6041	2.4
0.3	7531	1124	6.7
0.2	2800	105	26.7
0.1	109	49	2.2

**Table 1.** True links uncovered at decreasing false discovery rates for distance alone and distance assisted models.



**Figure 2.** Distribution of correlation of time series data (a,b,c,d) and genomic distance (e) for promoter-enhancer pairs and for non-interacting pairs. Here we define positive links as those confirmed by ChIA-PET experiments while negative links are defined as those not supported by ChIA-PET and involving the same set of enhancers. We observe that positive links tend to have higher correlations in the ChIP-Seq data compared to negative links, with the effect strongest for ER $\alpha$  and Pol-II.



**Figure 3.** An example of predictions with posterior probabilities above cut-off thresholds with FDR of 20%, 25%, 30% (indicated by different shades of green/gray). The green/gray colour of each link indicates whether the prediction is confirmed/unconfirmed by the ChIA-PET data.

409 beyond domain boundaries, the number of false positives is greatly reduced. Nevertheless, we see that  
 410 incorporation of the distance feature still improves classification performance within domains.

411 On the contrary (see Fig. 4e and Fig. 4f) focusing on the remaining inter-domain interactions we  
 412 notice that, in consequence of a large number of negative interactions, the correlation data alone is  
 413 insufficient for classification. The proximity data, despite being much better than the data-alone, also  
 414 does not offer the performance that we achieved for the intra-domain cases. However, distance-assisted  
 415 models perform much better than data-alone and distance-alone models and the top-ranked links have  
 416 similar precision than in the intra-domain case. Note however that the MAP results are much lower for the  
 417 inter-domain predictions, suggesting that many enhancers linking to promoters across TAD boundaries  
 418 according to the ChIA-PET data do not have this as their top-scoring interaction according to the model.

#### 419 **Testing alternative dataset design choices**

420 Our selection of data features involved some arbitrary choices and therefore we considered robustness  
 421 to varying some of the parameters used. We first investigated alternative promoter region sizes for  
 422 promoter-gene regions, their effect on test and training sets and the effect on the performance of the  
 423 model. The comparison between the distributions of features in Figures 2 and S4 and between PR curves  
 424 in Figures S5, S6 and S7, S8 show that increasing the promoter size up to 1500bp upstream from a  
 425 gene causes neither no changes to the distributions of features nor to the overall performance, and thus  
 426 the model is robust to changes in promoter region size. Similarly, Figures S9 and S10 show that using  
 427 alternative parametrisation of MACS in which we switched on  $\lambda_{local}$  parameter produces similar results to  
 428 our default parametrisation where we switched that parameter off. Figure S11 shows that the distributions  
 429 of features remain similarly unchanged.

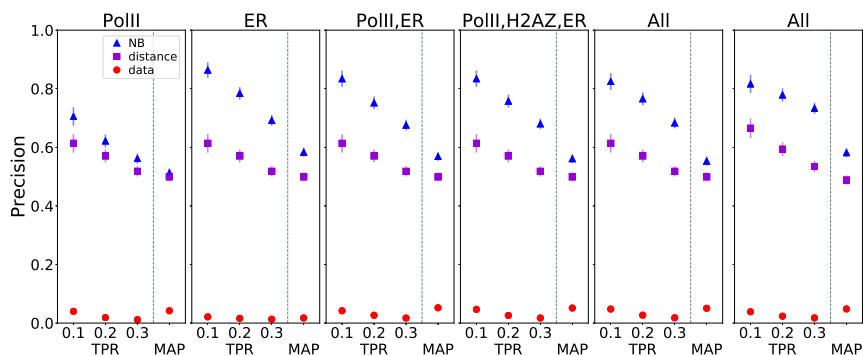
#### 430 **Validation of ER-regulated target gene predictions**

431 Finally, we used our method to provide a highly confident (FDR of 0.25) list of directly ER-regulated  
 432 target genes in this system. This list (Table S1) includes 1978 genes with at least one predicted enhancer  
 433 link. In Fig. 5 we compared our set of predicted distally regulated genes against a list of early differentially  
 434 expressed genes obtained from GRO-seq experiments (Hah et al., 2011). PR curves showed that the larger  
 435 the value of the score (see Materials and Methods), which is roughly proportional to the number of times  
 436 a gene is predicted to be a target of distal enhancer, the higher the chance that the gene is differentially  
 437 expressed. Using a score based only on proximity of ER- $\alpha$  binding events is much less predictive of early  
 438 differential expression.

## 439 **CONCLUSIONS**

440 We have developed a Bayesian method which is capable of integrating genomic distance with a correlation  
 441 of ChIP-seq time series in order to predict physical interactions between enhancers and promoters.  
 442 We evaluated the performance of our method against ChIA-PET predicted links and using different  
 443 combinations of features. Using complementary GRO-seq data from the same cell-line and experimental

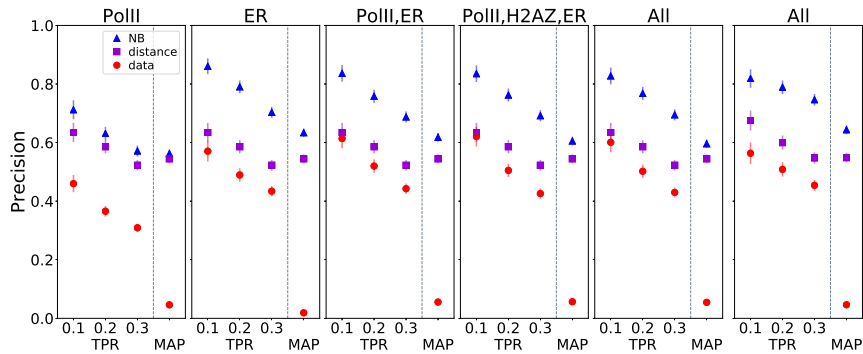
### Overall performance



(a) training

(b) test

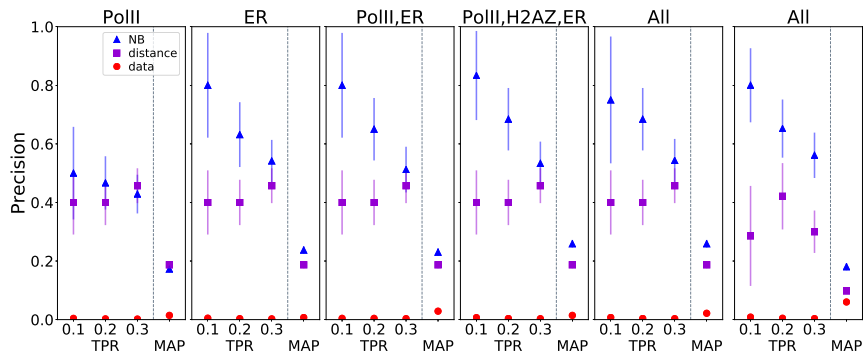
### Intra-domain Interactions



(c) training

(d) test

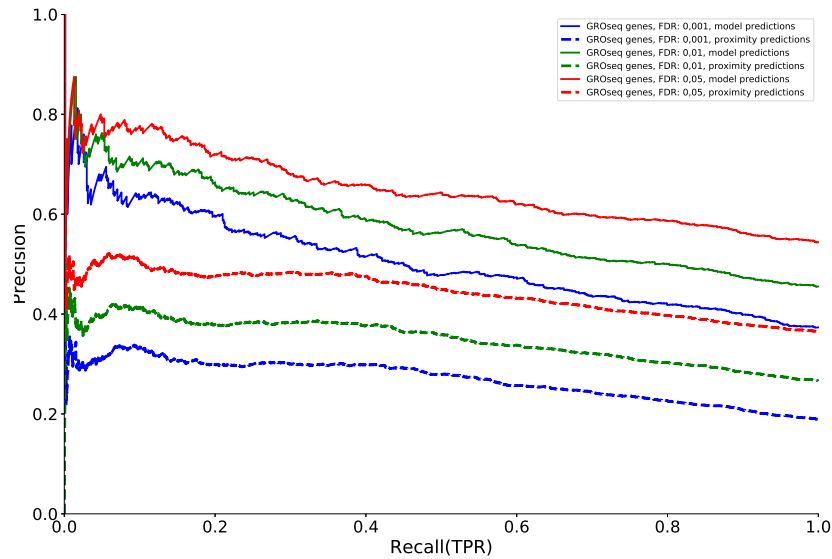
### Inter-domain Interactions



(e) training

(f) test

**Figure 4.** Graphs (a, b, c, d, e, f) show the performance of the model, measured by Precision-Recall and MAP scores. The precisions are plotted againsts TPR of 0.1, 0.2, 0.3. Each column shows the performance of the model with a variant of correlation-based feature/s (i.e. data, see header) and proximity-based feature (i.e. distance, see header). The first five columns of each row show the performance on the training data. The last column shows the performance on the test data.



**Figure 5.** The Precision-Recall curves assess the model on the ability to predict differentially expressed genes (as derived from GRO-Seq data), given a number of model-assigned regulators of each gene and the confidence of each prediction.

444 context we show that our model can accurately predict distally regulated, differentially expressed genes  
 445 under stimulation with estrogen. Our model can therefore serve as a complementary approach to  
 446 chromosome conformation capture techniques and offers insight into context-specific, and cell-type  
 447 specific transcriptional regulation.

## 448 REFERENCES

- 449 Andersson, R., Gebhard, C., Miguel-Escalada, I., Hoof, I., Bornholdt, J., Boyd, M., Chen, Y., Zhao, X.,  
 450 Schmidl, C., Suzuki, T., Ntini, E., Arner, E., Valen, E., Li, K., Schwarzfischer, L., Glatz, D., Raithel, J.,  
 451 Lilje, B., Rapin, N., Bagger, F. O., Jørgensen, M., Andersen, P. R., Bertin, N., Rackham, O., Burroughs,  
 452 a. M., Baillie, J. K., Ishizu, Y., Shimizu, Y., Furuhashi, E., Maeda, S., Negishi, Y., Mungall, C. J.,  
 453 Meehan, T. F., Lassmann, T., Itoh, M., Kawaji, H., Kondo, N., Kawai, J., Lennartsson, A., Daub, C. O.,  
 454 Heutink, P., Hume, D. a., Jensen, T. H., Suzuki, H., Hayashizaki, Y., Müller, F., Forrest, A. R. R.,  
 455 Carninci, P., Rehli, M., and Sandelin, A. (2014). An atlas of active enhancers across human cell types  
 456 and tissues. *Nature*, 507:455–461.
- 457 Bannister, A. J. and Kouzarides, T. (2011). Regulation of chromatin by histone modifications. *Cell*  
 458 *research*, 21(3):381–395.
- 459 Bernstein, B. E., Kamal, M., Lindblad-Toh, K., Bekiranov, S., Bailey, D. K., Huebert, D. J., McMahon, S.,  
 460 Karlsson, E. K., Kulbokas, E. J., Gingeras, T. R., Schreiber, S. L., and Lander, E. S. (2005). Genomic  
 461 maps and comparative analysis of histone modifications in human and mouse. *Cell*, 120(2):169–181.
- 462 Corradin, O., Saiakhova, A., Akhtar-Zaidi, B., Myeroff, L., Willis, J., Cowper-Sallari, R., Lupien, M.,  
 463 Markowitz, S., and Scacheri, P. C. (2014). Combinatorial effects of multiple enhancer variants in  
 464 linkage disequilibrium dictate levels of gene expression to confer susceptibility to common traits.  
 465 *Genome Research*, 24:1–13.
- 466 Dekker, J., Rippe, K., Dekker, M., and Kleckner, N. (2002). Capturing chromosome conformation.  
 467 *Science (New York, N.Y.)*, 295(5558):1306–1311.
- 468 Dixon, J. R., Selvaraj, S., Yue, F., Kim, A., Li, Y., Shen, Y., Hu, M., Liu, J. S., and Ren, B. (2012).  
 469 Topological domains in mammalian genomes identified by analysis of chromatin interactions. *Nature*,  
 470 485(7398):376–380.
- 471 Dostie, J., Richmond, T. a., Arnaout, R. a., Selzer, R. R., Lee, W. L., Honan, T. a., Rubio, E. D., Krumm,  
 472 A., Lamb, J., Nusbaum, C., Green, R. D., and Dekker, J. (2006). Chromosome Conformation Capture



473 Carbon Copy (5C): A massively parallel solution for mapping interactions between genomic elements.  
474 *Genome Research*, 16(10):1299–1309.

475 Drissen, R., Drissen, R., Palstra, R.-j., Palstra, R.-j., Gillemans, N., Gillemans, N., Splinter, E., Splinter,  
476 E., Grosveld, F., Grosveld, F., Philipsen, S., Philipsen, S., Laat, W. D., and Laat, W. D. (2004). The  
477 active spatial organization of the. *Genes & Development*, pages 2485–2490.

478 Ernst, J., Kheradpour, P., Mikkelson, T. S., Shores, N., Ward, L. D., Epstein, C. B., Zhang, X., Wang,  
479 L., Issner, R., Coyne, M., Ku, M., Durham, T., Kellis, M., and Bernstein, B. E. (2011). Mapping and  
480 analysis of chromatin state dynamics in nine human cell types. *Nature*, 473(7345):43–49.

481 Frey, B. and Dueck, D. (2007). Clustering by passing messages between data points. *science*,  
482 315(February):972–977.

483 Fullwood, M. J., Liu, M. H., Pan, Y. F., Liu, J., Xu, H., Mohamed, Y. B., Orlov, Y. L., Velkov, S., Ho, A.,  
484 Mei, P. H., Chew, E. G. Y., Huang, P. Y. H., Welboren, W.-J., Han, Y., Ooi, H. S., Ariyaratne, P. N.,  
485 Vega, V. B., Luo, Y., Tan, P. Y., Choy, P. Y., Wansa, K. D. S. A., Zhao, B., Lim, K. S., Leow, S. C.,  
486 Yow, J. S., Joseph, R., Li, H., Desai, K. V., Thomsen, J. S., Lee, Y. K., Karuturi, R. K. M., Herve, T.,  
487 Bourque, G., Stunnenberg, H. G., Ruan, X., Cacheux-Rataboul, V., Sung, W.-K., Liu, E. T., Wei, C.-L.,  
488 Cheung, E., and Ruan, Y. (2009). An oestrogen-receptor-alpha-bound human chromatin interactome.  
489 *Nature*, 462(7269):58–64.

490 Hagège, H., Klous, P., Braem, C., Splinter, E., Dekker, J., Cathala, G., de Laat, W., and Forné, T. (2007).  
491 Quantitative analysis of chromosome conformation capture assays (3C-qPCR). *Nature protocols*,  
492 2(7):1722–33.

493 Hah, N., Danko, C., Core, L., Waterfall, J., Siepel, A., Lis, J., and Kraus, W. (2011). A Rapid, Extensive,  
494 and Transient Transcriptional Response to Estrogen Signaling in Breast Cancer Cells. *Cell*, 145(7):1156.

495 He, B., Chen, C., Teng, L., and Tan, K. (2014). Global view of enhancer-promoter interactome in  
496 human cells. *Proceedings of the National Academy of Sciences of the United States of America*,  
497 111:E2191—9.

498 Honkela, A., Peltonen, J., Topa, H., Charapitsa, I., Matarese, F., Grote, K., Stunnenberg, H. G., Reid, G.,  
499 Lawrence, N. D., and Rattray, M. (2015). Genome-wide modeling of transcription kinetics reveals  
500 patterns of RNA production delays. *Proceedings of the National Academy of Sciences*, page 201420404.

501 Javierre, B. M., Sewitz, S., Cairns, J., Wingett, S. W., V??rnai, C., Thiecke, M. J., Freire-Pritchett, P.,  
502 Spivakov, M., Fraser, P., Burren, O. S., Cutler, A. J., Todd, J. A., Wallace, C., Wilder, S. P., Kreuzhuber,  
503 R., Kostadima, M., Zerbino, D. R., Stegle, O., Kreuzhuber, R., Burden, F., Farrow, S., Rehnstr??m,  
504 K., Downes, K., Grassi, L., Kostadima, M., Ouwehand, W. H., Frontini, M., Kreuzhuber, R., Burden,  
505 F., Farrow, S., Rehnstr??m, K., Downes, K., Kostadima, M., Ouwehand, W. H., Frontini, M., Hill,  
506 S. M., Wang, F., Wallace, C., Stunnenberg, H. G., Ouwehand, W. H., Frontini, M., Ouwehand, W. H.,  
507 Wallace, C., Martens, J. H., Kim, B., Sharifi, N., Janssen-Megens, E. M., Yaspo, M. L., Linsler, M.,  
508 Kovacovics, A., Clarke, L., Richardson, D., Datta, A., and Flicek, P. (2016). Lineage-Specific Genome  
509 Architecture Links Enhancers and Non-coding Disease Variants to Target Gene Promoters. *Cell*,  
510 167(5):1369–1384.e19.

511 Jin, F., Li, Y., Dixon, J. R., Selvaraj, S., Ye, Z., Lee, A. Y., Yen, C.-A., Schmitt, A. D., Espinoza, C. a.,  
512 and Ren, B. (2013). A high-resolution map of the three-dimensional chromatin interactome in human  
513 cells. *Nature*, 503(7475):290–294.

514 Li, G., Fullwood, M. J., Xu, H., Mulawadi, F. H., Velkov, S., Vega, V., Ariyaratne, P. N., Mohamed,  
515 Y. B., Ooi, H.-S., Tennakoon, C., Wei, C.-L., Ruan, Y., and Sung, W.-K. (2010). ChIA-PET tool  
516 for comprehensive chromatin interaction analysis with paired-end tag sequencing. *Genome biology*,  
517 11(2):R22.

518 Li, G., Ruan, X., Auerbach, R. K., Sandhu, K. S., Zheng, M., Wang, P., Poh, H. M., Goh, Y., Lim, J.,  
519 Zhang, J., Sim, H. S., Peh, S. Q., Mulawadi, F. H., Ong, C. T., Orlov, Y. L., Hong, S., Zhang, Z.,  
520 Landt, S., Raha, D., Euskirchen, G., Wei, C.-L., Ge, W., Wang, H., Davis, C., Fisher-Aylor, K. I.,  
521 Mortazavi, A., Gerstein, M., Gingeras, T., Wold, B., Sun, Y., Fullwood, M. J., Cheung, E., Liu, E.,  
522 Sung, W.-K., Snyder, M., and Ruan, Y. (2012). Extensive promoter-centered chromatin interactions  
523 provide a topological basis for transcription regulation. *Cell*, 148(1-2):84–98.

524 Liu, M. H. and Cheung, E. (2014). Estrogen receptor-mediated long-range chromatin interactions and  
525 transcription in breast cancer. *Molecular and cellular endocrinology*, 382(1):624–632.

526 Magnani, L. and Lupien, M. (2014). Chromatin and epigenetic determinants of estrogen receptor alpha  
527 (ESR1) signaling. *Molecular and cellular endocrinology*, 382(1):633–641.

528 Marstrand, T. T. and Storey, J. D. (2014). Identifying and mapping cell-type-specific chromatin program-  
529 ming of gene expression. *Proceedings of the National Academy of Sciences of the United States of*  
530 *America*, 111:E645—54.

531 Mifsud, B., Tavares-Cadete, F., Young, A. N., Sugar, R., Schoenfelder, S., Ferreira, L., Wingett, S. W.,  
532 Andrews, S., Grey, W., Ewels, P. a., Herman, B., Happe, S., Higgs, A., LeProust, E., Follows, G. a.,  
533 Fraser, P., Luscombe, N. M., and Osborne, C. S. (2015). Mapping long-range promoter contacts in  
534 human cells with high-resolution capture Hi-C. *Nature Genetics*, 47(6):598–606.

535 Nagano, T., Lubling, Y., Stevens, T. J., Schoenfelder, S., Yaffe, E., Dean, W., Laue, E. D., Tanay, A., and  
536 Fraser, P. (2013). Single-cell Hi-C reveals cell-to-cell variability in chromosome structure. *Nature*,  
537 502(7469):59–64.

538 Nagarajan, S., Hossan, T., Alawi, M., Najafova, Z., Indenbirken, D., Bedi, U., Taipaleenmäki, H., Ben-  
539 Batalla, I., Scheller, M., Loges, S., Knapp, S., Hesse, E., Chiang, C.-M., Grundhoff, A., and Johnsen,  
540 S. a. (2014). Bromodomain protein BRD4 is required for estrogen receptor-dependent enhancer  
541 activation and gene transcription. *Cell reports*, 8(2):460–469.

542 Quinlan, A. R. and Hall, I. M. (2010). BEDTools: A flexible suite of utilities for comparing genomic  
543 features. *Bioinformatics*, 26(6):841–842.

544 Ross-Innes, C. S., Stark, R., Teschendorff, A. E., Holmes, K. a., Ali, H. R., Dunning, M. J., Brown, G. D.,  
545 Gojis, O., Ellis, I. O., Green, A. R., Ali, S., Chin, S.-F., Palmieri, C., Caldas, C., and Carroll, J. S.  
546 (2012). Differential oestrogen receptor binding is associated with clinical outcome in breast cancer.  
547 *Nature*, 481(7381):389–393.

548 Roy, S., Siahpirani, A. F., Chasman, D., Knaack, S., Ay, F., Stewart, R., Wilson, M., and Sridharan,  
549 R. (2015). A predictive modeling approach for cell line-specific long-range regulatory interactions.  
550 *Nucleic Acids Research*, 43(18):gkv865.

551 Sanyal, A., Lajoie, B. R., Jain, G., and Dekker, J. (2012). The long-range interaction landscape of gene  
552 promoters. *Nature*, 489(7414):109–113.

553 Schoenfelder, S., Sexton, T., Chakalova, L., Cope, N. F., Horton, A., Andrews, S., Kurukuti, S., Mitchell,  
554 J. a., Umlauf, D., Dimitrova, D. S., Eskiw, C. H., Luo, Y., Wei, C.-L., Ruan, Y., Bieker, J. J., and Fraser,  
555 P. (2010). Preferential associations between co-regulated genes reveal a transcriptional interactome in  
556 erythroid cells. *Nature genetics*, 42(1):53–61.

557 Scott, D. W. (2015). *Multivariate density estimation: theory, practice, and visualization*. John Wiley &  
558 Sons.

559 Shen, Y., Yue, F., McCleary, D. F., Ye, Z., Edsall, L., Kuan, S., Wagner, U., Dixon, J., Lee, L., Lobanenko,  
560 V. V., and Ren, B. (2012). A map of the cis-regulatory sequences in the mouse genome. *Nature*,  
561 488(7409):116–120.

562 Shlyueva, D., Stampfel, G., and Stark, A. (2014). Transcriptional enhancers: from properties to genome-  
563 wide predictions. *Nature reviews. Genetics*, 15(4):272–86.

564 Simonis, M., Kooren, J., and de Laat, W. (2007). An evaluation of 3C-based methods to capture DNA  
565 interactions. *Nature methods*, 4(11):895–901.

566 Stasevich, T. J., Hayashi-Takanaka, Y., Sato, Y., Maehara, K., Ohkawa, Y., Sakata-Sogawa, K., Tokunaga,  
567 M., Nagase, T., Nozaki, N., McNally, J. G., and Kimura, H. (2014). Regulation of RNA polymerase II  
568 activation by histone acetylation in single living cells. *Nature*, 516(7530):272–275.

569 Thurman, R. E., Rynes, E., Humbert, R., Vierstra, J., Maurano, M. T., Haugen, E., Sheffield, N. C.,  
570 Stergachis, A. B., Wang, H., Vernot, B., Garg, K., John, S., Sandstrom, R., Bates, D., Boatman, L.,  
571 Canfield, T. K., Diegel, M., Dunn, D., Ebersol, A. K., Frum, T., Giste, E., Johnson, A. K., Johnson,  
572 E. M., Kutayavin, T., Lajoie, B., Lee, B.-K., Lee, K., London, D., Lotakis, D., Neph, S., Neri, F., Nguyen,  
573 E. D., Qu, H., Reynolds, A. P., Roach, V., Safi, A., Sanchez, M. E., Sanyal, A., Shafer, A., Simon, J. M.,  
574 Song, L., Vong, S., Weaver, M., Yan, Y., Zhang, Z., Zhang, Z., Lenhard, B., Tewari, M., Dorschner,  
575 M. O., Hansen, R. S., Navas, P. a., Stamatoyannopoulos, G., Iyer, V. R., Lieb, J. D., Sunyaev, S. R.,  
576 Akey, J. M., Sabo, P. J., Kaul, R., Furey, T. S., Dekker, J., Crawford, G. E., and Stamatoyannopoulos,  
577 J. a. (2012). The accessible chromatin landscape of the human genome. *Nature*, 489(7414):75–82.

578 Tolhuis, B., Palstra, R. J., Splinter, E., Grosveld, F., and De Laat, W. (2002). Looping and interaction  
579 between hypersensitive sites in the active  $\beta$ -globin locus. *Molecular Cell*, 10(6):1453–1465.

580 Vakoc, C. R., Letting, D. L., Gheldof, N., Sawado, T., Bender, M. a., Groudine, M., Weiss, M. J., Dekker,  
581 J., and Blobel, G. a. (2005). Proximity among distant regulatory elements at the  $\beta$ -globin locus requires  
582 GATA-1 and FOG-1. *Molecular Cell*, 17(3):453–462.

583 van Steensel, B. and Dekker, J. (2010). Genomics tools for unraveling chromosome architecture. *Nature*  
584 *Biotechnology*, 28(10):1089–1095.

585 wa Maina, C., Honkela, A., Matarese, F., Grote, K., Stunnenberg, H. G., Reid, G., Lawrence, N. D.,  
586 and Rattray, M. (2014). Inference of RNA Polymerase II Transcription Dynamics from Chromatin  
587 Immunoprecipitation Time Course Data. *PLoS Computational Biology*, 10(5):e1003598.

588 Zhang, Y., Liu, T., Meyer, C. A., Eeckhoute, J., Johnson, D. S., Bernstein, B. E., Nussbaum, C., Myers,  
589 R. M., Brown, M., Li, W., and Liu, X. S. (2008). Model-based Analysis of ChIP-Seq (MACS). *Genome*  
590 *Biology*, 9(9):R137.

591 Zhao, Z., Tavosoidana, G., Sjölander, M., Göndör, A., Mariano, P., Wang, S., Kanduri, C., Lezcano,  
592 M., Sandhu, K. S., Singh, U., Pant, V., Tiwari, V., Kurukuti, S., and Ohlsson, R. (2006). Circular  
593 chromosome conformation capture (4C) uncovers extensive networks of epigenetically regulated intra-  
594 and interchromosomal interactions. *Nature genetics*, 38(11):1341–7.

595 Zhu, Y., Sun, L., Chen, Z., Whitaker, J. W., Wang, T., and Wang, W. (2013). Predicting enhancer  
596 transcription and activity from chromatin modifications. *Nucleic Acids Research*, 41(22):10032–10043.