

Hi-C-LSTM: Learning representations of chromatin contacts using a recurrent neural network identifies genomic drivers of 3D genome conformation

Kevin Dsouza

University of British Columbia

Alexandra Maslova

Simon Fraser University

Ediem Al-Jibury

Imperial College London <https://orcid.org/0000-0003-2466-0989>

Matthias Merkschlager

Imperial College London

Vijay Bhargava

University of British Columbia

Maxwell Libbrecht (✉ maxwl@sfu.ca)

Simon Fraser University <https://orcid.org/0000-0003-2502-0262>

Article

Keywords: gene expression, Hi-C-LSTM, neural network

Posted Date: October 1st, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-878825/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Hi-C-LSTM: Learning representations of chromatin contacts using a recurrent neural network identifies genomic drivers of 3D genome conformation

Kevin B. Dsouza^{1*}, Alexandra Maslova², Ediem Al-Jibury^{3,4}, Matthias Merkschlager³, Vijay K. Bhargava¹, and Maxwell W. Libbrecht^{2*}

¹Department of Electrical and Computer Engineering, University of British Columbia

²School of Computing Science, Simon Fraser University

³MRC, London Institute of Medical Sciences, Institute of Clinical Sciences, Faculty of Medicine, Imperial College London

⁴Department of Computing, Imperial College London

Correspondence: *maxwl@sfu.ca*, *kevin@ece.ubc.ca*

September 5, 2021

Abstract

Despite the availability of chromatin conformation capture experiments, discerning the relationship between the 1D genome and 3D conformation remains a challenge, which limits our understanding of their affect on gene expression and disease. We propose Hi-C-LSTM, a method that produces low-dimensional latent representations that summarize intra-chromosomal Hi-C contacts via a recurrent long short-term memory (LSTM) neural network model. We find that these representations contain all the information needed to recreate the original Hi-C matrix with high accuracy, outperforming existing methods. These representations enable the identification of a variety of conformation-defining genomic elements, including nuclear compartments and conformation-related transcription factors. They furthermore enable in-silico perturbation experiments that measure the influence of cis-regulatory elements on conformation.

Introduction

The organisation of the genome in 3D space inside the nucleus is important to its function. Chromosome conformation capture (3C) techniques, developed in the last couple of decades, have enabled researchers to quantify the strength of interactions between loci that are nearby in space. Hi-C [1] uses a combination of chromatin conformation capture and high-throughput sequencing to assay pairwise chromatin interactions genome-wide. This rich source of data promises to help elucidate the influence of 3D structure on gene expression and thereby on development, evolution and disease. However, we lack a complete understanding of how the 1D genome influences 3D conformation.

The machine learning technique of representation learning [2] provides a way to link the 1D genome to 3D conformation. Representation learning aims to summarize high dimensional datasets into a low-dimensional representation. It has become a valuable tool for finding compact and informative representations that disentangle explanatory factors in diverse data types. Representation learning has recently driven advances in a variety of tasks including speech recognition [3], signal

37 processing [4], object recognition [5], natural language processing [6, 7] and domain adaptation
38 [8]. Representation learning has recently been applied to genomic sequences [9, 10] and Hi-C data
39 [11–14].

40 In order to understand the 1D-3D relationship and thereby link 3D conformation to genetic
41 variation and disease, we need representations for Hi-C data that can summarize the contact map
42 into a locus-level summary. Such a representation would encapsulate all the contacts from each
43 genomic position to the others into a small number of features per locus, such that the contacts
44 can be reproduced using just the features. Reducing the Hi-C map to locus-level representations
45 in this way would allow us to study the effect of sequence elements on chromatin conformation,
46 identify genomic drivers of 3D conformation and predict the effect of genetic variants.

47 Two methods for representation learning of Hi-C data have previously been developed, SNIPER
48 [11] and SCI [12]. SNIPER uses a fully-connected autoencoder [15] to transform the sparse Hi-C
49 inter-chromosomal matrix into a dense one row-wise, the bottleneck of which is assigned as the
50 representation for the corresponding row. SCI [12] treats the Hi-C matrix as a graph and performs
51 graph embedding [16], aiming to preserve the local and the global structures to form representations
52 for each node.

53 These existing methods for Hi-C representations have two weaknesses that limit their applicabil-
54 ity. First, SNIPER takes only inter-chromosomal contacts as input and therefore its representations
55 cannot incorporate intra-chromosomal contact patterns that are most important for the regulation
56 of gene expression, such as topological domains and promoter-enhancer looping. Second, the Hi-C
57 representations produced by both SNIPER and SCI do not account for the inherent sequential
58 nature of the genome. As we demonstrate below, these two weaknesses limits existing methods’
59 informativeness and makes them unable to accurately identify conformation-defining elements or
60 predict how those elements influence structure.

61 In this work, we propose a method called Hi-C-LSTM that produces low-dimensional representa-
62 tions of the Hi-C intra-chromosomal contacts, assigning a vector of features to each genomic position
63 that represents that position’s contact activity with all other positions in the given chromosome.
64 Hi-C-LSTM defines these representations using a sequential long short-term memory (LSTM) neu-
65 ral network model which, in contrast to existing methods like SNIPER and SCI, accounts for the
66 sequential nature of the genome. A second methodological innovation of Hi-C-LSTM is that, in-
67 stead of learning an encoder to create representations, we learn our representations directly through
68 iterative optimization. We find that this approach provides a large improvement in information
69 content relative to existing non-sequential methods, enables the use of intra-chromosomal interac-
70 tions, and enables the model to accurately predict the effects of genomic perturbations (Fig. 1,
71 Results).

72 We demonstrate the utility of Hi-C-LSTM’s representations through several analyses. First,
73 we show that our representations have information needed to recreate the Hi-C matrix and that
74 this recreation is more accurate using an LSTM than alternatives. Second, we show that our
75 representation captures cell type-specific functional activity, genomic elements and identifies ge-
76 nomic regions that drive conformation. Third, we show that feature attribution of Hi-C-LSTM can
77 identify sequence elements driving 3D conformation, such as binding sites of CTCF and cohesin
78 subunits [17, 18]. Fourth, we show Hi-C-LSTM can perform in-silico perturbation of CTCF and
79 cohesin binding sites. Fifth, we simulated a previously-assayed 2.1 Mbp structural variant at the
80 SOX9 locus and found that Hi-C-LSTM correctly reproduces experimentally-derived contacts.

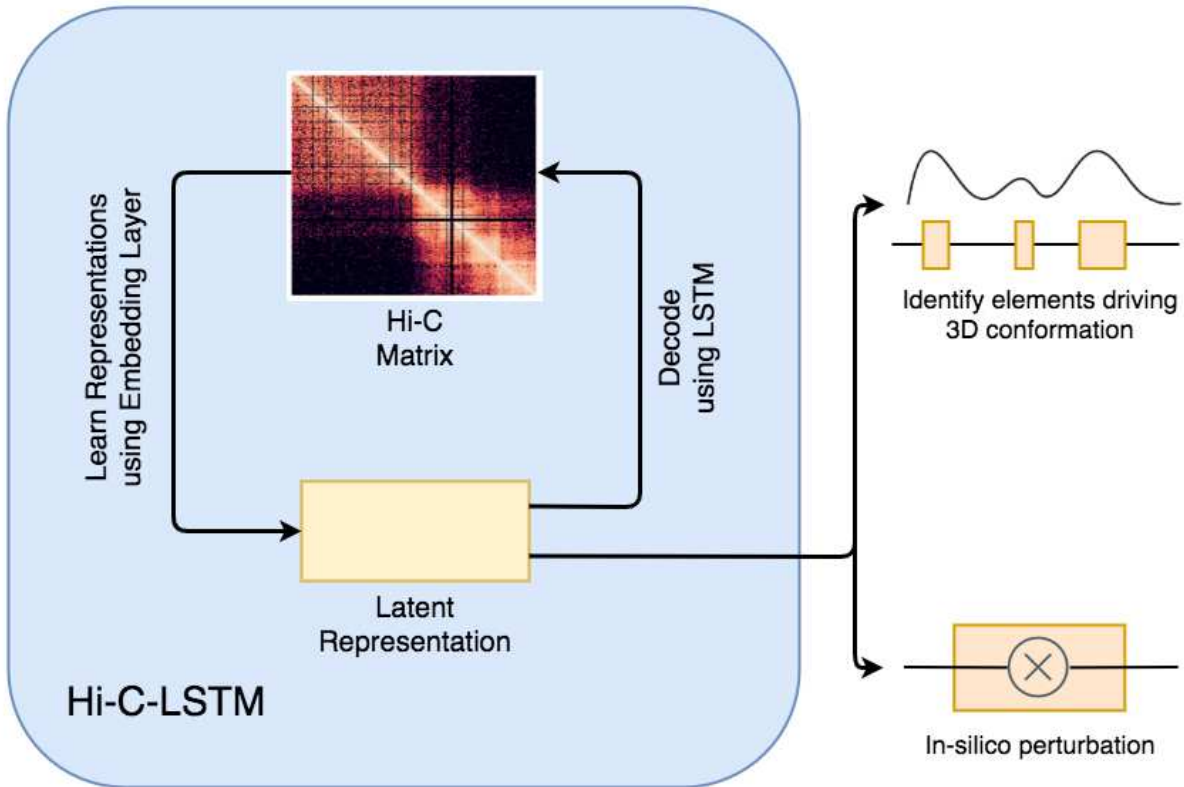


Figure 1: Overview of approach. Hi-C-LSTM learns a K -length vector representation of each genomic position that summarizes its chromatin contacts, using an LSTM embedding neural network. The representations and LSTM decoder are jointly optimized to maximize the accuracy with which the decoder can reproduce the original Hi-C matrix given just the representations.

The resulting representations identify sequence elements driving 3D conformation through Integrated Gradients (IG) analysis, and they enable a researcher to perform in-silico perturbation experiments by editing the representations and observing the effect on predicted contacts.

81 Related work

82 Hi-C-LSTM performs two main tasks; it forms Hi-C representations, and it predicts Hi-C contacts.
 83 Learning methods have been proposed that perform either of these tasks. SNIPER [11] and SCI [12]
 84 can form representations of Hi-C. SNIPER forms Hi-C representations using a feed-forward neural
 85 network autoencoder. While SNIPER predicts high-resolution Hi-C contacts using low-resolution
 86 contacts as input, Hi-C-LSTM predicts Hi-C contacts using just the genomic positions as input.
 87 SCI forms Hi-C representations by performing graph network embedding on the Hi-C data. SCI
 88 is similar to Hi-C-LSTM in that it can be used to identify elements, however, it differs in the
 89 underlying structure it uses to represent the genome. SCI represents the genome using a graph,
 90 whereas Hi-C-LSTM treats the genome as a sequence. We compare Hi-C-LSTM with these two
 91 methods as they are most similar to what we are trying to achieve.

92 The first Hi-C representations were formed using Principal component analysis (PCA) based
 93 methods, introduced in Lieberman-Aiden et al. [1]. These methods cluster the Hi-C matrix into
 94 A and B compartments based on the first principal component of the intra-chromosomal contact
 95 matrix. Imakaev et al. [19] later showed that PCA based reduction is inaccurate at classifying

96 compartments and Rao et al. [20] used a Gaussian hidden Markov model (HMM) to obtain latent
97 features that were better at locating compartments. We treat the PCA based method developed
98 in Lieberman-Aiden et al. [1] as a baseline.

99 Some methods form chromatin representations but are not directly comparable to ours. REACH-
100 3D [21] forms internal Hi-C representations using manifold learning combined with recurrent au-
101 toencoders, however, these are three dimensional and mainly used for 3D chromatin structure
102 inference. MATCHA [14] forms representations using hypergraph representation learning and uses
103 them to distinguish multi-way interactions from pairwise interaction cliques. We don't compare
104 Hi-C-LSTM with MATCHA because MATCHA works with multi-way interaction data (SPRITE
105 and ChIA-Drop) whereas we use pair-wise interaction data (Hi-C).

106 Many methods have been proposed for predicting Hi-C contacts. Some methods try to predict
107 the chromatin contacts by using either the nucleotide sequence or chromatin accessibility and
108 histone modifications or both [22–29]. Akita in particular [29], is a convolutional neural network
109 that predicts chromatin contacts from the nucleotide sequence alone, and can be used to perform in-
110 silico predictions. In addition to these, the maximum entropy genomic annotation from biomarkers
111 associated to structural ensembles (MEGABASE) coupled with an energy landscape model for
112 chromatin organization called minimal chromatin model (MiChroM), generates an ensemble of 3D
113 chromosome conformations [30]. Though these methods are similar to Hi-C-LSTM in that they
114 predict Hi-C contacts, we don't compare Hi-C-LSTM with them as none of them produce Hi-C
115 representations.

116 Results

117 Hi-C-LSTM representations capture the information needed to create the Hi-C 118 matrix

119 Hi-C-LSTM assigns a representation to each genomic position in the Hi-C contact map, such that a
120 LSTM [31] that takes these representations as input can predict the original contact map (Fig. 2).
121 The representation and the LSTM are jointly trained to optimize the reconstruction of the Hi-C
122 map. This process gives us position-specific representations genome-wide (see Methods for more
123 details).

124 We find that the Hi-C-LSTM achieves higher accuracy when constructing the Hi-C matrix com-
125 pared to existing methods (Fig. 3a). The inferred Hi-C map matches the original Hi-C map (Fig.
126 3c) closely, and differs from it by about 0.25 R-squared points on average. We adapt SNIPER to
127 our task by replacing the feed-forward decoder that converts low-resolution Hi-C to high-resolution
128 Hi-C with a decoder that reproduces the original input Hi-C. We call this SNIPER-FC. Hi-C-LSTM
129 outperforms SNIPER (SNIPER-FC) convincingly, by 10% higher R-squared on average (Fig. 3a).
130 Hi-C-LSTM also outperforms SCI (SCI-LSTM) by 12% higher R-squared on average (Fig. 3a).

131 Two hypotheses could explain Hi-C-LSTM's improved reconstructions: (1) that Hi-C-LSTM's
132 representation captures more information, or (2) that an LSTM is a more powerful decoder. We
133 found that both are true. To distinguish these hypotheses, we split each model respectively into
134 two components—its representation and decoder—and evaluated each possible pair of components.
135 We train the representations (Hi-C-LSTM, SCI, SNIPER) on all chromosomes and couple them
136 with selected decoders (LSTM, CNN, FC). Using the representations as input, we re-train these
137 decoders with a small subset of the chromosomes and test on the rest. (see Methods for more
138 details). We compute the average R-squared value for creating the Hi-C contact matrix using each
139 combination of selected representations and decoders

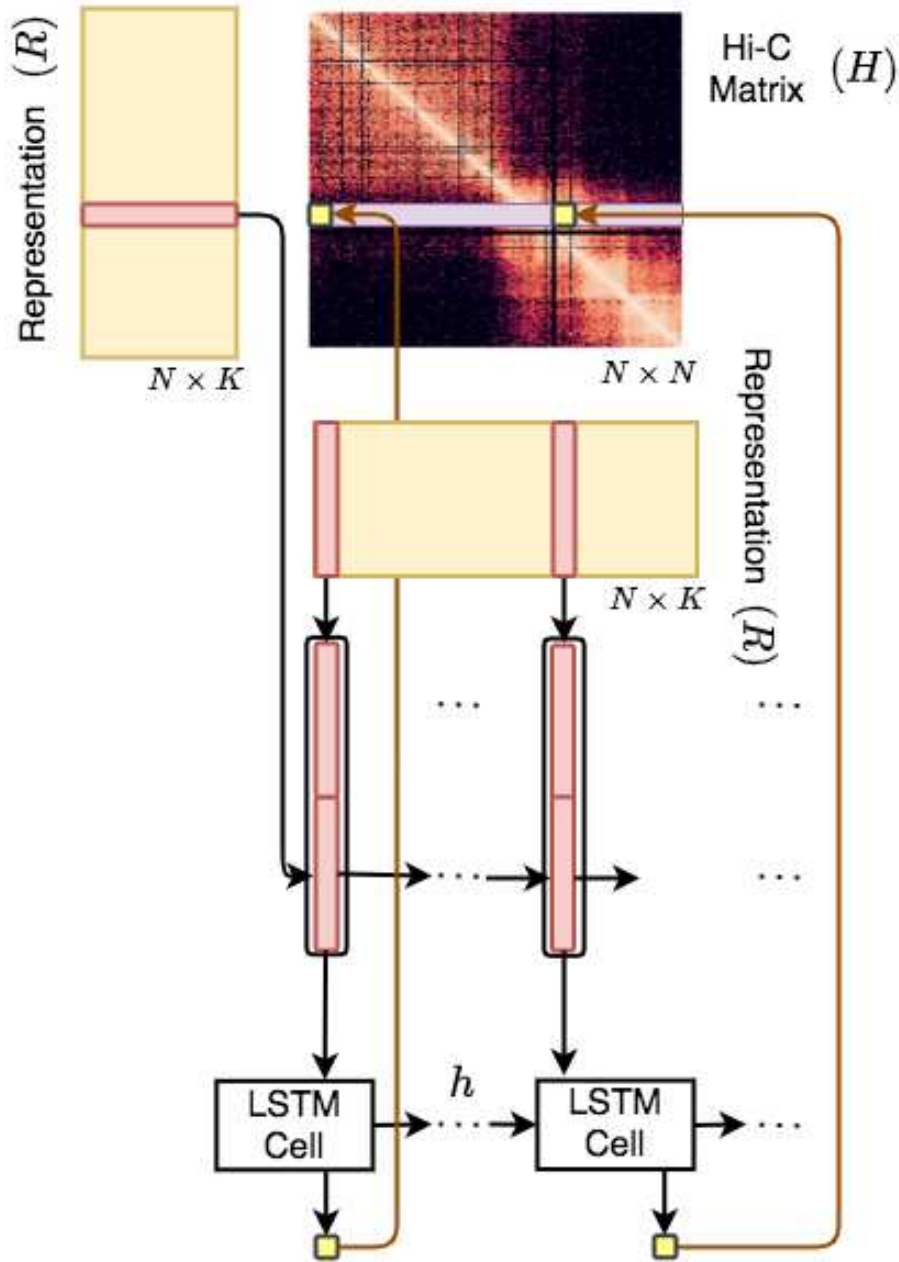


Figure 2: Overview of the Hi-C-LSTM model. A trained Hi-C-LSTM model consists of a K -length representation R_i for each genomic position i and LSTM connection weights (Methods). To predict the contact vector of a position i with all other positions, the LSTM iterates across the positions $j \in \{1 \dots N\}$. For each (i, j) pair, the LSTM takes as input the concatenated representation vector (R_i, R_j) and outputs the predicted Hi-C contact probability $H_{i,j}$. The LSTM hidden state h is carried over from (i, j) to $(i, j + 1)$. This process is repeated for all N rows of the contact map by reinitializing the LSTM states. The LSTM and the representation matrix are jointly trained to minimize the reconstruction error.

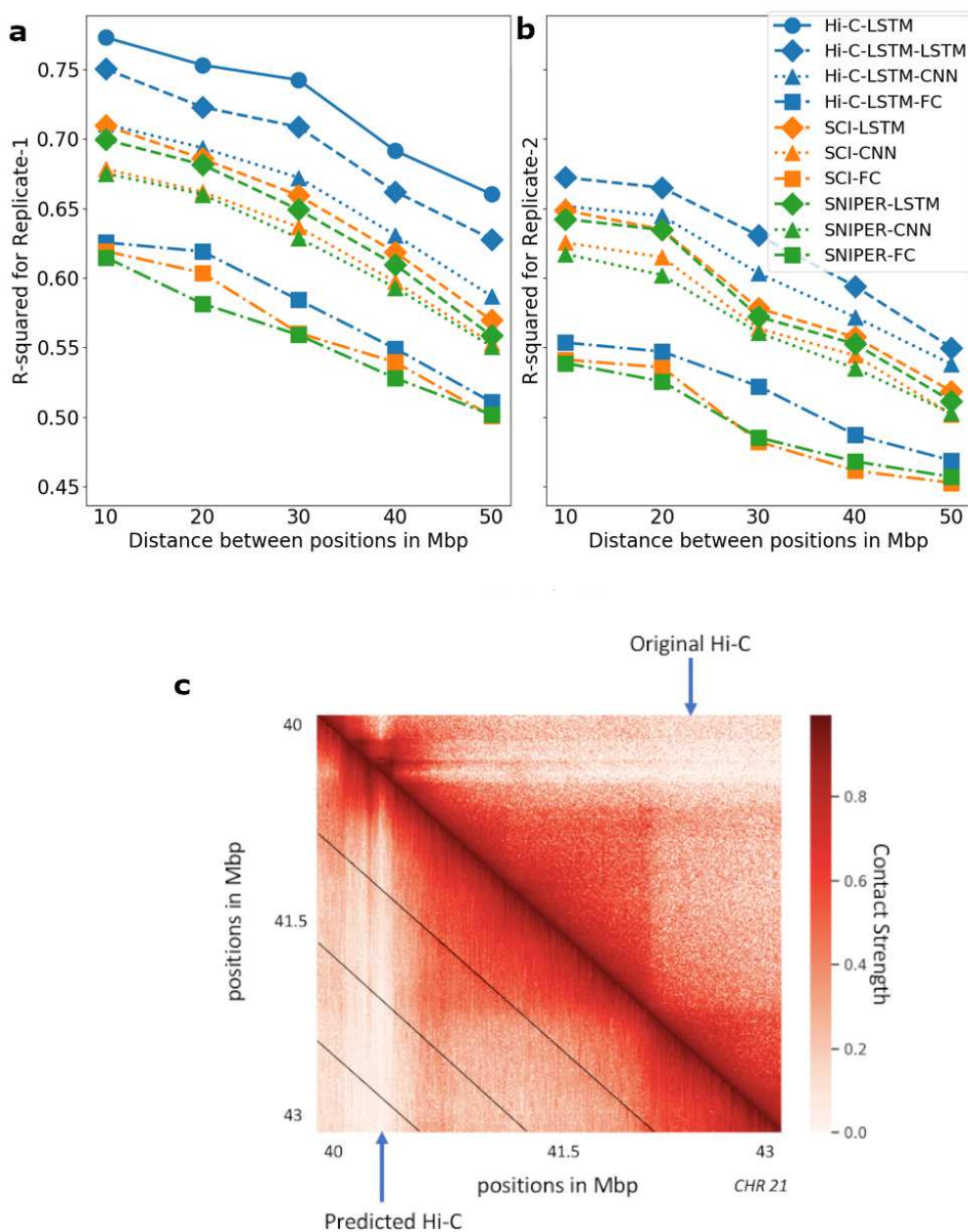


Figure 3: Accuracy with which representations reproduce the original Hi-C matrix. **a,b**) The Hi-C R-squared computed using the combinations of representations from different methods and selected decoders for replicate 1 and 2. The horizontal axis represents the distance between positions in Mbp. The vertical axis shows the average R-squared for the predicted Hi-C data. The R-squared was computed on a test set of chromosomes using selected decoders with the representations trained all chromosomes as input. The legend shows the different combinations of methods and decoders, read as $[representation]-[decoder]$. **c**) A selected portion of the original Hi-C map (upper-triangle) and the predicted Hi-C map (lower-triangle). The portion is selected from chromosome 21, between 40 Mbp to 43 Mbp. Diagonal black lines denote Hi-C-LSTM's frame boundaries (Methods).

140 We found that the choice of decoder has the largest influence on reconstruction performance.
141 Using a LSTM decoder performs best, even when using representations derived from SNIPER or
142 SCI (improvement of 0.14 and 0.12 R-squared points on average over fully-connected decoders
143 respectively, Fig. 3a). Furthermore, we found that Hi-C-LSTM’s representations are most infor-
144 mative, even when using decoder architectures derived from SNIPER or SCI (Fig. 3a).

145 Though the Hi-C-LSTM representations capture important information from a particular sam-
146 ple, we wanted to verify whether they capture real biological processes or irreproducible experimental
147 noise. To check the effectiveness of Hi-C-LSTM representations in creating the Hi-C contact map of
148 a biological replicate, we train the representations on one replicate (replicate 1), repeat the decoder
149 training process on replicate 2 (see Methods for more details), and compute the average R-squared
150 value for creating the Hi-C contact map of replicate 2 (Fig. 3b). The average R-squared reduces
151 slightly for inference of replicate 2 due to experimental variability; however, the performance trend
152 of the representation-decoder combinations is largely preserved (Fig. 3b). These results show that
153 Hi-C-LSTM’s improved performance is not merely driven by memorizing irreproducible noise.

154 **Hi-C-LSTM representations locate functional activity, genomic elements, and** 155 **regions that drive 3D conformation**

156 Considering that a good representation of Hi-C should contain information about the regulatory
157 state of genomic loci, we evaluated our model by checking whether these genomic phenomena
158 and regions are predictable from only the representation. Specifically, we test whether the position
159 specific representations learned via the Hi-C contact-generation process are useful for genomic tasks
160 that the model was not trained on, such as classifying genomic phenomena like gene expression
161 [32] and replication timing [33–36], locating nuclear elements like enhancers, transcription start
162 sites (TSSs) [37] and nuclear regions that are associated with 3D conformation like promoter-
163 enhancer interactions (PEIs) [38–40], frequently interacting regions (FIREs) [41,42], domains, loops
164 and subcompartments [20]. We used a boosted decision tree (XGBoost) model [43] to predict
165 binary genomic features from representations. (See Methods for more details regarding comparison
166 methods, baselines and classifier).

167 We find that the models built using the intra-chromosomal representations achieve higher predic-
168 tive accuracy overall relative to ones trained on inter-chromosomal representations when predicting
169 gene expression, enhancers and TSSs (Fig. 4a). This trend is likely due to the relatively close
170 range of the elements involved in prediction. In contrast, SNIPER is slightly better at predict-
171 ing replication timing when compared to the rest of the intra-chromosomal models except Hi-C
172 LSTM (SNIPER-INTER, Fig. 4a). While all methods achieve low absolute accuracy at predict-
173 ing promoter-enhancer interactions, Hi-C-LSTM performs best (0.5 mAP on average, 0.1 mAP
174 higher on average than SCI) (Fig. 4a, b). Both methods perform comparably in predicting the
175 other interacting genomic regions like FIREs, domains, loops, and subcompartments (Fig. 4a).
176 SNIPER-INTRA as well as SNIPER-INTER don’t perform as well as Hi-C-LSTM and SCI on
177 these tasks.

178 The only task on which other methods outperform Hi-C-LSTM is at predicting subcompart-
179 ments. Subcompartments were originally defined based on inter-chromosomal interactions, so rep-
180 resentations based on such interactions outperform those based on intra-chromosomal interactions
181 such as Hi-C-LSTM. Also subcompartment-ID (SBCID) (Methods) achieves perfect mAP by virtue
182 of its design (Fig. 4a). Among the rest of the methods, we find that methods which were designed
183 to predict subcompartments such as SCI and SNIPER-INTER, perform better than the others
184 (Fig. 4a). Hi-C-LSTM does perform marginally better than SNIPER-INTRA. Overall, although
185 Hi-C LSTM performs better than other models on most of the tasks, the performance of SCI and

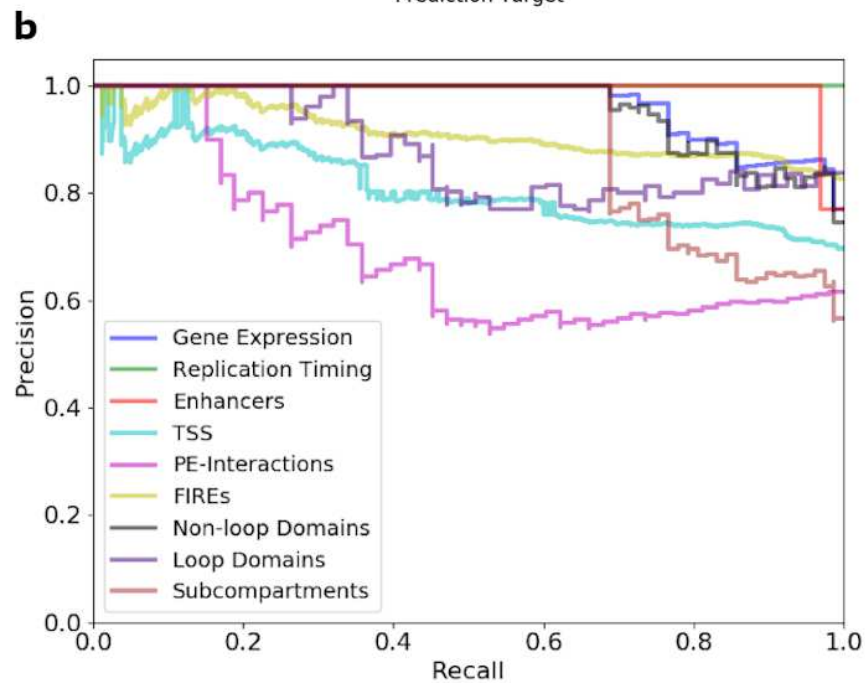
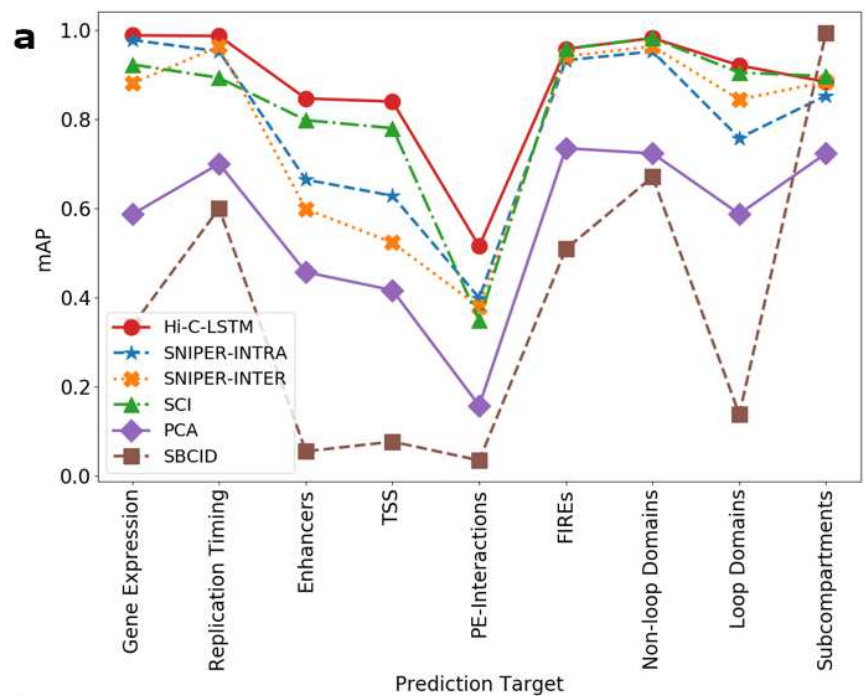


Figure 4: Important genomic phenomena and chromatin regions are classified using the Hi-C-LSTM representations as input. **a)** mAP for gene expression, replication timing, enhancers, transcription start sites (TSSs), promoter-enhancer interactions (PEIs), frequently interacting regions (FIREs), loop and non-loop domains, and subcompartments. The y-axis shows the mAP, the x-ticks refer to the prediction targets, and the legend shows the different methods compared with. **b)** The Precision-Recall curves of Hi-C-LSTM for the various prediction targets. The y-axis shows the Precision, the x-axis shows the Recall, and the legend shows the prediction targets.

186 SNIPER are comparable to Hi-C-LSTM and all three models perform significantly better than the
187 baselines on average (Fig. 4a).

188 **Feature attribution reveals association with genomic elements driving 3D con-** 189 **formation**

190 Given that our representations capture elements driving 3D conformation, we should be able to
191 identify those elements using our representations. To validate the ability of our representations to
192 locate genomic regions that drive chromatin conformation, we identified which genomic positions
193 have the largest impact on Hi-C contacts, using the technique of feature attribution. Feature
194 attribution is a technique that allows us to attribute the prediction of neural networks to their
195 input features. In this case, it identifies which genomic positions influence which Hi-C contacts.
196 We ran feature attribution analysis on the Hi-C-LSTM and aggregated the feature importance
197 scores across all the dimensions of the input representation to get a single score for each genomic
198 position (see Methods for more details). We expected to see higher feature attribution for the
199 elements, regions, and domains that are crucial for chromatin conformation.

200 We found that the CTCF and cohesin binding sites as given by ChIP-seq have a large influence
201 on contacts given their high feature importance score. The genome folds to form “loop domains”,
202 which are found to be a result of tethering between two loci bound by CTCF and cohesin subunits
203 RAD21 and SMC3 [18]. Among the many models of genome folding, cohesin ring-associated com-
204 plex that extrudes chromatin fibers and is delimited by CTCF is most promising. This extrusion
205 model explains why loops don’t overlap [17]. We found that CTCF sites show 10% higher mean im-
206 portance score than RAD21 and SMC3 sites and all three sites have a spread that is predominantly
207 positive (Fig. 5c). The high feature importance scores observed at CTCF and cohesin binding sites
208 validates the crucial role they play in loop formation [17, 18].

209 The importance of CTCF is further validated by the aggregated feature importance (Fig. 5c),
210 showing a markedly positive score near CTCF binding sites given by Segway [44], particularly
211 the strong ones (mean importance score of 0.45). Moreover, we see that the model places high
212 importance on regulatory elements, particularly enhancers (mean importance score of 0.4) (Fig.
213 5c). The active domain types have a higher mean score and a spread that largely occupies the
214 positive portion of the feature importance plot when compared to the inactive regions (Fig. 5c).
215 This suggests that active regions may play a dominant role in nuclear organization, where the
216 movement of repressed regions to the periphery is a side-effect.

217 Aggregated feature importance also demonstrates the largely positive feature attribution of
218 genomic regions that are an integral part of 3D conformation like FIREs, topologically associating
219 domain (TAD) boundaries with and without CTCF sites, loop and non-loop domains (Fig. 5c).
220 TAD boundaries enriched with CTCF show a 20% higher mean importance score compared to
221 TAD boundaries not associated with CTCF, pointing to the importance of CTCF sites at domain
222 boundaries in conformation (Fig. 5c). Moreover, loop domains show a 20% higher mean importance
223 score compared to non-loop domains, which is expected because of the increased contact strength
224 on average and the presence of CTCF sites (Fig. 5c).

225 The variation of the aggregated feature importance across interesting genomic regions helps us
226 distinguish boundaries of domains and genomic regulatory elements (Fig. 5a,b). We observe the
227 variation of the feature importance signal across TADs and a selected portion of chromosome 21
228 (28 Mbp to 29.2 Mbp) [45] to check if we can isolate the boundaries of domains, genes and other
229 regulatory elements. To deal with TADs of varying sizes, we partition the interior of all TADs
230 into 10 equi-spaced bins and average the feature importance signal within these bins. We plot this
231 signal along with the signal outside the TAD boundary 50Kbp upstream and downstream, averaged

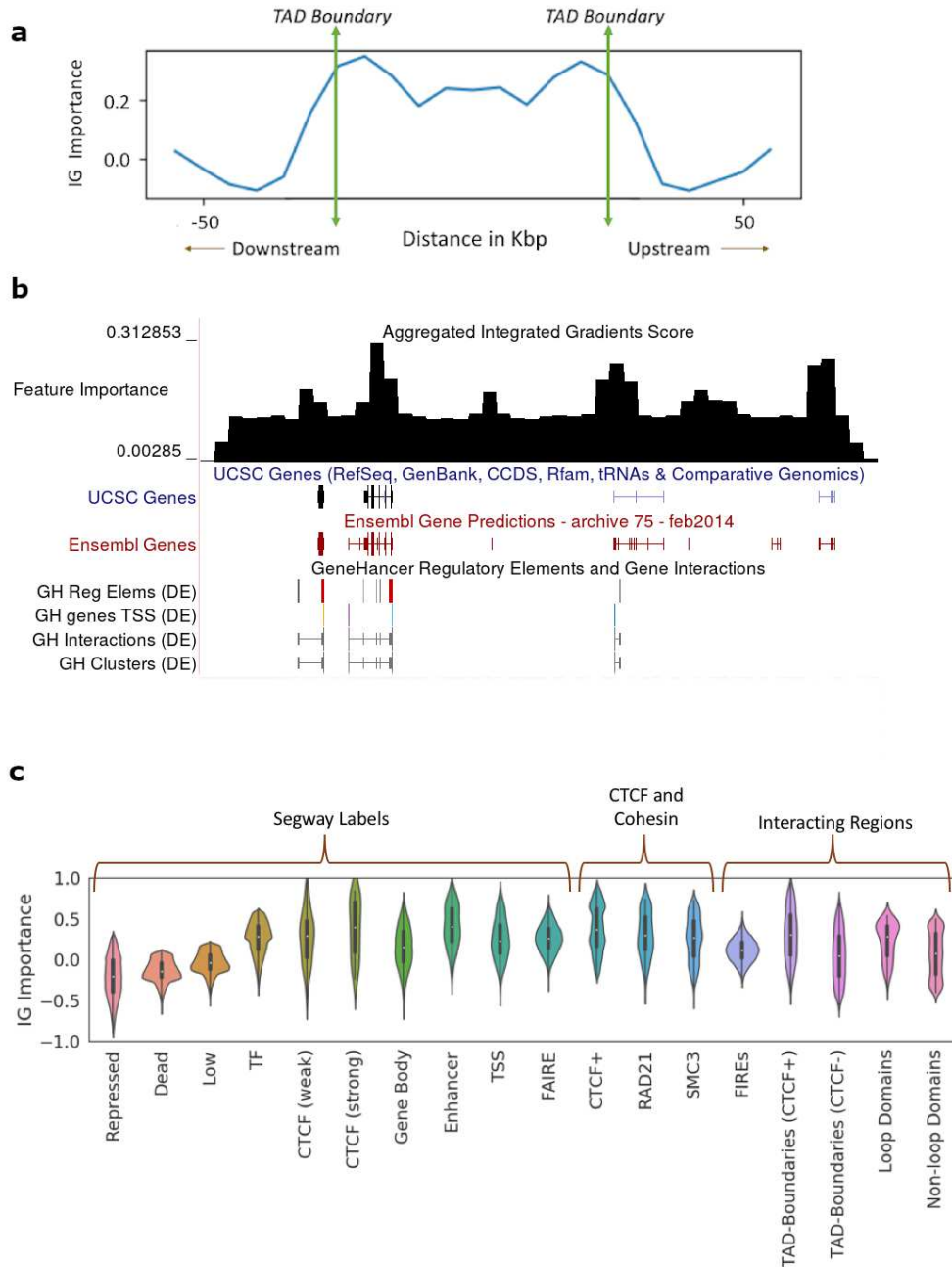


Figure 5: Hi-C-LSTM representations identify genomic elements involved in conformation through Integrated Gradients (IG) feature importance analysis. **a**) The IG feature importance averaged across different TADs of varying sizes. The vertical axis indicates the average IG importance at each position and the horizontal axis refers to relative distance between positions in Kbp, upstream/downstream of the TADs. **b**) The IG feature importance for a selected genomic locus (chr21 28-29.2Mbp) along with genes, regulatory elements and Hi-C. We see that the feature importance scores peak at known regulatory elements. **c**) Violin plots of aggregated feature attribution scores for selected elements. The x-axis shows the labels/elements and the y-axis displays the log plus z normalized feature importance scores from Integrated Gradients.

232 across all TADs (Fig. 5a). The feature importance has largely similar values in the interior of the
233 TAD, noticeably peaks at the TAD boundaries, and slopes downward in the immediate exterior
234 vicinity of the TAD (Fig. 5a). This trend validates the importance of TADs and TAD boundaries
235 in chromatin conformation, which we saw in (Fig. 5c). We also consider a candidate region in
236 chromosome 21 that is referred to in [45] to observe the variation of feature importance across
237 active genomic elements (Fig. 5b). For this selected region in chromosome 21, as we don't have to
238 deal with domains of varying sizes, we just average the feature importance signal within a specified
239 number of bins and plot this in the UCSC Genome Browser (<https://genome.ucsc.edu>) along with
240 genes and regulatory elements. The feature importance peaks around genes, regulatory elements
241 and domain boundaries (Fig. 5b), showing that they play a more important role in conformation
242 than other functional elements.

243 **Hi-C-LSTM enables in-silico knockout experiments**

244 As Hi-C-LSTM models the dependence of sequence on 3D conformation, it enables us to perform
245 in-silico deletion, insertion and reversal of certain genomic loci and observe changes in the result-
246 ing Hi-C contact map. In-silico knockout experiments have gained prominence lately, mainly in
247 intercepting signal flows in signaling pathways [46] and drug discovery [47–49]. A Hi-C in-silico
248 manipulation tool is of great value it enables researchers to identify the importance and influence
249 of any genomic locus of interest to 3D chromatin conformation.

250 Hi-C-LSTM enables a researcher to perform two types of experiments. First, one can simulate
251 the knockout of a locus by deleting a portion of the representation or replacing it with a null
252 representation. As a null, we use the average local features within 0.2 Mbp. Second, one can
253 simulate the replacement or translocation of an element by replacing or removing the corresponding
254 representation (see Methods).

255 Previous work showed that inserting even a single base pair near the loop anchors can make
256 many loops and domains vanish, altering chromatin conformation at the megabase scale [17]. Given
257 the crucial role played by CTCF and cohesin subunits in conformation at loop anchors (See Classi-
258 fication, Attribution), we hypothesized that knocking out CTCF and cohesin subunit binding sites
259 will change the Hi-C contact map noticeably. The average difference in predicted contact strength
260 between no knockout and knockout at the site under consideration as a function of genomic dis-
261 tance is observed (Fig. 6c). After CTCF and cohesin knockouts, the average contact strength
262 reduces by >15% when compared to the no knockout case (Fig. 6c). CTCF knockout is seen to
263 affect insulation at about 100 Kbp and reflect possible loss of loops at 200 Kbp (Fig. 6c). The
264 knockout of cohesin subunits SMC3 and RAD21 binding sites is observed to be independent of
265 CTCF knockout with 5% higher average inferred strength over distance, hinting at their relative
266 importance (Fig. 6c).

267 The CTCF sites at loop anchors occur mainly in a convergent orientation, with the forward
268 and reverse motifs together, suggesting that this formation maybe required for loop formation
269 [20, 50–55]. To check how important the orientation of CTCF motifs is to conformation, we con-
270 ducted CTCF orientation replacement experiments at loop boundaries. The average difference in
271 predicted contact strength between no replacement and replacement at the site under consideration
272 as a function of genomic distance is observed (Fig. 6c). The replacement of convergent with the
273 divergent orientation around loops is seen to behave similar to the case of CTCF knockout thereby
274 validating observations made in [56] (Fig. 6c). On the other hand, replacement of divergent with
275 the convergent orientation is seen to preserve loops at 200 Kbp and behave similar to the control,
276 although with reduced inferred contact strength (5% on average) (Fig. 6c).

277 The difference in inferred Hi-C between the CTCF (Fig. 6a) and cohesin (Fig. 6b) knockout

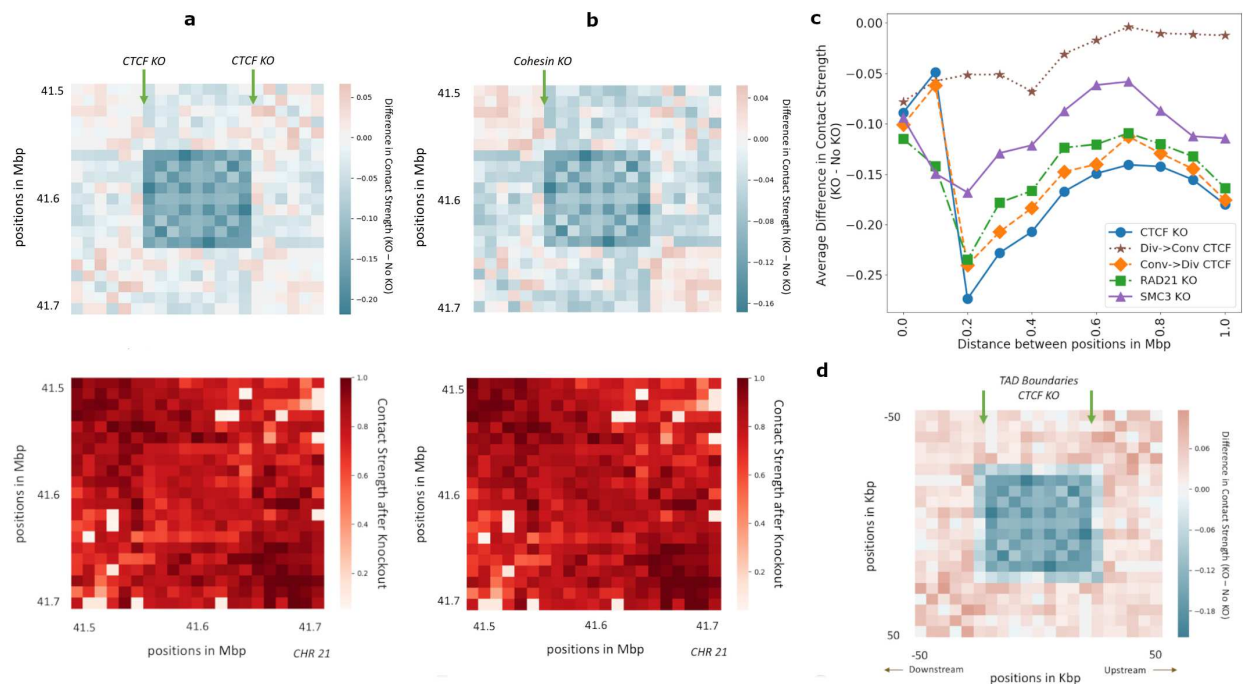


Figure 6: In-silico deletion and orientation replacement of CTCF and cohesin subunits is performed and changes in the resulting Hi-C contact matrix is observed. **a)** The difference in predicted Hi-C contact strength between CTCF knockout and no knockout, performed on chromosome 21, between 41.5 Mbp to 41.7Mbp. The bottom figure shows just the predicted Hi-C after knockout. **b)** Same as A, but for Cohesin knockout. **c)** Average difference in contact strength of the inferred Hi-C matrix between knockout and no knockout (y-axis) for varying distance between positions in Mbp (x-axis). The knockout experiments include CTCF and cohesin knockout and convergent/divergent CTCF replacements (legend). **d)** The genome-wide average difference in predicted Hi-C contact strength between CTCF knockout at TAD boundaries and no knockout.

278 and the no knockout for a selected portion of chromosome 21 (41.5 Mbp to 41.7Mbp), shows the
279 importance of CTCF and cohesin sites in conformation. The CTCF knockout at both the edges
280 of the loop results in decrease in contact strength (0.18 lower on average) within the loop (Fig.
281 6a). Cohesin knockout at the start of the loop also results in decrease in contact strength within
282 the loop (0.12 lower on average), but not as strongly as the CTCF knockout (Fig. 6b). Around
283 the loop, CTCF and cohesin knockout results in patches of decreased (0.05 lower on average) as
284 well as increased contacts (0.05 higher on average) (Fig. 6a, b). The predicted Hi-C after CTCF
285 and cohesin knockout (Fig. 6a, b:Bottom) validates the fading of loops. The average difference
286 in inferred Hi-C between the CTCF knockout at TAD boundaries and the no knockout (Fig.
287 6d) shows similar trends, with decreased contacts (0.2 lower on average) within the TAD and
288 increased contacts (0.08 higher on average) outside the TAD. The symmetry of the Hi-C matrix is
289 largely preserved after the knockouts, validating the capability of Hi-C-LSTM to perform knockout
290 experiments.

291 **Hi-C-LSTM accurately predicts effects of a 2.1 Mbp duplication at the SOX9** 292 **locus**

293 To further validate Hi-C-LSTM as a tool for in-silico genome alterations, we simulated a structural
294 variant at the SOX9 locus that was previously assayed by Melo et al. [57]. This variant was
295 observed in an individual with Cook’s syndrome and comprises the tandem duplication of a 2.1 Mbp
296 region on chromosome 17 that includes regulatory elements of SOX9 (chr17:67,958,880–70,085,143;
297 GRCh37/hg19, Fig. 7a). To simulate a Hi-C experiment on a genome with this variant, we made
298 a new Hi-C-LSTM representation matrix that includes a tandem copy of the representation at the
299 locus in question and passed this representation matrix through the original Hi-C-LSTM decoder
300 to produce a simulated Hi-C matrix on a post-duplication genome (Fig. 7b). Because Hi-C reads
301 cannot be disambiguated between the two duplicated loci, we simulated mapping reads to the
302 original hg19 reference by summing reads originating from the two copies (see Methods). We
303 evaluated Hi-C-LSTM’s predictions according to the agreement between this predicted matrix and
304 a Hi-C experiment performed by Melo et al. [57] (Fig. 7c).

305 We found that Hi-C-LSTM accurately predicted the effect of the duplication. The domains that
306 existed pre-duplication (D_1 , D_2 , D_3 , Fig. 7a) are correctly captured post-duplication. In addition,
307 a new chromatin domain (D_{New}) that was introduced by the duplication is correctly predicted
308 by Hi-C-LSTM (Fig. 7b). To quantitatively evaluate our predictions, we compared them to a
309 baseline that predicts the original pre-duplication Hi-C for the interactions between the upstream,
310 downstream and duplicated regions, and the genomic average for the interactions of the duplicated
311 region with itself (see Methods). We found that Hi-C-LSTM’s predictions significantly outperform
312 this baseline overall (Fig. 7d). Note the baseline is a slightly better predictor of contacts between
313 the upstream and downstream regions.

314 Hi-C-LSTM’s predictions have the advantage that they describe contacts on the true post-
315 duplication genome, in contrast to the reference genome used to map reads (Fig. 7c). Hi-C-LSTM’s
316 contacts recapitulate the post-duplication topological domain structure hypothesized by Melo et
317 al. These duplication experiments further validate the ability of Hi-C-LSTM to perform in-silico
318 mutagenesis.

319 **Discussion**

320 In this work we have proposed a deep LSTM model that uses intra-chromosomal contacts to form
321 position-specific representations of chromatin conformation. These representations are able to

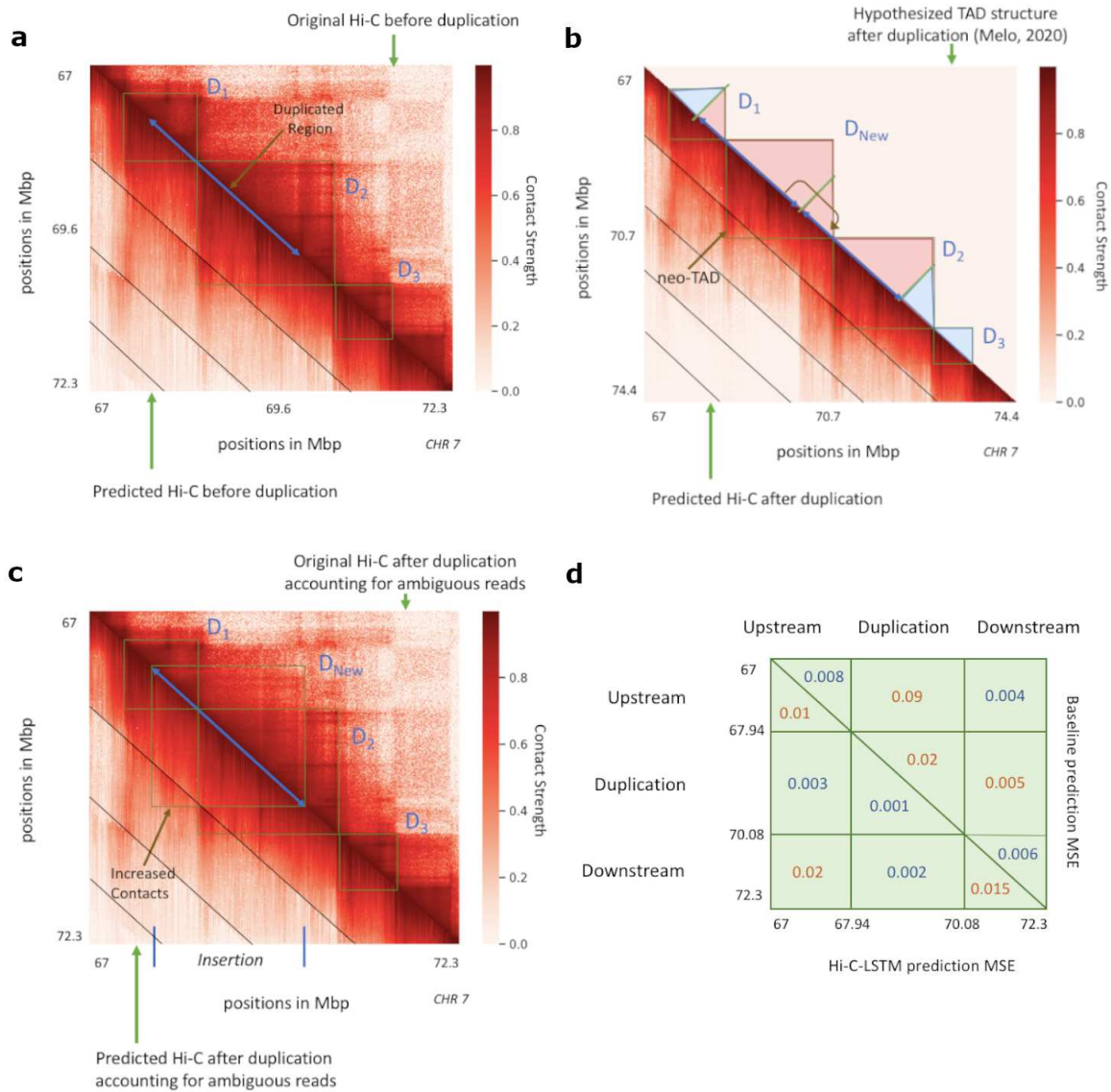


Figure 7: In-silico duplication of a 2.1 Mbp region on Chromosome 7 [57]. In all subplots, upper and lower triangles denote observed and predicted Hi-C contact probabilities respectively, and diagonal black lines denote Hi-C-LSTM frame boundaries. **a)** Original and predicted Hi-C before duplication. D_1 , D_2 and D_3 indicate the three pre-duplication topological domains. **b)** Predicted Hi-C after duplication on a simulated reference genome that includes both copies. Lower triangle indicates Hi-C-LSTM predicted contacts. The true Hi-C contact matrix on this reference genome is not observable because the read mapper cannot disambiguate between the two copies. The upper triangle depicts the post-duplication topological domain structure hypothesized by Melo et al, which includes a novel topological domain D_{New} . **c)** Original and predicted Hi-C on the original pre-duplication reference genome. Upper triangle shows observed post-duplication Hi-C data assayed by Melo et al. Lower triangle shows Hi-C-LSTM predictions, mapped to the pre-duplication reference by summing the contacts for the two copies (Results). **d)** Average mean-squared error in predicting the observed data by (lower triangle) Hi-C-LSTM, and (upper triangle) a simple baseline (Results) at the upstream, duplicated, and downstream regions.

322 capture a variety of genomic phenomena and elements and at the same time distinguish genomic
323 regions, transcription factors and domains that are known to play an important role in chromatin
324 conformation. They also elucidate the interplay between genome structure and function. The
325 classification and feature attribution results validate the ability of the representations to locate
326 vital regions such as CTCF and cohesin binding sites.

327 The primary contribution of this work is the application of a deep LSTM to the problem of
328 forming representations for intra-chromosomal interactions. The Hi-C-LSTM not only outperforms
329 the existing models like SCI and SNIPER that form representations in predicting genomic phenom-
330 ena but also locates elements driving 3D conformation as revealed by feature importance analysis.
331 In addition to these, the Hi-C-LSTM has few distinct advantages over its counterparts. One, it
332 can be used as a contact generation model. It’s observed that the Hi-C-LSTM representations are
333 more informative in this regard and that sequential models like the LSTM perform much better
334 at contact generation. Two, a low-dimensional Hi-C-LSTM representation is powerful enough to
335 reasonably recreate the Hi-C matrix (see Hyperparameters). Three, the Hi-C-LSTM framework
336 allows us to conduct in-silico experiments like insertion, deletion and reversal of elements driving
337 3D conformation and observe changes in contact generation. This would be extremely useful in
338 fully understanding the role of CTCF and cohesin binding sites and other transcription factors in
339 chromatin conformation.

340 An important limitation of Hi-C-LSTM’s *in silico* experiment is that they can simulate only *cis*
341 effects. Variation in chromatin structure can be caused either by *cis* or *trans* effects. *Cis* effects are
342 caused by genetic variants on the same DNA molecule, whereas *trans* effects arise from diffusible
343 elements like transcription factors. Hi-C-LSTM can model only *cis* effects because *trans*-acting
344 cellular machinery is captured within the Hi-C-LSTM decoder, which cannot be easily modified.
345 An example of a *cis*-effect is the duplication at the SOX9 locus, in which case we showed Hi-C-
346 LSTM correctly models the resulting neo-TAD (see Duplication) [57]. Hi-C-LSTM cannot model
347 *trans* effects such as recent investigation of the removal of RAD21 [18] and CTCF [58, 59].

348 The good performance of Hi-C-LSTM suggests several avenues for future work. First, extending
349 the mode to incorporate data from multiple cell types and the resulting representations may yield
350 insights into differences in chromatin organization across development. Second, the success of a
351 LSTM model suggests trying other recurrent neural network models such as Transformers [60].
352 Third, a modified version of Hi-C-LSTM may be able to infer a 3D structure of chromatin. The Hi-
353 C representations that we form currently are embedded on a lower-dimensional manifold that does
354 not have any direct physical significance. However, a Hi-C-LSTM-like model trained to produce
355 three-dimensional representations may be able to reproduce the true nuclear positions of chromatin.

356 **Methods**

357 The code and data repository for this project, including training, evaluation, data handling, and
358 generated data can be found in our GitHub repository (<https://github.com/smaslova/HiCLSTM>).

359 **Data sets**

360 The Hi-C data for the GM12878 B-Lymphocyte cell line was acquired using the GEO accession
361 number GSE63525 (<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE63525>) [20].

362 We generated a intrachromosomal Hi-C data set on the hg19 human reference genome assembly
363 [61] at 10Kb resolution with KR (Balanced) normalization [62] using juicer tools (<https://github.com/aidenlab/juicer/wiki/Juicer-Tools-Quick-Start>) with the command
364

365 `java -jar juicer_tools.jar dump observed KR data/chr.hic chr chr BP 10000 chr.txt,`
366 where `chr` refers to the chromosome being extracted.

367 Following SCI [12], to mitigate the extreme range of magnitudes present in Hi-C read counts,
368 we transformed Hi-C values into contact probabilities between 0 and 1. We calculated contact
369 probabilities according to the exponential transformation (Eq. 1)

$$cf = \frac{1}{v + \delta} \tag{1}$$
$$CP = \exp(-a * cf),$$

370 where v is the raw input contact strength, δ is a very small positive real number (we set δ to
371 be 10^{-10}), cf is the coefficient obtained, a is the coefficient multiplier, and CP is the resulting
372 contact probability. We chose $a = 8$ because it appeared to provide a good separation of low and
373 high contact values.

374 RNA-seq data for 57 cell types was obtained from the Roadmap Consortium (https://egg2.wustl.edu/roadmap/web_portal/processed_data.html).

376 For the classification task, each gene was considered to be active if its log mean expression value
377 across the gene was greater than 0.5 [63, 64].

378 We defined promoter-enhancer interactions as the ones that were used to train TargetFinder
379 (<https://github.com/shwhalen/targetfinder>) [65].

380 Frequently interacting region (FIRE) scores at 40Kbp resolution were downloaded from the
381 additional material of [41] and were converted to binary indicators using 0.5 as a threshold following
382 [66].

383 The replication timing data given by Repli-Seq [67] was downloaded from Replication Domain
384 (<https://www2.replicationdomain.com>) at 40Kbp resolution.

385 Locations of known enhancers and transcription start sites (TSSs) were obtained from FANTOM
386 (<https://fantom.gsc.riken.jp/5>) and ENCODE (<https://www.encodeproject.org/files/ENCFF140P>
387 CA) respectively.

388 Domain, loop and subcompartment annotations were obtained from the results of Rao et al. [20]
389 using the GEO accession number GSE63525 (<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE63525>).

391 Segway and Segway-GBR labels were obtained from Hoffmanlab (<https://segway.hoffmanlab.org>)
392 and Noblelab (<https://noble.gs.washington.edu/proj/gbr>) respectively.

393 CTCF, RAD21 and SMC3 peak calls were downloaded from ENCODE (<https://www.encodeproject.org>). The CTCF orientations were obtained by using the CTCF motif from the MEME suite
394 (<https://meme-suite.org/meme/doc/fimo.html>) (version 5.3.3) and running FIMO [68] to get the
395 motif instances using the command

397 `fimo -oc output_directory motif_file.meme sequence_file.fna.`

398 We use all default options while running `fimo` including the p-value threshold (`--thresh`) of 10^{-4} .
399 We ran FIMO after obtaining the human genome sequence file under mammals and the hg19
400 genome assembly.

401 Topologically-associating domain (TAD) annotations were downloaded from TADKB [69].

402 LSTM

403 Long short-term memory (LSTM) networks were proposed as a solution to the vanishing gradient
404 problem [70] in recurrent neural networks (RNNs) [71]. They are known to be a good candidate
405 for modelling sequential data and have been widely used for sequential tasks [72–74]. An LSTM
406 is made up of a memory state (h_t), a cell state (c_t), and three gates that control the flow of data:

407 input (i_t), forget (f_t) and output (o_t) gates. The input and the forget gates together regulate the
 408 effect of a new input on the cell state. The output gate determines the contribution of the cell state
 409 on the output of the LSTM.

410 Let matrices W and U be the weights of the input and recurrent connections, and b refer to the
 411 biases. There are four sets of weight matrices and biases in the LSTM. These include one for each
 412 of the three gates—forget gate (W_f, U_f, b_f), input gate (W_i, U_i, b_i) and output gate (W_o, U_o, b_o)—
 413 and one to form the cell state (W_c, U_c, b_c). The current cell state (c_t) is formed by the modulation
 414 of the previous cell state (c_{t-1}) by the forget gate (f_t) and combining it with the modulation of
 415 the current input (x_t) and the previous memory state (h_{t-1}) by the input gate (i_t). Finally, the
 416 current memory state (h_t) is formed by the modulation of the current cell state (c_t) by the output
 417 gate (o_t).

An LSTM’s output is determined by the following series of operations [31].

$$\begin{aligned}
 f_t &= \sigma(W_f x_t + U_f h_{t-1} + b_f) \\
 i_t &= \sigma(W_i x_t + U_i h_{t-1} + b_i) \\
 o_t &= \sigma(W_o x_t + U_o h_{t-1} + b_o) \\
 c_t &= f_t \circ c_{t-1} + i_t \circ \sigma(W_c x_t + U_c h_{t-1} + b_c) \\
 h_t &= o_t \circ \sigma(c_t)
 \end{aligned} \tag{2}$$

418 where \circ is the Hadamard product and σ refers to the sigmoid activation function.

419 Hi-C-LSTM

420 Hi-C-LSTM creates a representation given a pair of genomic positions in the Hi-C contact matrix
 421 using an embedding neural network layer and predicts the contact strength at that particular pair
 422 via a deep LSTM [31] that takes these representations as input (Fig. 2). Hi-C-LSTM takes as input
 423 a $N \times N$ intra-chromosomal Hi-C contact matrix ($\mathbb{R}^{N \times N}$), for each chromosome, where N is the
 424 chromosome length.

425 A trained Hi-C-LSTM model consists of LSTM parameters (section LSTM) and a representation
 426 matrix $R \in \mathbb{R}^{N \times M}$, where M is the representation size. At each genomic position, (i, j) pair is given
 427 as input to an embedding layer, which indexes the row and column representations $R_i, R_j \in \mathbb{R}^M$
 428 and feeds these two vectors as input to the LSTM. The output of the LSTM is the predicted Hi-C
 429 contact probability $\hat{H}_{i,j}$ for the given (i, j) pair.

430 The hidden states of the LSTM are carried over from preceding columns thereby maintaining a
 431 memory for the row. For the sake of memory usage, the hidden states are reinitialized after every
 432 each frame of 1.5 Mbp or 150 resolution bins (see LSTM). This process is repeated for each row of
 433 the Hi-C matrix (Eq. 3).

$$\begin{aligned}
 \hat{H}_{i,j} &= LSTM(R_i, R_j, h_j) \quad \text{for } j = 1, 2, \dots, N \\
 &\quad \text{for } i = 1, 2, \dots, N
 \end{aligned} \tag{3}$$

434 where h_j is the same as h_{j-1} within the frame and is reinitialized at the beginning of each new
 435 frame.

436 The LSTM and the embedding neural network layer are jointly trained using the mean squared
 437 error (MSE) loss function which facilitates the faithful construction of the Hi-C intra-chromosomal
 438 matrix (Eq. 4).

$$MSE = \frac{1}{N} \left[\sum_{i=j}^N (H_{i,j} - \hat{H}_{i,j})^2 \right] \quad \text{for } i = 1, 2, \dots, N \tag{4}$$

439 At the end of all the training iterations, the output of the embedding neural network layer at
440 each row i (R_i) is treated as the representation for that row. The Hi-C-LSTM framework infers
441 the Hi-C contact matrix from pairs of position IDs and therefore is a transformation from linear
442 sequential space to the Hi-C space. The linear position IDs are a convenient and useful modeling
443 assumption which builds a framework that doesn't make any other transfer function assumptions.

444 Modeling choices and training

445 The LSTM model required us to make a few design choices. As layer normalization can significantly
446 reduce the training time and is effective at stabilizing the hidden state dynamics in LSTMs, we
447 used a unidirectional layer norm LSTM [75] with one hidden layer. We found that variants such
448 as the bidirectional LSTM [76] and LSTM with multiple layers provided a marginal increase in
449 test performance (Supplementary: Fig. 1). The variants were also prone to overfitting. Therefore,
450 we chose the single-layer unidirectional model over these variants accounting for computational
451 efficiency and good generalization. Gradient clipping [70] and the *softsign* activation [77] were
452 used at all nodes owing to their mitigating effect on hidden state saturation. The design choices
453 were made after conducting ablation experiments which are elaborated in the following section
454 (Hyperparameters). We used a batch size of 300 and a sequence length 150 bins, both of which
455 were observed to be data dependent and the best fit for our data. We used a learning rate of 0.01
456 for 5 epochs and 0.001 for 5 more epochs. We reinitialized the hidden states of the LSTM after
457 every frame of length 150 and predicted each diagonal block of length 150 with fresh hidden states
458 (Figure 3c). The prediction error improved towards the end of the frame and increased at the
459 start of the next frame (Supplementary: Fig. 1). We tried passing the hidden states across frames
460 and saw that the convergence time significantly increased as the training graph had to be retained
461 across iterations. So we chose to reinitialize the hidden states in each window instead.

462 We employed PyTorch (<https://pytorch.org>), a Python-based deep learning framework and
463 trained Hi-C-LSTM on GeForce GTX 1080 Ti GPUs with ADAM as the optimizer [78]. All
464 parameters in PyTorch were set to their default values while training. As our primary goal was
465 not to infer values for unseen positions but to form reliable representations for every chromosome,
466 we trained our model on all the chromosomes. For our Hi-C reproduction evaluation, we trained
467 the representations on all chromosomes but the decoders only on a random subset. We chose to
468 train the decoders on a random subset of chromosomes to prevent the decoder from overpowering
469 the representations.

470 Hyperparameter selection

471 To choose the representation size of our model, we performed an ablation analysis. We computed
472 the average mAP across all downstream tasks with the Hi-C-LSTM model which consists of a single
473 layer, unidirectional LSTM with layer norm in the absence of dropout [79] for odd chromosomes
474 and used the even chromosomes to validate whether the choice of hyperparameters remained the
475 same irrespective of chromosome set. We observed the mAP (section Classification Method) of the
476 Hi-C-LSTM vs. increasing representation size along with Hi-C-LSTM that is bidirectional, in the
477 presence of dropout, without layer norm and 2 layers (Supplementary: Fig. 2). While both the
478 presence of dropout and the absence of layer norm adversely affected mAP, the addition of a layer
479 and a complimentary direction did not yield significant improvements in downstream performance.
480 We conducted a similar ablation experiment and computed the average Hi-C R-squared for the
481 predictions with increasing representation size (Supplementary: Fig. 2) and observed that the
482 performance trend is preserved, which was indicative of the fact that recreating the Hi-C matrix

483 faithfully aids in doing well across downstream tasks. These results were verified to be true for even
484 chromosomes as well (Supplementary: Fig. 2). For both odd and even chromosomes, even though
485 the Hi-C prediction accuracy increased substantially with hidden size, we noticed the elbow at a
486 representation size of 16 for average mAP and therefore set our representation size to that value as
487 a trade-off.

488 **Hi-C reproduction evaluation**

489 We investigated three hypotheses with following analysis. First, we asked whether the Hi-C-LSTM
490 representations faithfully construct the Hi-C matrix. Second, whether the Hi-C-LSTM representa-
491 tion and the decoder are both powerful in generating the Hi-C map. Third, we evaluated the utility
492 of the representations to infer a replicate map. In all cases, we computed the average prediction ac-
493 curacy in reconstructing the Hi-C contact matrix, measured using R-squared, which represents the
494 proportion of the variance of the original Hi-C value that's explained by the Hi-C value predicted
495 by the Hi-C-LSTM.

496 In our first experiment, we trained both the representations and decoders on replicate 1 (Figure
497 3a). We took representations trained using all chromosomes from Hi-C-LSTM, SCI and SNIPER
498 and coupled these with some selected decoders, namely, a LSTM, a convolutional neural network
499 (CNN) and a fully connected (FC) feed-forward neural network (used by SNIPER). We compared
500 LSTM with CNN and FC decoders mainly because CNNs provided us with an alternative way of
501 incorporating structure (using moving filters) and FC networks did not include any information
502 about underlying structure. We re-trained these decoders using either of the representations as
503 input, with a small subset of the chromosomes and tested on the rest. All the decoders were
504 configured to have the same number of layers and hidden size per layer. As the decoders were
505 separately trained, this process allowed us to check the power of the representations alone, moreover,
506 as a small subset of chromosomes were used to train the decoder, we reduced the possibility of the
507 decoders overfitting.

508 In our second experiment (Figure 3b), we trained the representations on replicate 1 using all
509 chromosomes, and repeated the aforementioned decoder training process on replicate 2.

510 **Comparison methods**

511 We compared our downstream classification results with five alternatives: two variations of SNIPER
512 (<https://github.com/ma-compbio/SNIPER>), one with inter-chromosomal (SNIPER-INTER) and
513 the other with intra-chromosomal contacts (SNIPER-INTRA), SCI (<https://github.com/TheJackso>
514 [nLaboratory/sci](https://github.com/TheJacksonLaboratory/sci)) and two baselines, namely, the subcompartment-ID (SBCID) and principal com-
515 ponent analysis (PCA). SNIPER-INTRA was the same as the original SNIPER-INTER, modified
516 to take the intra-chromosomal row as input instead of the inter-chromosomal row. All the param-
517 eters for the two SNIPER versions and SCI were set as given in their respective papers [11], [12].
518 The SBCID baseline used the one-hot-encoded vector of the subcompartment as the representation
519 at the position under contention. The PCA baseline assigned the principal components from the
520 PCA of the Hi-C matrix as the representations.

521 **Element identification evaluation**

522 We used the following analysis to evaluate the ability of a representation to identify genomic
523 phenomena and chromatin regions.

524 For each type of element, a boosted decision tree classifier called XGBoost [43] was trained on
525 the representations. We employed tree boosting as it is shown to outperform other classification

526 models with respect to accuracy when ample data is available. Following Avocado [66], we used
527 XGBoost with a maximum depth of 6 and a maximum of 5000 estimators and these parameters
528 were chosen following ablation experiments with odd chromosomes as the training set and even
529 chromosomes as the test set (Supplementary: Fig. 3). N-fold cross-validation, with $n = 20$, was
530 used to validate our training with and an early stopping criterion of 20 epochs. The rest of the
531 XGBoost parameters were set to their default values.

532 For each task, the genomic loci under contention were assigned labels. All tasks were treated
533 as binary classification tasks, except the subcompartments task, which was treated as a multi-class
534 classification task. For tasks without preassigned negative labels, negative labels were created by
535 randomly sampling genome-wide, excluding the regions with positive labels.

536 The XGBoost classifier was given the representations at these genomic loci as input and the
537 assigned labels as targets. The classifier was evaluated using the metric of mean average precision
538 (mAP), which is a standard metric for classification tasks and is defined as the average of the
539 maximum precision scores achieved at varying recall levels.

540 Sequence attribution

541 We validated the utility of the Hi-C-LSTM representations in locating genomic regions important for
542 conformation using feature attribution analysis. Feature attribution was carried out on the intra-
543 chromosomal representations using Integrated Gradients [80]. Integrated Gradients is a feature
544 attribution technique that follows an axiomatic approach to attribution, adhering to the axioms of
545 sensitivity and implementation invariance. Sensitivity implies that if the input and baseline differs
546 in one feature and have different predictions, then the differing feature should be assigned a non-
547 zero attribution. Implementation invariance requires that two networks, whose output is equal for
548 every input despite having different implementations, should have the same attributions. We used
549 Captum (<https://captum.ai>), a Integrated Gradients feature attribution framework in PyTorch
550 that is generic and works with sequential models. The resulting feature attributions were summed
551 across all features, giving us one importance score for every position in the genome. The feature
552 importance scores were then subjected to log normalization followed by min-max normalization
553 (Eq. 5). Specifically, let IG be to the integrated gradients (IG) score, IG_{min} and IG_{max} be the
554 minimum and maximum IG scores. The normalized IG score IG_{norm} is defined as

$$IG_{norm} = \frac{\log IG - \log IG_{min}}{\log IG_{max} - \log IG_{min}}. \quad (5)$$

555 In-silico perturbation

556 The Hi-C-LSTM enables us to perform in-silico deletion, orientation replacement and reversal of
557 genomic loci and predict changes in the resulting Hi-C contact map. We performed three types
558 of experiments:: knockout, CTCF orientation replacement, and duplication. In a knockout experi-
559 ment, we chose certain genomic sites (such as CTCF and cohesin binding sites) and replaced their
560 representations with a null representation. As a null representation, we used the average repre-
561 sentations in a window of 0.2 Mbp around the site in question, because this captures the genomic
562 neighborhood while removing the features specific to site. The knockout of the representation
563 at a particular row affects not just the Hi-C inference at columns corresponding to that row but
564 also the succeeding rows because of Hi-C-LSTM’s sequential behavior. The LSTM weights remain
565 unchanged, but as the input to the LSTM is modified, the inferred Hi-C contact probability is al-
566 tered based on the information retained by the LSTM about the relationship between the sequence
567 elements under contention and chromatin structure.

568 In a CTCF orientation replacement experiment, we replaced the representations of downstream-
569 facing CTCF motifs with the genome-wide average of the upstream-facing motifs and vice versa.
570 This was done under the assumption that the average representation of the given orientation would
571 encapsulate the important information regarding the role played by the orientation in chromatin
572 conformation.

573 Our duplication experiment was carried out by creating a tandem duplication the representa-
574 tions from the 2.1 Mbp region between 67.95 Mbp to 70.08 Mbp in chromosome 7 region [57] and
575 then passing the resulting representation matrix to the LSTM to infer contacts. Given our Hi-C res-
576 olution of 10 kbp, the duplicated region corresponds 214 bins, i.e., bin 6795 to bin 7008. Specifically,
577 the duplicated representation matrix \hat{R}_i is defined as $\hat{R} := [R_{1:6794}, R_{6795:7008}, R_{6795:7008}, R_{7009:N}]$.

578 To enable comparison to Hi-C data mapped to the original pre-duplication reference genome,
579 we combined inferred contacts from both copies. This combination is required because Hi-C reads
580 cannot be disambiguated between the two duplicated copies when they are mapped to the reference
581 genome. Specifically, we passed the predicted contact probability cp through the inverse exponential
582 transformation to define predicted read counts $CS = \frac{1}{-\log cp/a} - \delta$ (see Eq. 1). We summed
583 predicted read counts from the two duplicated copies to simulate mapping reads from both copies
584 to the same reference genome CS' , then re-applied the exponential transform to obtain predicted
585 contact probability cp' .

586 Our baseline for the quantitative evaluation was the original pre-duplication Hi-C for the in-
587 teractions between the upstream, downstream and duplicated regions, and the genomic average for
588 the interactions of the duplicated region with itself. We considered a window of 214 bins (length
589 of the duplicated region), and computed the average genomic contact strength for the bins with
590 themselves in a window of this size.

591 Data Availability

592 The data that support the findings of this study are publicly available to download and are referenced in the
593 bibliography. Refer to Methods for more details. The data and representations generated from the project
594 can be found in the GitHub Repository (<https://github.com/smaslova/HiCLSTM>).

595 Code Availability

596 The code repository for this project, including training, evaluation, data handling can be found in our GitHub
597 repository (<https://github.com/smaslova/HiCLSTM>).

598 Supplementary Results

599 The additional file contains supplementary figures of salient features of Hi-C-LSTM predictions, ablation
600 experiments with Hi-C-LSTM, parameter search for the XGBoost classifier, confusion matrix for classification
601 of subcompartments, and feature importance for Segway-GBR labels.

602 References

- 603 [1] Van Berkum, N. L. et al. Hi-C: a method to study the three-dimensional architecture of genomes. *JoVE*
604 (*Journal of Visualized Experiments*) **39**, e1869 (2010).
- 605 [2] Bengio, Y., Courville, A. & Vincent, P. Representation learning: A review and new perspectives. *IEEE*
606 *transactions on pattern analysis and machine intelligence* **35**, 1798-828 (2013).
- 607 [3] Seide, F., Li, G. & Yu, D. Conversational speech transcription using context-dependent deep neural
608 networks. In Twelfth annual conference of the *International speech communication association* (2011).

- 609 [4] Boulangier-Lewandowski, N., Bengio, Y. & Vincent, P. Modeling temporal dependencies in high-
610 dimensional sequences: Application to polyphonic music generation and transcription. Preprint at
611 arXiv:1206.6392 (2012).
- 612 [5] Krizhevsky, A., Sutskever, I. & Hinton, G. E. Imagenet classification with deep convolutional neural
613 networks. *Communications of the ACM* **60**, 84-90 (2017).
- 614 [6] Schwenk, H., Rousseau, A. & Attik, M. Large, pruned or continuous space language models on a gpu
615 for statistical machine translation. In Proceedings of the *NAACL-HLT 2012 Workshop: Will We Ever*
616 *Really Replace the N-gram Model? On the Future of Language Modeling for HLT* (2012).
- 617 [7] Le, H. S., Oparin, I., Allauzen, A., Gauvain, J. L. & Yvon, F. Structured output layer neural network
618 language models for speech recognition. *IEEE Transactions on Audio, Speech, and Language Processing*
619 **21**, 197-206 (2012).
- 620 [8] Glorot, X., Bordes, A. & Bengio, Y. Domain adaptation for large-scale sentiment classification: A deep
621 learning approach. In *ICML* (2011).
- 622 [9] Koo, P. K. & Eddy, S. R. Representation learning of genomic sequence motifs with convolutional neural
623 networks. *PLoS computational biology* **15**, e1007560 (2019).
- 624 [10] Agarwal, V., Reddy, N. & Anand, A. Unsupervised Representation Learning of DNA Sequences. Preprint
625 at arXiv:1906.03087 (2019).
- 626 [11] Xiong, K. & Ma, J. Revealing Hi-C subcompartments by imputing inter-chromosomal chromatin inter-
627 actions. *Nature communications* **10**, (2019).
- 628 [12] Ashoor, H. et al. Graph embedding and unsupervised learning predict genomic sub-compartments from
629 HiC chromatin interaction data. *Nature communications* **11**, 1 (2020).
- 630 [13] Zhang, R., Zou, Y. & Ma J. Hyper-SAGNN: a self-attention based graph neural network for hypergraphs.
631 Preprint at arXiv:1911.02613 (2019).
- 632 [14] Zhang, R. & Ma, J. Probing multi-way chromatin interaction with hypergraph representation learning.
633 *Cell Systems* **10**, 397-407 (2020).
- 634 [15] Sutskever, I., Vinyals, O. & Le, Q. V. Sequence to sequence learning with neural networks. In *Advances*
635 *in neural information processing systems* (2014).
- 636 [16] Tang, J., Qu, M., Wang, M., Zhang, M., Yan, J. & Mei, Q. Line: Large-scale information network
637 embedding. In Proceedings of the *24th international conference on world wide web* (2015).
- 638 [17] Sanborn, A. L. et al. Chromatin extrusion explains key features of loop and domain formation in wild-
639 type and engineered genomes. Proceedings of the *National Academy of Sciences* **112**, E6456-65 (2015).
- 640 [18] Rao, S. S. et al. Cohesin loss eliminates all loop domains. *Cell* **171**, 305-20 (2017).
- 641 [19] Imaikaev, M. et al. Iterative correction of Hi-C data reveals hallmarks of chromosome organization.
642 *Nature methods* **9**, 999-1003 (2012).
- 643 [20] Rao, S. S. et al. A 3D map of the human genome at kilobase resolution reveals principles of chromatin
644 looping. *Cell* **159**, 1665-80 (2014).
- 645 [21] Cristescu, B. C, Borsos, Z., Lygeros, J., Martínez, M. R. & Rapsomaniki, M. A. Inference of the three-
646 dimensional chromatin structure and its temporal behavior. Preprint at arXiv:1811.09619 (2018).
- 647 [22] Zhu, Y. et al. Constructing 3D interaction maps from 1D epigenomes. *Nature communications* **7**, 1
648 (2016).
- 649 [23] Al Bkhetan, Z. & Plewczynski, D. Three-dimensional epigenome statistical model: genome-wide chro-
650 matin looping prediction. *Scientific reports* **8**, 1 (2018).
- 651 [24] Zhang, S., Chasman, D., Knaack, S. & Roy, S. In silico prediction of high-resolution Hi-C interaction
652 matrices. *Nature communications* **10**, 1 (2019).

- 653 [25] Li, W., Wong, W. H. & Jiang, R. DeepTACT: predicting 3D chromatin contacts via bootstrapping deep
654 learning. *Nucleic acids research* **47**, e60 (2019).
- 655 [26] Schreiber, J., Libbrecht, M., Bilmes, J. & Noble, W. S. Nucleotide sequence and DNaseI sensitivity are
656 predictive of 3D chromatin architecture. Preprint at <https://doi.org/10.1101/103614> (2017).
- 657 [27] Farré, P. & Emberly, E. A maximum-entropy model for predicting chromatin contacts. *PLoS computa-*
658 *tional biology* **14**, e1005956 (2018).
- 659 [28] Farré, P., Heurteau, A., Cuvier, O. & Emberly, E. Dense neural networks for predicting chromatin
660 conformation. *BMC bioinformatics* **19**, 1-12 (2018).
- 661 [29] Fudenberg, G., Kelley, D. R. & Pollard, K. S. Predicting 3D genome folding from DNA sequence with
662 Akita. *Nature Methods* **17**, 1111-1117 (2020).
- 663 [30] Di Pierro, M., Cheng, R. R., Aiden, E. L., Wolynes, P. G. & Onuchic, J. N. De novo prediction of
664 human chromosome structures: Epigenetic marking patterns encode genome architecture. *Proceedings*
665 *of the National Academy of Sciences* **114**, 12126-31 (2017).
- 666 [31] Hochreiter, S. & Schmidhuber, J. Long short-term memory. *Neural computation* **9**, 1735-80 (1997).
- 667 [32] Qi, H. Y, Zhang, Z. J., Li, Y. J. & Fang, X. D. Role of chromatin conformation in eukaryotic gene
668 regulation. *Yi chuan= Hereditas* **33**, 1291-9 (2011).
- 669 [33] Rhind, N. & Gilbert, D. M. DNA replication timing. *Cold Spring Harbor perspectives in biology* **5**,
670 a010132 (2013).
- 671 [34] Ryba, T. et al. Evolutionarily conserved replication timing profiles predict long-range chromatin inter-
672 actions and distinguish closely related cell types. *Genome research* **20**, 761-70 (2010).
- 673 [35] Dileep, V., Ay, F., Sima, J., Vera, D. L., Noble, W. S. & Gilbert, D. M. Topologically associating domains
674 and their long-range contacts are established during early G1 coincident with the establishment of the
675 replication-timing program. *Genome research* **25**, 1104-13 (2015).
- 676 [36] Du, Q. et al. Replication timing and epigenome remodelling are associated with the nature of chromo-
677 somal rearrangements in cancer. *Nature communications* **10**, 1-5 (2019).
- 678 [37] Zheng, H. & Xie, W. The role of 3D genome organization in development and cell differentiation. *Nature*
679 *Reviews Molecular Cell Biology* **13**, 1 (2019).
- 680 [38] Mora, A., Sandve, G. K., Gabrielsen, O. S. & Eskeland, R. In the loop: promoter–enhancer interactions
681 and bioinformatics. *Briefings in bioinformatics* **17**, 980-95 (2016).
- 682 [39] Krivega, I. & Dean, A. Enhancer and promoter interactions—long distance calls. *Current opinion in*
683 *genetics & development* **22**, 79-85 (2012).
- 684 [40] Dong, X., Li, C., Chen, Y., Ding, G. & Li, Y. Human transcriptional interactome of chromatin contribute
685 to gene co-expression. *BMC genomics* **11**, 1-5 (2010).
- 686 [41] Schmitt, A. D. et al. A compendium of chromatin contact maps reveals spatially active regions in the
687 human genome. *Cell reports* **17**, 2042-59 (2016).
- 688 [42] Beagan, J. A. & Phillips-Cremins, J. E. On the existence and functionality of topologically associating
689 domains. *Nature Genetics* **10**, 1-9 (2020).
- 690 [43] Chen, T. & Guestrin, C. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm*
691 *sigkdd international conference on knowledge discovery and data mining* 785-794 (2016).
- 692 [44] Hoffman, M. M., Buske, O. J., Wang, J., Weng, Z., Bilmes, J. A. & Noble, W. S. Unsupervised pattern
693 discovery in human chromatin structure through genomic segmentation. *Nature methods* **9**, 473 (2012).
- 694 [45] Bintu, B. et al. Super-resolution chromatin tracing reveals domains and cooperative interactions in
695 single cells. *Science* **362**, 6413 (2018).

- 696 [46] Hannig, J., Giese, H., Schweizer, B., Amstein, L., Ackermann, J. & Koch, I. isiKnock: in silico knockouts
697 in signaling pathways. *Bioinformatics* **35**, 892-4 (2019).
- 698 [47] Verma, R., Pradhan, D., Maseet, M., Singh, H., Jain, A. K. & Khan, L. A. Genome-wide screening
699 and in silico gene knockout to predict potential candidates for drug designing against *Candida albicans*.
700 *Infection, Genetics and Evolution* **80**, 104196 (2020).
- 701 [48] Bintener, T., Pacheco, M. P. & Sauter, T. Towards the routine use of in silico screenings for drug
702 discovery using metabolic modelling. *Biochemical Society Transactions* **5**, BST20190867 (2020).
- 703 [49] Scheidel, J., Amstein, L., Ackermann, J., Dikic, I. & Koch, I. In silico knockout studies of xenophagic
704 capturing of salmonella. *PLoS computational biology* **12**, e1005200 (2016).
- 705 [50] Cuddapah, S., Jothi, R., Schones, D. E., Roh, T. Y., Cui, K. & Zhao, K. Global analysis of the insulator
706 binding protein CTCF in chromatin barrier regions reveals demarcation of active and repressive domains.
707 *Genome research* **19**, 24-32 (2009).
- 708 [51] Dixon, J. R. et al. Topological domains in mammalian genomes identified by analysis of chromatin
709 interactions. *Nature* **485**, 376-80 (2012).
- 710 [52] Xie, X., Mikkelsen, T. S., Gnirke, A., Lindblad-Toh, K., Kellis, M. & Lander, E. S. Systematic discovery
711 of regulatory motifs in conserved regions of the human genome, including thousands of CTCF insulator
712 sites. *Proceedings of the National Academy of Sciences* **104**, 7145-50 (2007).
- 713 [53] Hou, C., Zhao, H., Tanimoto, K. & Dean, A. CTCF-dependent enhancer-blocking by alternative chro-
714 matin loop formation. *Proceedings of the National Academy of Sciences* **105**, 20398-403 (2008).
- 715 [54] Phillips, J. E. & Corces, V. G. CTCF: master weaver of the genome. *Cell* **137**, 1194-211 (2009).
- 716 [55] Splinter, E. et al. CTCF mediates long-range chromatin looping and local histone modification in the
717 β -globin locus. *Genes & development* **20**, 2349-54 (2006).
- 718 [56] Guo, Y. et al. CRISPR inversion of CTCF sites alters genome topology and enhancer/promoter function.
719 *Cell* **162**, 900-10 (2015).
- 720 [57] Melo, U. S. et al. Hi-C identifies complex genomic rearrangements and TAD-shuffling in developmental
721 diseases. *The American Journal of Human Genetics* **106**, 872-884 (2020).
- 722 [58] Nora, E. P. et al. Targeted degradation of CTCF decouples local insulation of chromosome domains
723 from genomic compartmentalization. *Cell* **169**, 930-944 (2017).
- 724 [59] Kubo, N. et al. Promoter-proximal CTCF binding promotes distal enhancer-dependent gene activation.
725 *Nature structural & molecular biology* **28**, 152-161 (2021).
- 726 [60] Vaswani, A. et al. Attention is all you need. In *Advances in neural information processing systems*
727 (2017).
- 728 [61] Genome Reference Consortium Human Build 37 (GRCh37). BioProject: PRJNA31257. Accessed Jan
729 2020.
- 730 [62] Knight, P. A. & Ruiz, D. A fast algorithm for matrix balancing. *IMA Journal of Numerical Analysis*
731 **33**, 1029-1047 (2013).
- 732 [63] Friedman, N., Linial, M., Nachman, I. & Pe'er, D. Using Bayesian networks to analyze expression data.
733 *Journal of computational biology* **7**, 601-620 (2000).
- 734 [64] Pe'er, D., Regev, A., Elidan, G. & Friedman, N. Inferring subnetworks from perturbed expression
735 profiles. *Bioinformatics* **17**, S215-S224 (2001).
- 736 [65] Whalen, S., Truty, R. M. & Pollard, K. S. Enhancer-promoter interactions are encoded by complex
737 genomic signatures on looping chromatin. *Nature genetics* **48**, 488-96 (2016).
- 738 [66] Schreiber, J., Durham, T., Bilmes, J. & Noble, W. S. Avocado: a multi-scale deep tensor factorization
739 method learns a latent representation of the human epigenome. *Genome biology* **21**, 1-18 (2020).

- 740 [67] Marchal, C. et al. Genome-wide analysis of replication timing by next-generation sequencing with E/L
741 Repli-seq. *Nature protocols* **13**, 819-39 (2018).
- 742 [68] Grant, C. E., Bailey, T. L. & Noble, W. S. FIMO: scanning for occurrences of a given motif. *Bioinform-*
743 *atics* **27**, 1017-1018 (2011).
- 744 [69] Liu, T. et al. TADKB: Family classification and a knowledge base of topologically associating domains.
745 *BMC genomics* **20**, 1-17 (2019).
- 746 [70] Pascanu, R., Mikolov, T. & Bengio, Y. On the difficulty of training recurrent neural networks. In
747 *International conference on machine learning* (2013).
- 748 [71] Elman, J. L. Finding structure in time. *Cognitive science* **14**, 179-211 (1990).
- 749 [72] Sutskever, I., Vinyals, O. & Le, Q. V. Sequence to sequence learning with neural networks. In *Advances*
750 *in neural information processing systems* (2014).
- 751 [73] Lu, L., Zhang, X., Cho, K. & Renals, S. A study of the recurrent neural network encoder-decoder
752 for large vocabulary speech recognition. In Sixteenth Annual Conference of the *International Speech*
753 *Communication Association* (2015).
- 754 [74] Young, T., Hazarika, D., Poria, S. & Cambria, E. Recent trends in deep learning based natural language
755 processing. *IEEE Computational intelligence magazine* **13**, 55-75 (2018).
- 756 [75] Ba, J. L., Kiros, J. R. & Hinton, G. E. Layer normalization. Preprint at arXiv:1607.06450 (2016).
- 757 [76] Schuster, M. & Paliwal, K. K. Bidirectional recurrent neural networks. *IEEE transactions on Signal*
758 *Processing* **45**, 2673-81 (1997).
- 759 [77] Glorot, X. & Bengio, Y. Understanding the difficulty of training deep feedforward neural networks. In
760 *Proceedings of the thirteenth international conference on Artificial intelligence and statistics* (2010).
- 761 [78] Kingma, D. P. & Ba, J. Adam: A method for stochastic optimization. Preprint at arXiv:1412.6980
762 (2014).
- 763 [79] Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I. & Salakhutdinov, R. Dropout: a simple way to
764 prevent neural networks from overfitting. *The journal of machine learning research* **15**, 1929-58 (2014).
- 765 [80] Sundararajan, M., Taly, A. & Yan, Q. Axiomatic attribution for deep networks. Preprint at
766 arXiv:1703.01365 (2017).

767 **Authors' contributions**

768 K.B.D. led ideation, genomic data processing, building and validating the machine learning model, wrote the
769 first draft of the manuscript, and edited the manuscript. A.M. contributed towards ideation, data processing,
770 parallelization of the model and model validation. V.K.B. partially funded the project. M.W.L. supervised
771 the project. All authors participated in the design of the study, the interpretation of results, and editing the
772 manuscript. All authors read and approved the final manuscript.

773 **Competing interests**

774 The authors declare that they have no competing interests.

775 **Materials and Correspondence**

776 Any correspondence and material requests can be addressed to the corresponding authors, Kevin B. Dsouza
777 or Maxwell W. Libbrecht.

778 **Authors' information**

779 K.B.D. is a PhD student at UBC where he works on computational genomics jointly under the information
780 theory group at UBC and computational biology group at The Simon Fraser University. His work focuses on
781 building machine learning tools to aid in the understanding of structural and functional genomic data. A.M.
782 holds an MSc in Bioinformatics from the University of British Columbia and is currently a PhD student at
783 Simon Fraser University where she works in the computation biology group. Her research focuses on the use
784 of machine learning approaches in the analysis for genomics data. E.A-J. is a PhD student at Imperial College
785 London studying 3D genome organisation and gene expression. M.M. is a Career Scientist at the MRC's
786 Clinical Sciences Centre at Imperial College. He is a molecular immunologist whose work has been central to
787 the understanding of development, and cellular reprogramming. V.K.B. is a Fellow of the IEEE, The Royal
788 Society of Canada, and currently a Professor in the Department of Electrical and Computer Engineering at
789 the University of British Columbia in Vancouver. M.W.L. is an Assistant Professor in Computing Science
790 at Simon Fraser University where his research focuses on developing machine learning methods applied to
791 high-throughput genomics data sets.

792 **Figures**

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SupplementaryResultsHiCLSTM.pdf](#)
- [rs.pdf](#)