

1 **DeepC: predicting 3D genome folding using megabase-scale transfer**

2 **learning**

3

4 Ron Schwessinger^{1,2,3}, Matthew Gosden¹, Damien Downes¹, Richard C Brown³, A. Marieke
5 Oudelaar^{1,2}, Jelena Telenius², Yee Whye Teh⁴, Gerton Lunter^{*2,3} and Jim R. Hughes^{*1,2}

6

7 ¹MRC Molecular Haematology Unit & ²MRC WIMM Centre for Computational Biology

8 MRC Weatherall Institute of Molecular Medicine, University of Oxford, Oxford, UK

9 ³Wellcome Centre for Human Genetics, University of Oxford, Oxford, UK

10 ⁴Department of Statistics, University of Oxford, Oxford, UK

11

12 *Corresponding Authors are Gerton Lunter gerton.lunter@well.ox.ac.uk and Jim R. Hughes

13 jim.hughes@imm.ox.ac.uk

14 **Abstract**

15 Predicting the impact of non-coding genetic variation requires interpreting it in the context of 3D
16 genome architecture. We have developed deepC, a transfer learning based deep neural network that
17 accurately predicts genome folding from megabase-scale DNA sequence. DeepC predicts domain
18 boundaries at high-resolution, learns the sequence determinants of genome folding and predicts the
19 impact of both large-scale structural and single base pair variations.

20 **Introduction**

21 Most genetic variants associated with common diseases affect gene regulatory regions distal to target
22 genes^{1,2}. Genome 3D structure is central to mediating these functional interactions, but its intricately
23 convoluted and large-scale nature renders it challenging to understand and predict. Proposed machine
24 learning and polymer modelling approaches to predict 3D genome structure have produced promising
25 results, but none effectively integrates across resolutions. Methods that use information at the base
26 pair level focus on window-to-window-based predictions³⁻⁵, while methods that incorporate a large
27 genomic context do so by coarse segregation into genomic features⁶⁻⁸ or polymer beads^{9,10}, thus
28 compromising their ability to predict the impact of variation at base pair resolution.

29 We propose that to accurately predict topologically associated domains (TADs) a model needs to
30 capture sequence patterns across large genomic distances. Regulatory elements can interact over
31 megabase distances and boundary elements lying in between may alter chromatin contacts
32 significantly. A chromatin interaction model should thus integrate information at the megabase-scale.
33 However, to predict the impact of genetic variation the model must also learn to interpret DNA
34 sequence at base pair resolution. Since 3D genome interactions are determined by genomic regulatory
35 elements such as CTCF bound domain boundaries¹¹, a model that has learned the grammar of
36 regulatory elements could help guide the prediction of 3D genome structure.

37 Based on these ideas, we developed deepC, a deep neural network that bridges the gap from base
38 pairs to TADs. DeepC uses a transfer learning approach and tissue-specific Hi-C data to train models
39 that predict genome folding from megabase (Mb) windows of DNA sequence (Fig. 1a). The trained
40 models can then be used to predict chromatin domain boundaries at high resolution and to identify
41 the sequence determinants of genome folding. Importantly, they allow us to predict the impact of
42 genetic variants from large structural variations down to single-nucleotide polymorphisms.

43

44 Results

45 A deep learning model for predicting chromatin interactions from megabase-scale DNA

46 We encode Hi-C data as a vector of pairwise interaction values between 5 kilobase (kb) genomic bins
47 at distances of up to ~1 Mb (Fig. 1b). DeepC learns to predict these contact frequencies taking as
48 input the underlying ~1 Mb window of DNA sequence. The deepC network architecture is
49 constructed from a convolutional module with max-pooling that has proven powerful for predicting
50 chromatin features from DNA sequence^{12,13}. This is followed by a dilated convolutional module that
51 excels at incorporating large-scale context while maintaining resolution^{14–16}. Finally, a fully
52 connected layer integrates the detected patterns over a megabase of DNA sequence to predict
53 chromatin folding.

54 We found two factors to be crucial for deepC’s effective learning and generalization. First, we
55 percentile-normalize the raw contact frequency signal in Hi-C data by genomic distance (Extended
56 Data Figure 1, Methods), termed the “skeleton”. This normalization reveals informative longer-range
57 interactions and enhances the contrast at domain boundaries. Second, we employ transfer learning¹⁷
58 (Fig. 1a, Extended Data Figure 2), a concept that has proven powerful in deep learning applications
59 for image analysis and natural language processing. In a first phase of training the initial
60 convolutional module learns to predict a compendium of chromatin features such as open chromatin
61 regions and CTCF binding sites across cell types^{12,18,19}. Next, the convolutional module is stripped of
62 the fully connected layer responsible for interpretation. Only the learned sequence patterns are
63 transferred to the second training phase where they are refined, and the dilated module and fully
64 connected layer are trained *ab initio* to predict chromatin interactions. The same weights pre-trained
65 on a chromatin feature compendium across cell types are used for transfer learning irrespective of the
66 cell type of the Hi-C data source.

67 We trained deepC models on seven human²⁰ and one mouse²¹ Hi-C data sets with different sequencing
68 depths and at different resolutions (Supplementary Figure 1 - 4). We focused our analysis on the
69 primary GM12878 (~3.6 B reads) and K562 (~1.3 B reads) data, training models at 5 kb resolution.
70 DeepC yields smooth but detailed predictions that resolve the hierarchical nature of TADs and
71 insulated domains (Fig. 1c, Supplementary Figure 1). In a cross-validation scheme across all
72 chromosomes in GM12878 (Fig. 1d), deepC achieves an average, distance-stratified Pearson
73 correlation between predictions and Hi-C skeleton of ~0.36 on raw skeleton data and ~0.57 when
74 applying a small smoothing filter to the discrete and noisy skeleton (~0.28 and ~0.51 in K562,
75 Supplementary Figure 5). We compared deepC to a recently proposed random-forest-based method
76 HiC-Reg⁸, that predicts chromatin interactions up to 1 Mb distance using chromatin features of the
77 interacting windows and the window in between rather than using DNA sequence as input. We
78 observed that deepC generalizes better in predicting the domain structure of unseen chromosomes
79 (Supplementary Figure 6 and 7).

80 Although we focused our main analysis on the deeply sequenced Hi-C data from Rao et al.²⁰ we
81 hypothesized that deepC is able to predict chromatin interactions after training on data with
82 significantly lower sequencing depth. To this end we trained GM12878 models with Hi-C data
83 downsampled from originally ~2.6 B to 1 B, 100 M and 10 M valid Hi-C contacts respectively
84 (Supplementary Figure 8). The 1 B and 100 M contact models still learned to predict chromatin
85 structure, with a mean Pearson correlation of 0.46 and 0.4 between the 1 B and 100 M model on hold
86 out chromosomes 16 and 17 respectively. In contrast the 10 M model failed to learn chromatin
87 structure. While deepC can predict chromatin interactions from less deeply sequenced samples, we
88 do note that dedicated methods have been proposed for increasing the resolution of Hi-C maps that
89 take as input low resolution maps directly^{22,23} rather than predicting from DNA sequence.

90 We train separate models for each Hi-C data set derived from a different cell type. These distinct
91 models learn tissue-specific chromatin interactions (Extended Data Figure 3, Supplementary Figure

92 9). We also trained a deepC model to predict chromatin interactions in multiple cell types jointly, but
93 we found that the jointly trained network captured the tissue-specific Hi-C patterns less well
94 compared to the individually trained models (Supplementary Figure 10).

95 **Validating deepC with high sensitivity chromosome confirmation capture**

96 We next sought to validate deepC predictions with an independent set of chromatin interactions. To
97 this end, we utilized NG Capture-C²⁴ (Methods), which generates high-resolution interaction data
98 from targeted viewpoints and identifies chromatin interactions at higher sensitivity than Hi-C. We
99 captured the interactions of 220 viewpoints in two cell types (GM12878 and K562) covering 81 CTCF
100 sites and 139 sites lying within insulated domains, not overlapping with regulatory elements to
101 capture the domain structure. We observed good agreement between the predicted domain structure
102 and interaction peaks in the NG Capture-C tracks (Extended Data Figure 4), showing deepC is
103 capable of predicting true biophysical boundaries that are evident in these sensitive 3C assays but
104 poorly captured by the original Hi-C, especially at lower sequencing depth.

105 Although deepC effectively captures the positions of boundary elements some aspects of interactions
106 are not captured fully by the model. When comparing the virtual 4C track from the Hi-C skeleton,
107 the predictions and distance-normalized NG Capture-C tracks from CTCF viewpoints
108 (Supplementary Figure 11) we saw that the NG Capture-C correlated better with the Hi-C skeleton
109 than with the predictions (Fig. 2a, Pearson correlation in GM12878: 0.59 vs 0.30; K562: 0.55 vs
110 0.37). This was due to the tendency of deepC predictions to de-emphasize the characteristic punctate
111 nature of signal at the apex of interacting CTCF elements. This may be explained by an inability of
112 deepC to model detailed characteristics of the loop extrusion mechanism such as cohesin processivity,
113 which may not be encoded in the local DNA sequence and more dependent on factors such as nuclear
114 concentration of extruding factors²⁵. In contrast, for intra domain viewpoints deepC correlates equally
115 well with NG Capture-C data as NG Capture-C correlates with the Hi-C skeleton (Fig. 2a, GM12878:

116 0.30 vs 0.36; K562: 0.46 vs 0.42). Therefore, even though deepC is predicting interactions from
117 sequence in these instances it performs as well as comparing two different experimental sources of
118 3C data.

119 Taken together, these analyses suggested that deepC is capable of modelling the DNA encoded
120 signals that determine the activity and position of boundary elements. To test this, we called
121 boundaries within 1 Mb from the NG Capture-C viewpoints using the Hi-C data, the skeleton and the
122 deepC predictions respectively (Supplementary Figure 12 and 13) using the established insulation-
123 score-based approach²⁶ with parameters for high-resolution calling (Methods). We then compared the
124 called boundaries from these three sources with the high-sensitivity NG Capture-C 3C data to
125 quantify the enrichment of chromatin interactions (Fig. 2b, Supplementary Figure 14). We observed
126 clear enrichment over the deepC predicted boundaries indicating that they on average represent
127 biophysical barriers to genome interactions. In contrast, the boundaries called directly from the Hi-C
128 data and from the Hi-C skeleton showed less pronounced enrichment suggesting that calling directly
129 from the data, on average, captures boundaries less effectively at the available sequencing depths.
130 We confirmed these results with boundaries called using TopDom²⁷ (Supplementary Figure 15), a
131 TAD caller that showed best overall robustness in a recent bench marking study²⁸.

132 To visualize the coherence of the deepC predictions and called boundaries across specific loci at a
133 sensitivity higher than the available Hi-C data, we utilized Tiled-C²⁹ (Methods), which generates Hi-
134 C like data for specific loci at high sensitivity and resolution. We performed Tiled-C for a selection
135 of loci where deepC predicted fine-grained boundaries (Fig. 2c, Supplementary Figure 16) or cell-
136 type-specific patterns (Extended Data Figure 3, Supplementary Figure 9). In line with the enrichment
137 analysis, we confirmed that when called at high resolution, deepC boundaries are evident and align
138 well with the boundaries in this highly sensitive 3C data. The added benefit for boundary calling is
139 particularly striking in the comparatively lower-coverage K562 Hi-C data (~1.3 B reads) (Fig. 2c,
140 Supplementary Figure 16).

141 When comparing the overall structure of deepC predictions to the Hi-C and Tiled-C data we observed
142 that deepC tends to predict inter domain interactions in the form of stripes and dots more pronounced,
143 some of which are only faintly detectable in Hi-C and Tiled-C data and some appear novel (Fig. 1c,
144 2c, Supplementary Figure 9). This suggests that deepC tends to underestimate the insulation between
145 domains. Future refinements to the model architecture might be able to better capture the inter domain
146 insulation. The effect appears amplified when comparing skeleton transformed to raw data as
147 necessary for Tiled-C.

148 **Dissecting the sequence determinants of genome folding**

149 DeepC allows us to dissect the sequence determinants of genome folding at base pair resolution. To
150 estimate the relative importance of every base pair for predicting chromatin interactions we employed
151 the saliency score as a computationally efficient method adapted from image analysis³⁰. The saliency
152 score estimates how much the interaction prediction depends on each single base pair by calculating
153 the gradient of the model output with respect to the sequence input (Methods). The saliency score
154 predicts important regions and highlights transcription factor motifs within them (Extended Data
155 Figure 5).

156 Genome-wide we identify sharp saliency peaks at CTCF sites and broader saliency peaks at active
157 promoters (Fig. 3a). As bases with high saliency scores mark positions predicted to be important for
158 chromatin architecture, we hypothesized that mutations within these regions would be enriched for
159 those affecting gene expression. To test this, we retrieved 6607 GM12878 cell-type-specific eQTLs
160 (GTEx v7) that are located in open chromatin (DNase-seq) or CTCF sites (CTCF ChIP-seq,
161 ENCODE) and are thus likely to lie in regulatory elements. We found that these eQTLs have
162 significantly higher saliency scores than SNPs randomly re-sampled from the same regions (p-value
163 $< 1e-85$ using a two-sample Kolmogorov-Smirnow test) (Supplementary Figure 17). This suggests
164 that the deepC saliency score can be used to fine map eQTLs when expression changes are mediated

165 through an impact on chromatin architecture.

166 A long-standing question has been which functional elements within the genome underlie the patterns
167 of genome folding. To investigate this, we performed an *in silico* deletion screen of all active elements
168 genome-wide and used deepC to assess their importance for chromatin interactions (Fig. 3b). As
169 expected, we find that deleting CTCF sites as well as enhancers and promoters with proximal CTCF
170 binding has the strongest average predicted impact. We also find promoter and enhancers without
171 proximal CTCF binding sites to be important, with deletions of promoters on average having a
172 stronger effect. In addition, deletions of promoters and enhancers with strong activity-associated
173 histone marks have a higher predicted impact than those without such marks.

174 Our analysis suggests that in addition to known factors such as CTCF binding and orientation, active
175 regulatory elements, in particular promoters, are critical elements for effectively predicting genome
176 interactions. Furthermore, deepC predicts boundaries not associated with CTCF sites. Taken together
177 this indicates that deepC has learned aspects of a complex grammar of genome folding beyond CTCF
178 motifs and their relative orientation and suggests a causal role for regulatory elements in defining and
179 stabilizing 3D genome structure.

180 **Predicting the impact of sequence variation on genome folding**

181 To test deepC's ability to predict the impact of sequence variation on genome folding we utilized two
182 well-characterized examples from the literature. Hnisz et al. showed that a ~30 kb CRISPR-mediated
183 deletion in HEK293T cells, encompassing four CTCF sites at the *LMO2* locus, leads to a local
184 rearrangement of the chromatin structure as confirmed by 5C³¹. Computationally reproducing this
185 deletion in GM12878 (Fig. 3c) and K562 cells (Supplementary Figure 18) recapitulates the domain
186 fusion observed by Hnisz and colleagues. Furthermore, using deepC we computationally deleted the
187 CTCF sites individually, predicting that no single deletion alone is sufficient for causing the
188 rearrangement (Supplementary Figure 19), suggesting multiple redundant boundaries at this region.

Crucially, our sequence-based model can predict the impact of single base pair variants. To demonstrate this, we tested two asthma-risk associated SNPs³² shown to impact *ORMDL3* expression in immune cells. Schmiedel et al. demonstrated that the SNPs impact CTCF binding sites, disrupting enhancer-promoter interactions of *ORMDL3* in CD4-positive T-cells. DeepC predictions recapitulate the loss of a boundary element, insulating *ORMDL3* from downstream interactions (Fig. 3d). When testing the individual SNPs Supplementary Figure 20 a and b), deepC predicts the rs12936231 risk allele to have a strong effect on genome folding (mean absolute interaction difference 0.176). In contrast, although the rs4065275 risk-allele suggests a boundary creating effect, the predicted strength is weak (0.006). To put these predicted SNP effects into context, we compared them to the *in silico* deletion screen effects (Supplementary Figure 20c). The effect of rs12936231 lies above the 25th percentile of the 500 bp, weak CTCF site deletions. In addition, we sampled 1000 SNPs from CTCF sites (Fig. 3d) as well as from promoters, enhancers and background sequences (Supplementary Figure 20d). In comparison to sampled 1000 CTCF SNPs, rs4065275 lies within the top 11 % and rs12936231 within the top 1 %. Taken together, deepC prioritizes rs12936231 as the likely causal variant. Interestingly, deepC predicts this effect in GM12878 (immortalized B-cells) but not in K562 (myeloid leukemia cell line with erythroid characteristics) or IMR90 (human embryonic lung fibroblasts), pointing to a potential lymphoid specific effect (Supplementary Figure 21).

207 **Discussion**

208 Mammalian chromatin architecture folds at the megabase and sub-megabase scale constraining distal
209 regulatory interactions within TADs and smaller insulated domains^{11,33,34}. Ultimately, chromatin
210 interactions are encoded in the DNA sequence through an intricate interplay of protein binding sites
211 and other sequence determinants. Understanding the link between individual sequences and large-
212 scale chromatin interactions at base pair resolution is a key challenge for understanding chromatin
213 architecture and its role in gene regulation^{35,36}. We developed deepC to traverse the gap between base
214 pair sequences and megabase structures. DeepC is the first sequence-based deep learning model that
215 predicts chromatin interactions from DNA sequence while integrating a context of megabase scale.
216 This scale of analysis is necessary for accurate prediction of chromatin interactions, which in turn
217 allowed for the determination of the elements driving these interactions and assessment of mutations
218 disrupting them.

219 We found deepC models to yield substantially better predictions when we pre-seeded the model with
220 hidden layers optimized to predict a compendium of chromatin features. This allows deepC to predict
221 intricate chromatin interactions even when trained on low depth, low resolution Hi-C data. By
222 validating the results with NG Capture-C and Tiled-C, each 3C methods capable of extreme depth
223 and sensitivity, we showed that the deepC approach effectively increases the resolution of Hi-C data.
224 We demonstrated that deepC can be used to fine-map the sequence determinants of chromatin
225 architecture at base pair resolution and link these with effects on gene expression. Additionally, our
226 genome wide deletion screen of potential regulatory elements shed light on the mechanics of
227 chromatin interactions. It confirmed that CTCF binding site deletions are most likely to cause strong
228 chromatin interaction changes. Importantly, deepC also indicates that both promoters and enhancers
229 contribute to genome folding, in addition to CTCF. Generally, promoter deletions have a higher
230 predicted effect on genome organization than enhancer deletions, and we find that deletions of

231 enhancer and promoters associated with active chromatin marks have a higher predicted impact than
232 those without such marks. Our observations are in line with findings from orthologous methods, that
233 find CTCF binding, open chromatin, active histone marks and RNA-seq to be most predictive^{4-8,10}.
234 The finding that identifying active promoters and enhancers, in addition to CTCF binding sites, is
235 required to accurately predict 3D genome structures suggests that these elements play an important
236 role in establishing these structures, possibly via recruitment of its components and by actively
237 stabilizing certain loops.

238 We believe deep learning-based genome folding predictions will facilitate chromatin architecture
239 research. In a parallel study, Fudenberg et al.³⁷ have developed an alternative model (Akita) to
240 accurately predict interaction in megabase scale loci. DeepC and Akita have a similar convolutional
241 module as network base but vary significantly in the remaining network structure as well as the data
242 encoding and training scheme. We believe that future comparative study and consolidation between
243 these advances will bring further insights into genome function.

244 Here we present deepC as a valuable tool for dissecting the functional elements that shape chromatin
245 architecture and for predicting the impact of sequence changes from single base pair to structural
246 variants. Furthermore, deepC represents a step towards predictive models of gene regulation that
247 integrate the intricate and long-ranged chromatin landscape of mammalian genomes.

248 **Acknowledgments**

249 The authors would like to thank Dr. Robert Beagrie for help in refining the manuscript. This work
250 was supported by the MRC (MC_UU_12009/14 to J.R.H.) and the Wellcome Trust via Strategic
251 Award (106130/Z/14/Z to J.R.H.) and Institutional Strategic Support Fund (reference 105605/Z/14/Z
252 to J.R.H.). The Wellcome Trust Genomic Medicine and Statistics PhD Programme (203728/Z/16/Z
253 to R.S. & 203141/Z/16/Z to R.B.). The Stevenson Junior Research Fellowship at University College,
254 Oxford (to A.M.O). G.L. is supported by the Wellcome Trust supporting award (090532/Z/09/Z to
255 G.L.). Y.W.T. is supported by the European Research Council under the European Union's Seventh
256 Framework Programme (FP7/2007-2013 to Y.W.T.) ERC grant agreement no. 617071.
257 The Genotype-Tissue Expression (GTEx) Project was supported by the Common Fund of the Office
258 of the Director of the National Institutes of Health, and by NCI, NHGRI, NHLBI, NIDA, NIMH, and
259 NINDS.

261 **Author contributions**

262 R.S., G.L. and J.R.H. conceived the project. R.S, R.B., Y.W.T. and G.L designed the neural network
263 architectures. R.S. optimized and trained the neural networks and performed downstream analysis.
264 R.S. M.G. D.D. A.M.O. and J.R.H. designed and evaluated the validation strategy. M.G. performed
265 NG Capture-C experiments. D.D. performed Tiled-C experiments. R.S., A.M.O. and J.T. performed
266 bioinformatic analysis of NG Capture-C and Tiled-C. R.S. performed integrative analysis and
267 prepared the figures. R.S., G.L and J.R.H. wrote the manuscript with inputs from all authors.

269 **Competing financial interests**

270 The authors declare no competing financial interests.

271 **References**

- 272 1. Hindorff, L. A. *et al.* Potential etiologic and functional implications of genome-wide
273 association loci for human diseases and traits. *Proc. Natl. Acad. Sci. U. S. A.* **106**, 9362–9367
274 (2009).
- 275 2. Maurano, M. T. *et al.* Systematic localization of common disease-associated variation in
276 regulatory DNA. *Science* **337**, 1190–1195 (2012).
- 277 3. Schreiber, J., Libbrecht, M., Bilmes, J. & Noble, W. S. Nucleotide sequence and DNaseI
278 sensitivity are predictive of 3D chromatin architecture. Preprint at:
279 <https://www.biorxiv.org/content/10.1101/103614v5> (2017).
- 280 4. Whalen, S., Truty, R. M. & Pollard, K. S. Enhancer–promoter interactions are encoded by
281 complex genomic signatures on looping chromatin. *Nat. Genet.* **48**, 488–496 (2016).
- 282 5. Li, W., Wong, W. H. & Jiang, R. DeepTACT: predicting 3D chromatin contacts via
283 bootstrapping deep learning. *Nucleic Acids Res.* **47**, e60 (2019).
- 284 6. Qi, Y. & Zhang, B. Predicting three-dimensional genome organization with chromatin states.
285 *PLOS Comput. Biol.* **15**, e1007024 (2019).
- 286 7. Belokopytova, P. S., Nuriddinov, M. A., Mozheiko, E. A., Fishman, D. & Fishman, V.
287 Quantitative prediction of enhancer–promoter interactions. *Genome Res.* **30**, 72–84 (2020).
- 288 8. Zhang, S., Chasman, D., Knaack, S. & Roy, S. In silico prediction of high-resolution Hi-C
289 interaction matrices. *Nat. Commun.* **10**, 5449 (2019).
- 290 9. Buckle, A., Brackley, C. A., Boyle, S., Marenduzzo, D. & Gilbert, N. Polymer Simulations
291 of Heteromorphic Chromatin Predict the 3D Folding of Complex Genomic Loci. *Mol. Cell*
292 **72**, 786-797.e11 (2018).
- 293 10. Bianco, S. *et al.* Polymer physics predicts the effects of structural variants on chromatin
294 architecture. *Nat. Genet.* **50**, 662–667 (2018).
- 295 11. Hnisz, D., Day, D. S. & Young, R. A. Insulated Neighborhoods: Structural and Functional

296 Units of Mammalian Gene Control. *Cell* **167**, 1188–1200 (2016).

297 12. Zhou, J. & Troyanskaya, O. G. Predicting effects of noncoding variants with deep learning–
 298 based sequence model. *Nat. Methods* **12**, 931–934 (2015).

299 13. Kelley, D. R., Snoek, J. & Rinn, J. L. Basset: learning the regulatory code of the accessible
 300 genome with deep convolutional neural networks. *Genome Res.* **26**, 990–999 (2016).

301 14. Kelley, D. R. *et al.* Sequential regulatory activity prediction across chromosomes with
 302 convolutional neural networks. *Genome Res.* **28**, 739–750 (2018).

303 15. Yu, F. & Koltun, V. Multi-Scale Context Aggregation by Dilated Convolutions. Preprint at:
 304 <http://arxiv.org/abs/1511.07122> (2015).

305 16. Oord, A. van den *et al.* WaveNet: A Generative Model for Raw Audio. *2009 IEEE Int. Conf.*
 306 *Acoust. Speech Signal Process.* 3437–3440 (2016).

307 17. Yosinski, J., Clune, J., Bengio, Y. & Lipson, H. How transferable are features in deep neural
 308 networks? *Adv. Neural Inf. Process. Syst.* **4**, 3320–3328 (2014).

309 18. Bernstein, B. E. *et al.* An integrated encyclopedia of DNA elements in the human genome.
 310 *Nature* **489**, 57–74 (2012).

311 19. Kundaje, A. *et al.* Integrative analysis of 111 reference human epigenomes. *Nature* **518**, 317–
 312 330 (2015).

313 20. Rao, S. S. P. *et al.* A 3D Map of the Human Genome at Kilobase Resolution Reveals
 314 Principles of Chromatin Looping. *Cell* **159**, 1665–1680 (2014).

315 21. Bonev, B. *et al.* Multiscale 3D Genome Rewiring during Mouse Neural Development. *Cell*
 316 **171**, 557–572.e24 (2017).

317 22. Zhang, Y. *et al.* Enhancing Hi-C data resolution with deep convolutional neural network
 318 HiCPlus. *Nat. Commun.* **9**, 750 (2018).

319 23. Liu, Q., Lv, H. & Jiang, R. hicGAN infers super resolution Hi-C data with generative
 320 adversarial networks. *Bioinformatics* **35**, i99–i107 (2019).

321 24. Davies, J. O. J. *et al.* Multiplexed analysis of chromosome conformation at vastly improved
322 sensitivity. *Nat. Methods* **13**, 74–80 (2016).

323 25. Fudenberg, G. *et al.* Formation of Chromosomal Domains by Loop Extrusion. *Cell Rep.* **15**,
324 2038–2049 (2016).

325 26. Crane, E. *et al.* Condensin-driven remodelling of X chromosome topology during dosage
326 compensation. *Nature* **523**, 240–244 (2015).

327 27. Shin, H. *et al.* TopDom: An efficient and deterministic method for identifying topological
328 domains in genomes. *Nucleic Acids Res.* **44**, e70 (2015).

329 28. Zufferey, M., Tavernari, D., Oricchio, E. & Ciriello, G. Comparison of computational
330 methods for the identification of topologically associating domains. *Genome Biol.* **19**, 217
331 (2018).

332 29. Oudelaar, A. M. *et al.* Dissection of the 4D chromatin structure of the α -globin locus through
333 in vivo erythroid differentiation with extreme spatial and temporal resolution. Preprint at:
334 <https://www.biorxiv.org/content/10.1101/763763v2> (2019).

335 30. Simonyan, K., Vedaldi, A. & Zisserman, A. Deep Inside Convolutional Networks:
336 Visualising Image Classification Models and Saliency Maps. *2nd Int. Conf. Learn.*
337 *Represent. ICLR 2014 - Work. Track Proc.* (2013).

338 31. Hnisz, D. *et al.* Activation of proto-oncogenes by disruption of chromosome neighborhoods.
339 *Science* **351**, 1454–1458 (2016).

340 32. Schmiedel, B. J. *et al.* 17q21 asthma-risk variants switch CTCF binding and regulate IL-2
341 production by T cells. *Nat. Commun.* **7**, 13426 (2016).

342 33. Robson, M. I., Ringel, A. R. & Mundlos, S. Regulatory Landscaping: How Enhancer-
343 Promoter Communication Is Sculpted in 3D. *Mol. Cell* **74**, 1110–1122 (2019).

344 34. Dixon, J. R., Gorkin, D. U. & Ren, B. Chromatin Domains: The Unit of Chromosome
345 Organization. *Mol. Cell* **62**, 668–680 (2016).

346 35. Marti-Renom, M. A. *et al.* Challenges and guidelines toward 4D nucleome data and model
347 standards. *Nat. Genet.* **50**, 1352–1358 (2018).

348 36. Spielmann, M., Lupiáñez, D. G. & Mundlos, S. Structural variation in the 3D genome. *Nat.*
349 *Rev. Genet.* **19**, 453–467 (2018).

350 37. Fudenberg, G., Kelley, D. R. & Pollard, K. S. Predicting 3D genome folding from DNA
351 sequence. Preprint at: <https://www.biorxiv.org/content/10.1101/800060v1> (2019).

352

353

354

355 **Figure legends**

356

357 **Figure 1** | Predicting Hi-C interactions from DNA sequence. a) Overview of the deepC architecture
358 and training workflow. b) Encoding of Hi-C data as target vector for prediction given a 1 Mb window
359 of DNA sequence. c) Comparison of Hi-C data, the derived Hi-C skeleton and the interactions
360 predicted from DNA sequence using deepC. Shown is a ~ 7 Mb region on hold out chromosome 17.
361 d) Distance-stratified Pearson correlation between the Hi-C skeleton and the deepC predictions in a
362 cross-validation scheme across all chromosomes. Solid lines indicate the mean correlation value and
363 the area indicates the space between the maximum and the minimum values over all chromosomes.
364 Red shows the correlation with the raw and blue with the (5x5) mean filter smoothed skeleton values.
365 Dotted lines at 0.3, 0.4 and 0.5.

366

367 **Figure 2** | Validation of deepC predictions. a) Comparing the correlation between the validation NG
368 Capture-C profiles and the virtual 4C profiles derived from the Hi-C skeleton (red); and the deepC
369 prediction map (blue) from all viewpoints in two cell types. Compared are n = 81 CTCF and n = 139
370 intra domain viewpoints. Boxplots: median middle thick line, 25th and 75th percentile left and right
371 hinge respectively, whiskers stretch up 1.5 times the IQR (inter quartile range). b) Meta-profiles of
372 the average NG Capture-C signal over domain boundaries called at high resolution from the Hi-C
373 data, the skeleton and the deepC predicted interaction map respectively. Shown is the mean distance-
374 normalized NG Capture-C signal relative to the boundary centre. The labels "CTCF" and "Intra
375 domain" refer to the NG Capture-C viewpoint fragments. These were designed to overlap either
376 CTCF sites or to lie within insulated domains but not overlap with regulatory genomic elements, so
377 as to capture the domain structure and not the interactions of specific genomic elements. c) Shown
378 are Hi-C data, the deepC predicted interaction map and the Tiled-C high sensitivity map over a locus
379 on chr17, a hold-out chromosome. Boundaries called at high resolution from Hi-C (green) and deepC
380 predictions (purple) are aligned under the respective map and the Tiled-C map. Cell-type-specific

381 CTCF ChIP-seq tracks are visualized below. For contrast, Hi-C and Tiled-C data were bounded
382 between 5 and 95 % of coverage and deepC predictions were bounded between 2 and 8 predicted
383 regression score.

384

385 **Figure 3** | DeepC for dissecting the determinants of genome folding and predicting the impact of
386 variation. a) Mean saliency enrichment in 10 bp bins over different classes of regulatory elements
387 genome-wide in GM12878 cells. b) Mean predicted chromatin interaction difference between
388 reference and variant when deleting all open chromatin and CTCF sites and shuffled background
389 sequences individually, genome-wide in GM12878 cells. Differences were quantified as the mean
390 absolute interaction difference per bin to bin interaction. Each deletion was classified based on the
391 overlap with genome segmentation classes. A total of $n = 96805$ was analysed. Boxplots show the
392 median as thick line, the 25th and 75th percentile as left and right hinge respectively. The whiskers
393 stretch up 1.5 times the IQR (inter quartile range) away from the respective whisker. Outliers outside
394 this range were omitted from plotting. c) Effect of deleting a ~ 30 kb fragment containing four CTCF
395 sites. Shown is the effect as validated by 5C data in wild-type (WT) and mutant in HEK293 cells by
396 Hnisz et al. and the predicted effect using the GM12878 model. CTCF ChIP-seq tracks from HEK293
397 and GM12878 cells are visualized below. Dotted lines mark the deleted fragment. The mutant deepC
398 prediction was shifted to be centred on the deleted site. d) Predicting the effect of two Asthma
399 associated SNPs with the GM12878 model. Shown is the prediction for non risk and the risk allele,
400 the differential map (non risk – risk) and a GM12878 CTCF ChIP-seq track. Location of the SNPs is
401 indicated by triangles (red – rs12936231, blue – rs4065275). Black lines indicate the location of
402 *ORMDL3* and *IKZF3* and only those two genes are highlighted for clarity. Comparing the predicted
403 SNP effects against 1000 randomly sampled SNPs within CTCF sites places rs4065275 in the top 11
404 % and rs12936231 in the top 1% of predicted mean absolute interaction difference.

405 Online Methods

406 **Chromatin feature data.** As human chromatin feature compendium the ENCODE¹⁸ and Roadmap¹⁹
407 chromatin data utilized in DeepSEA¹² were used. Narrow peak calls (hg19) for 918 experiments were
408 downloaded. The data was supplemented with additional erythroid lineage data. Five sets of ATAC-
409 seq data from Corces et al.³⁸, two DNase-seq experiments³⁹ and ten ATAC-seq and one CTCF ChIP-
410 seq experiment from Downes et al. and in house erythroid differentiation⁴⁰ were used. All data used
411 are listed in Supplementary Table 1. All additional data was aligned to hg19 using the NGseqBasic
412 pipeline⁴¹. Peaks were called with macs2⁴² (default parameters, -q 0.01). The peak signals were
413 aggregated following the procedure described in Zhou et al.¹². In brief, the genome was split into 200
414 bp bins. All peak calls were intersected with these bins. If a bin overlaps a peak call to at least 50 %
415 (100 bp), the bin was labelled as belonging to that dataset class. All genomic bins that do not intersect
416 with at least one peak call were discarded. Only autosomes were used for all analysis.

417 Mouse chromatin data were retrieved from ENCODE¹⁸. Histone modification peak calls were
418 downloaded from the ENCODE data portal. For DNase-seq, ATAC-seq and transcription factor
419 ChIP-seq data the aligned bam files were downloaded and peak called with macs2⁴² as described
420 above. Replicates were collapsed into unions. All mouse data used are listed in Supplementary Table
421 1.

422

423 **Hi-C data.** Publicly available, deeply sequenced Hi-C data from Rao et al.²⁰ was used. The available
424 5 and 10 kb resolution intra chromosomal contacts maps of 7 cell lines (and 1kb data from GM12878)
425 were downloaded and normalized using the provided KRnorm factors. Four replicates of mouse ES
426 cell data²¹ were retrieved as raw fastqs from (GSE96107).

427

428 **Hi-C encoding for deep learning.** The genome was divided into bins matching the bin size of the
429 respective Hi-C data resolution used for training (1 kb, 5 kb, 10kb). For every stretch of DNA of size

430 1 Mb + bin size bps (e.g. 1005000 for 5kb bins), the chromatin interactions associated with the
431 window were assigned as squares in a vertical, zig-zag pole over the centre of the sequence window
432 (see Fig. 1b). Every square encodes the Hi-C interactions observed between two bin sized windows
433 of increasing distance (up to 1 Mb away). By sliding the large DNA stretch over a chromosome with
434 a bin sized increment, this encoding recovers the chromosome wide Hi-C map up to an interaction
435 distance of 1 Mb. Regions with a median interaction count along this pole of 0 were excluded from
436 training. The Hi-C data was percentile normalized across individual chromosomes for every
437 interaction distance in bin sizes. It is of particular interest to resolve high levels of Hi-C interactions
438 at high resolution and only a low percentage of chromatin interactions is expected to yield strong
439 interactions at larger distances for example the corners of TAD triangular structures. Thus, the
440 percentile normalization was designed to better resolve these high interaction levels at larger distances
441 by using uneven percentiles in a pyramid like scheme (from low to high: 2 x 20 %, 4 x 10 %, 4 x 5
442 %, see Extended Data Figure 1). The identifier of the respective pyramid percentile (1 – 10) was
443 stored. The chromatin interaction network was then trained to predict the percentile identifier as a
444 regression problem (see below).

445

446 **Deep neural network architectures and training.** A two-step training process with transfer learning
447 (Extended Data Figure 2) was used. First a convolutional neural network was trained to predict
448 chromatin features from 1 kb of DNA sequence, using the compendium of 936 datasets described
449 above. The principle network architecture was adapted from DeepSEA¹². Five convolutional layers
450 (hidden units: 300, 600, 600, 900, 900; filter widths: 8, 8, 8, 4, 4, 4) with ReLU activation, max
451 pooling (widths: 4, 5, 5, 5, 2) and dropout (rate: 0.2) were used followed by a fully connected layer
452 with sigmoid activation to output individual probabilities for each chromatin feature class (multi-
453 label classification). The network parameters were trained by minimizing the sum of the binary cross

454 entropies using the ADAM optimizer (epsilon 0.1) in batches of 100. Batch size, dropout rate,
455 learning rate and filter size were optimized by grid search.

456 Second, a chromatin interaction network was trained to predict Hi-C interaction from DNA sequence.
457 The chromatin interaction network takes as input 1 Mb + 1x Hi-C bin size [bp], (e.g. 1005000 for a
458 5 kb bin network). The first module consists of five convolutional layers, with ReLU, max pooling
459 and dropout with the dimensions and hyperparameters matching the chromatin feature network. The
460 hidden weights were initialized by seeding with the weights of the trained chromatin feature network
461 from step one. All chromatin features were used for pre-training and the same weights were used for
462 seeding the chromatin interaction network training independent of the Hi-C data cell type.

463 The second module is a series of ten dilated, gated 1D convolutional layers with residuals¹⁶. Gated
464 convolutional layers require training double the amount of filter parameters but have the potential of
465 modelling more complex functions through their multiplicative units. The residual units allow
466 information to propagate more easily through the network without having to necessarily pass through
467 convolutions⁴³. 100 hidden units were used, and dilation rates were increased exponentially to reach
468 the full sequence context in the last layer (1, 2, 4, 8, 16, 32, 64, 128, 256, 1). The dilated layers were
469 followed by a fully connected layer. Output are the predicted interaction strengths (in units matching
470 the percentile normalization). The model was trained with ADAM (epsilon 0.1) to minimize the sum
471 of squares error between the outputs and the true percentiles. GPU memory limited us to using a batch
472 size of 1. Hidden units (for dilated layers), dropout rate, learning and ADAM epsilon were optimized
473 using grid search.

474

475 **Network training, computational resources and limitations.** For both training procedures the data
476 were split into training, validation and test set based on chromosomes. For the chromatin feature
477 network chr11 and 12 were used for validation and chr15, 16 and 17 for testing. For the chromatin

478 interaction network, to increase the number of training examples the same validation chromosomes
479 were used but only chr16 and 17 were used as test chromosomes.

480 All models were trained on NVIDIA Titan V cards with 12 GB of video memory. Training on smaller
481 cards is possible but slower. The final models have ~ 60 M parameters. Scaling the models to larger
482 DNA inputs will likely benefit from network pruning or a refined architecture.

483 Fully training the chromatin feature network required 14 epochs with about 8 hours per epoch. The
484 training set order was reshuffled after every epoch. To minimize the amount of times large chunks of
485 DNA sequence had to be loaded into memory the network was trained on one chromosome at a time.
486 Within a chromosome, the order in which training batches were drawn was random. Interestingly, we
487 observed that the chromatin interaction network, when seeded with the pre-trained weights in the first
488 convolutional filters, converged quickly, after training on $\sim 3 - 6$ chromosomes and only marginally
489 improved after training for an entire epoch or longer. Models were trained for one full epoch as we
490 have not observed significant improvement after training for longer and the limited batch size as well
491 as the network complexity make training slow. For cross-validation, we trained multiple iterations
492 holding out different chromosomes from training.

493 While training networks is only feasible with GPU support, predictions with trained models can be
494 run on CPU only. For example, predicting the impact of a variant requires ~ 5 min with GPU and \sim
495 2h with only CPU support.

496

497 **Predicting changes in chromatin interactions.** For calculating differences in chromatin
498 interactions, the interactions over the reference sequence were predicted for every position that is
499 within 1 Mb (plus 1x Hi-C bin size) of the sequence variant. This matches the respective models
500 spatial reach. The reference sequence was then modified to match the sequence variant of interest.
501 After predicting the chromatin interactions over the variant sequence, the difference was quantified

502 by calculating the absolute difference between reference and variant prediction at each interaction bin
503 and summarized as the mean absolute difference over all covered interactions.

504

505 **Distance-stratified correlation.** The Pearson correlation coefficient was calculated between the Hi-
506 C skeleton and the deepC predicted regression score. Note that the skeleton percentiles are discrete
507 (percentile tag 1 – 10), while the regression score is continuous. The Hi-C skeleton is noisy even at
508 very deep sequencing depths (e.g. GM12878). Therefore, a small mean filter was employed using a
509 5x5 window to smooth the skeleton and the distance-stratified correlation was calculated between the
510 prediction and the raw or the smoothed skeleton respectively.

511

512 **Comparison against HiC-Reg⁸.** The available CrossChrom predictions, trained on chr14 and
513 predicted on chr17, were downloaded from the supplementary material. HiC-Reg predictions were
514 distance-normalized as described above.

515

516 **Re-aligning and downsampling Hi-C data.** The primary GM12878 replicate was realigned from
517 raw fastqs using HiCPro⁴⁴. The valid Hi-C interactions were downsampled to achieve 1 billion, 100
518 million and 10 million valid interactions, respectively. Mouse ES cell Hi-C data were aligned and
519 processed from raw fastq data using HiCPro⁴⁴.

520

521 **Selection of validation capture probes.** A total of 220 viewpoints were selected for validating the
522 deepC predictions, specifically selecting genomic locations where the Hi-C data and deepC
523 predictions differed in detail or where the deepC predicted structures were only very faintly noticeable
524 in the Hi-C data. Two sets were designed, one targeting 81 CTCF sites and one targeting 139 intra
525 domain viewpoints that lie within a distinct Hi-C/deepC domain but are not intersecting with any
526 potential functional elements. Capture probes were designed using CapSequim

527 (<http://apps.molbiol.ox.ac.uk/CaptureC/cgi-bin/CapSequim.cgi>), filtering out repetitive probe regions
528 as described in the online documentation. For the final probe design see Supplementary Table 2.

529

530 **Cell culture and fixation.** Human GM12878 lymphocyte cell line, were obtained from the NIGMS
531 Human Genetic Cell Repository at the Coriell Institute for Medical Research and cultured in RPMI
532 1640 supplemented with 15% FBS, 2mM L-Glutamine and 100U/ml Pen-Strep at 37 °C in a 5% CO₂
533 incubator. K562 cells were supplied by the WIMM transgenics facility. Cells were maintained in
534 RPMI 1640 media supplemented with 10 % FCS at 37 °C in a 5% CO₂ incubator. Both cell types
535 were fixed and processed using the same protocol. Cells were resuspended at 1x10⁶ cells per ml and
536 fixed at room temperature with 2% v/v formaldehyde for 10 minutes. Fixation was quenched with
537 120 mM glycine. Cells were washed with ice cold PBS. Cells resuspended in cold lysis buffer (10
538 mM Tris, 10 mM NaCl, 0.2% Igepal CA-630 and complete proteinase inhibitor (Roche)) and snap
539 frozen to -80 °C. See Life Science Reporting Summary for additional details.

540

541 **3C library preparation.** 3C libraries were prepared as described previously²⁴ with the following
542 modifications: Centrifugation's were performed at 500 rcf, thermomixer incubations were set to 500
543 rpm, and following ligation chromatin was pelleted by centrifugation (15 min, 4°C, 500 rcf) and the
544 supernatant discarded. To increase sequencing depth and minimize PCR duplicates experiments were
545 performed in technical triplicates and four unique adapters were used per replicate.

546

547 **NG Capture-C.** Double capture was performed as described previously²⁴ with biotinylated
548 oligonucleotides (IDT xGen Lockdown Probes) in two pools (Supplementary Table 2) with 3 pg of
549 each oligonucleotide per 3C library. The generated NG Capture-C libraries were sequenced using
550 Illumina sequencing platforms (V2 chemistry; 150-bp paired-end reads) and data collected using the
551 NextSeq System Suite (v2). To resolve even subtle changes in chromatin interaction domains at high

552 resolution, the libraries were deeply sequenced (GM12878 CTCF - 128 M; GM12878 intra domain -
553 118 M reads; K562 CTCF – 302 M; K562 intra domain – 289 M reads). All technical replicates were
554 merged for the analysis.

555

556 **Tiled-C.** Tiled-C generates Hi-C like data focused on loci of interest at greater depth by using an
557 oligonucleotide capture enriching a 3C library for viewpoints tiled over the regions of interest. Tiled
558 oligonucleotides were designed using the design approach and tool described in Oudelaar and Beagrie
559 *et al.*²⁹ (<https://oligo.readthedocs.io/en/latest/>). A panel of double-stranded capture oligonucleotides
560 from Twist Bioscience (Custom probes for NGS target enrichment) was used. The Tiled-C procedure
561 was performed as described in Oudelaar and Beagrie *et al.* using the 3C libraries from GM12878 and
562 K562 cells. All biological and technical replicates were merged for the analysis. For a summary of
563 the regions of interest and designed probes see Supplementary Table 3.

564

565 **NG Capture-C analysis.** NG Capture-C data were mapped, quality controlled and visualized using
566 CcseqBasicS⁴⁵ following the procedure described previously²⁴.

567

568 **Tiled-C analysis.** Tiled-C data were mapped and quality controlled using the tiled mode of
569 CcseqBasicS described above. Each region was ICE normalized using the ICE implementation of
570 HiCPro with default parameters.

571

572 **Distance-normalize NG Capture-C tracks.** To compare them to the Hi-C skeleton, NG Capture-C
573 tracks were normalized for distance dependence per viewpoint. The number of interactions with each
574 restriction enzyme fragment was extracted and the distance to the viewpoint recorded. The
575 interactions were then normalized for the total number of *cis* interactions for the respective viewpoint.
576 Pooling this information across all viewpoints we observed that the distance decay approximately

577 follows a log – log linear trend when we split the data into three distance bins (close, intermediate
578 and far). The distance thresholds for these bins were empirically optimized for every NG Capture-C
579 set (see Supplementary Table 4), excluding all interactions closer than 2.5 kb to the respective
580 viewpoint. The distance decay is then approximated by three linear regression fits, one for each
581 distance bin. The distance-normalized interactions were calculated per viewpoint by dividing the
582 observed *cis* normalized interactions with the expected interactions from the linear fit at the respective
583 distance.

584

585 **Insulation score boundary calling.** Interaction domain boundaries were called using the Hi-C data,
586 the Hi-C skeleton and the deepC predicted interactions using an insulation score-based approach that
587 was adapted from Crane et al.²⁶ . Using the 5 kb bin sized data, the mean insulation score profile was
588 calculated based on a 25 kb window to allow for a more intricate boundary call. The first derivative,
589 or delta vector, of the insulation score profile was approximated using a 1D Sobel operator. Zero
590 crossings in this delta vector represent local minima and maxima of the insulation score. Maxima
591 were discarded. The remaining boundaries were further filtered by calculating the approximation of
592 the second derivative of the insulation profile using the same procedure described above. The height
593 of this delta2 vector reflects the change in delta, with sharper boundaries having a higher delta2 score.
594 Boundaries with a delta2 score smaller than 0.1 were removed and the remaining boundaries were
595 stratified based on their delta2 score.

596

597 **TopDom.** TopDom was retrieved from the gitHub implementation
598 (<https://github.com/HenrikBengtsson/TopDom>). Boundaries were called using the window
599 parameters 5, 10 and 20. Boundaries between all types of called domains were used.

600

601 **Distance-normalized NG Capture-C signal over boundaries.** The mean, distance-normalized NG
602 Capture-C signals over boundaries were calculated. In NG Capture-C tracks from single viewpoints,
603 domain boundaries can be subtle and get harder to detect the further away from the viewpoint they
604 are located. Therefore, boundaries further then 1 Mb away from a viewpoint were excluded. The
605 mean normalized Capture-C signal over boundaries relative to their centre was calculated.

606
607 **Virtual4C from Hi-C skeleton and deepC maps.** By extracting all interacting windows with a
608 viewpoint of interest, Hi-C data can be transformed into virtual 4C profiles. For this work, virtual 4C
609 profiles were derived from the Hi-C skeleton and the deepC predictions yielding distance- normalized
610 profiles. Virtual 4C profiles from the Hi-C skeleton and deepC predictions were compared to
611 distance-normalized NG Capture-C tracks by calculating the respective Pearson correlation of all
612 interactions within 1 Mb from a given viewpoint. Because the skeleton percentiles are discrete and
613 punctuate, a running mean smoothing window of 25 kb was applied. In contrast, the deepC
614 predictions are smooth and therefore no additional smoothing was applied.

615
616 **Chromatin segmentation.** GM12878 and K562 chromatin data were downloaded from the
617 ENCODE data portal (see Supplementary Table 5). Filtered alignments to hg19 were downloaded
618 and replicates were merged. Peaks were called using macs2⁴² with default settings and -q 0.01.
619 Deeptools⁴⁶ was used to create bigwig coverage tracks. DNase-seq and CTCF ChIP-seq peaks were
620 merged to a union set merging peaks within 10 bp of each other using bedtools⁴⁷ (bedtools merge -d
621 10). Union peaks were formatted to 600 bp elements centred on the peaks. Deeptools was then used
622 to extract the read coverage for each chromatin dataset over each peak union element. For this,
623 elements were extended to 1000 bps to better capture flanking histone modifications. Using the
624 derived count matrix, chromatin classes were segmented using GenoSTAN⁴⁸ running on the elements
625 rather than entire chromosome stretches. The HMM model was trained using the Poisson log-normal

626 distributions. Twelve classes were fitted and merged into eleven classes based on similarity of the
627 chromatin signatures. The classes were manually curated and classified into promoter, enhancer and
628 CTCF sites with varying activity levels based on H3K27ac coverage.

629

630 **Saliency score.** Adapted from image analysis³⁰ the saliency score serves as a proxy for the importance
631 of every base pair to the interaction predictions. Explicitly, the saliency score was calculated as the
632 dot product of the gradient of the model output with respect to the sequence input and the one-hot
633 encoded DNA sequence input. This effectively masks the impact of non-present bases. For a given
634 window the saliency score relates to the interaction pole on the center. To visualize saliency tracks,
635 the sequence window was moved in bin sized steps and the saliency per base pair was averaged over
636 all sequence windows (sized 1 Mb + bin size) that include the respective base pair. To simplify
637 visualization and interpretation the absolute value of the saliency score was used. Metaplots were
638 computed with deepTools.

639

640 **eQTL data analysis.** EBV transformed lymphocyte specific eQTLs were retrieved from GTEx (v7
641 accessed from the GTEx portal 01/03/2019). A union of DNase-seq and CTCF ChIP-seq peaks was
642 created using bedtools merge. The eQTL SNPs were filtered for intersection with the union of
643 GM12878 open chromatin and CTCF peaks. Indels were removed. A background SNP set was
644 constructed by shuffling the eQTL SNPs on the respective same chromosome and forcing them to
645 stem from within the union peaks (bedtools shuffle -chrom -incl). Absolute saliency scores of the
646 SNP bases derived from the 5 kb resolution GM12878 model were extracted. Empirical cumulative
647 distributions were derived and tested for significance using a two sample Kolmogorov-Smirnov test
648 (R, ks.test, reshuffled SNP saliency (n = 6607) vs. eQTL set (n = 6607) saliency, alternative
649 hypothesis: “less”).

650

651 **Deletion screen.** Separately, GM12878 DNase-seq and CTCF ChIP-seq peaks were merged if
652 multiple peaks were found within 1.5 kb of each other (bedtools merge -d 1500). Peaks were extended
653 to at least 300 bp. All DNase peaks that overlapped with CTCF peaks were removed. For every
654 remaining CTCF (n=45635) and DNase (n=47320) site the impact on chromatin interactions upon
655 deleting the respective site was predicted the 5 kb GM12878 model. Chromatin classes were assigned
656 based on overlap with the GenoSTAN chromatin segmentation described above. In addition, n = 3850
657 background sites were selected by shuffling all CTCF and open chromatin sites on chr16 forcing no
658 overlap (bedtools shuffle -chrom -noOverlapping).

659

660 **5C data.** Processed 5C data from Hnisz et al.³¹ were downloaded, binned into 5kb bins and
661 visualized.

662

663 **SNP sampling.** 1000 random SNPs each were sampled from strong CTCF sites, strong promoters
664 and strong enhancers as classified by the chromatin segmentation procedure described above. For a
665 background set, 1000 SNPs were sampled from the 400 bp regions flanking these regulatory elements
666 while avoiding any overlap with other elements. SNPs positions were sampled using bedtools shuffle
667 and variant bases were randomly selected from the three bases not present in hg19 at the respective
668 position.

669

670 **Statistics and Replication.** Statistical analysis was performed in R. Statistical tests are described in
671 the relevant subsection of the Online Methods. NG Capture-C and Tiled-C experiments were
672 performed once, technical replicates were pooled for maximum read depth (see Life Science
673 Reporting Summary for additional details).

674

675 **Additional software and packages.** All neural networks were implemented in python (v3.5) and
676 tensorflow⁴⁹ (developed under 1.8.0).

677

678 *Additional Tools*

- 679 • samtools⁵⁰ (v1.3)
- 680 • FastQC (v0.11.4) (<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>)
- 681 • Bowtie⁵¹ (v1.1.2)

682 *Additional R packages*

- 683 • cowplot (v0.6.2, <https://github.com/wilkelab/cowplot>)
- 684 • GenomicRanges⁵² – (v1.30.3)
- 685 • ggplot2⁵³ (v3.1.0)
- 686 • RcolorBrewer (v1.1.1-2,
687 <https://cran.r-project.org/web/packages/RColorBrewer/index.html>)
- 688 • rtracklayer⁵⁴ (v1.30.4)
- 689 • tidyverse (v1.3.0) (<https://www.tidyverse.org>)
- 690 • zoo⁵⁵ (v1.8.1)

691 *Additional Python libraries*

- 692 • numpy⁵⁶ (1.16.4)
- 693 • h5py (v2.9.0, <http://www.h5py.org>)
- 694 • pysam (0.15.2, <https://github.com/pysam-developers/pysam>)

695

696 **Data availability**

697 Hi-C data from Rao et al. is available under GSE63525. Chromatin feature data from ENCODE,
698 Roadmap and other publicly available data are listed in detail with accession numbers in
699 Supplementary Table 1. Additional ENCODE data used for chromatin segmentation and visualization
700 are listed with accession numbers in Supplementary Table 5. Tiled-C and NG Capture-C validation
701 data are available under the GEO super series GSE137437.

702

703 **Code availability**

704 All code for training and employing deepC networks as well as trained models are available under:
705 <https://github.com/rschwess/deepC>; All code for training and employing chromatin feature networks
706 is available under: <https://github.com/rschwess/deepHaem>

707

708 **Methods-only References**

- 709 38. Corces, M. R. *et al.* Lineage-specific and single-cell chromatin accessibility charts human
710 hematopoiesis and leukemia evolution. *Nat. Genet.* **48**, 1193–1203 (2016).
- 711 39. Schwessinger, R. *et al.* Sasquatch: predicting the impact of regulatory SNPs on transcription
712 factor binding from cell- and tissue-specific DNase footprints. *Genome Res.* **27**, 1730–1742
713 (2017).
- 714 40. Downes, D. J. *et al.* An integrated platform to systematically identify causal variants and
715 genes for polygenic human traits. Preprint at:
716 <https://www.biorxiv.org/content/10.1101/813618v1> (2019).
- 717 41. Telenius, J., Consortium, T. W. & Hughes, J. R. NGseqBasic - a single-command UNIX tool
718 for ATAC-seq, DNaseI-seq, Cut-and-Run, and ChIP-seq data mapping, high-resolution
719 visualisation, and quality control. Preprint at:

720 <https://www.biorxiv.org/content/10.1101/393413v1> (2018).

721 42. Zhang, Y. *et al.* Model-based analysis of ChIP-Seq (MACS). *Genome Biol.* **9**, R137 (2008).

722 43. He, K., Zhang, X., Ren, S. & Sun, J. Deep Residual Learning for Image Recognition.

723 *Multimed. Tools Appl.* **77**, 10437–10453 (2015).

724 44. Servant, N. *et al.* HiC-Pro: an optimized and flexible pipeline for Hi-C data processing.

725 *Genome Biol.* **16**, 259 (2015).

726 45. Telenius, J. M. *et al.* CaptureCompendium: a comprehensive toolkit for 3C analysis. Preprint

727 at: <http://biorxiv.org/content/early/2020/02/18/2020.02.17.952572> (2020).

728 46. Ramírez, F., Dündar, F., Diehl, S., Grüning, B. A. & Manke, T. deepTools: a flexible

729 platform for exploring deep-sequencing data. *Nucleic Acids Res.* **42**, W187–W191 (2014).

730 47. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic

731 features. *Bioinformatics* **26**, 841–842 (2010).

732 48. Zacher, B. *et al.* Accurate Promoter and Enhancer Identification in 127 ENCODE and

733 Roadmap Epigenomics Cell Types and Tissues by GenoSTAN. *PLoS One* **12**, e0169249

734 (2017).

735 49. Abadi, M. *et al.* TensorFlow: A System for Large-Scale Machine Learning TensorFlow: A

736 system for large-scale machine learning. *12th USENIX Symp. Oper. Syst. Des. Implement.*

737 *(OSDI '16)* 265–284 (2016).

738 50. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–

739 2079 (2009).

740 51. Langmead, B., Trapnell, C., Pop, M. & Salzberg, S. L. Ultrafast and memory-efficient

741 alignment of short DNA sequences to the human genome. *Genome Biol.* **10**, R25 (2009).

742 52. Lawrence, M. *et al.* Software for Computing and Annotating Genomic Ranges. *PLoS*

743 *Comput. Biol.* **9**, e1003118 (2013).

744 53. Wickham, H. *ggplot2: Elegant Graphics for Data Analysis*. (Springer New York, 2009).

- 745 54. Lawrence, M., Gentleman, R. & Carey, V. rtracklayer: An R package for interfacing with
746 genome browsers. *Bioinformatics* **25**, 1841–1842 (2009).
- 747 55. Zeileis, A. & Grothendieck, G. Zoo: S3 infrastructure for regular and irregular time series. *J.*
748 *Stat. Softw.* **14**, 1–27 (2005).
- 749 56. van der Walt, S., Colbert, S. C. & Varoquaux, G. The NumPy Array: A Structure for
750 Efficient Numerical Computation. *Comput. Sci. Eng.* **13**, 22–30 (2011).
- 751

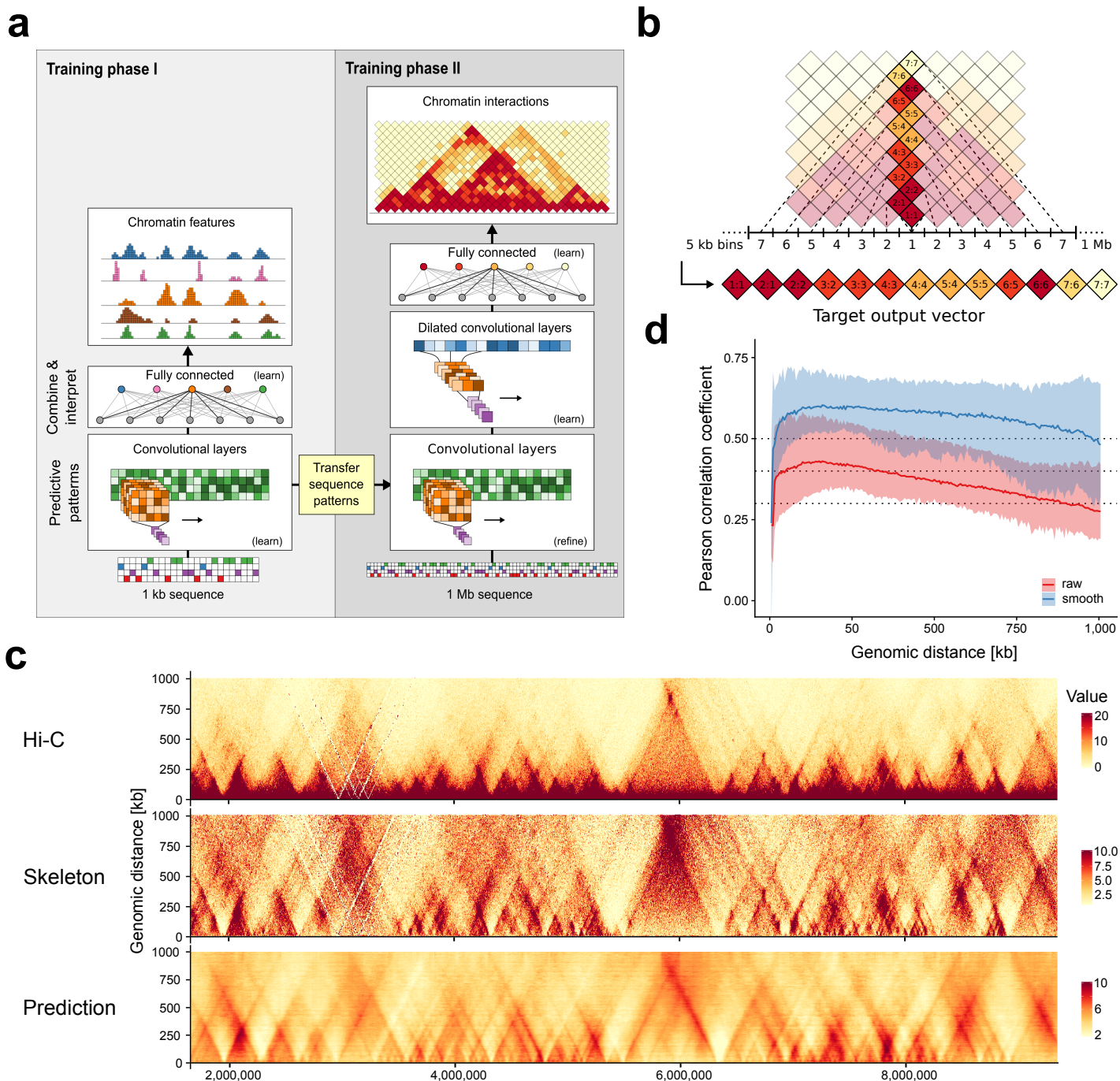


Figure 1 | Predicting Hi-C interactions from DNA sequence. **a)** Overview of the deepC architecture and training workflow. **b)** Encoding of Hi-C data as target vector for prediction given a 1 Mb window of DNA sequence. **c)** Comparison of Hi-C data, the derived Hi-C skeleton and the interactions predicted from DNA sequence using deepC. Shown is a ~ 7 Mb region on hold out chromosome 17. **d)** Distance-stratified Pearson correlation between the Hi-C skeleton and the deepC predictions in a cross-validation scheme across all chromosomes. Solid lines indicate the mean correlation value and the area indicates the space between the maximum and the minimum values over all chromosomes. Red shows the correlation with the raw and blue with the (5x5) mean filter smoothed skeleton values. Dotted lines at 0.3, 0.4 and 0.5.

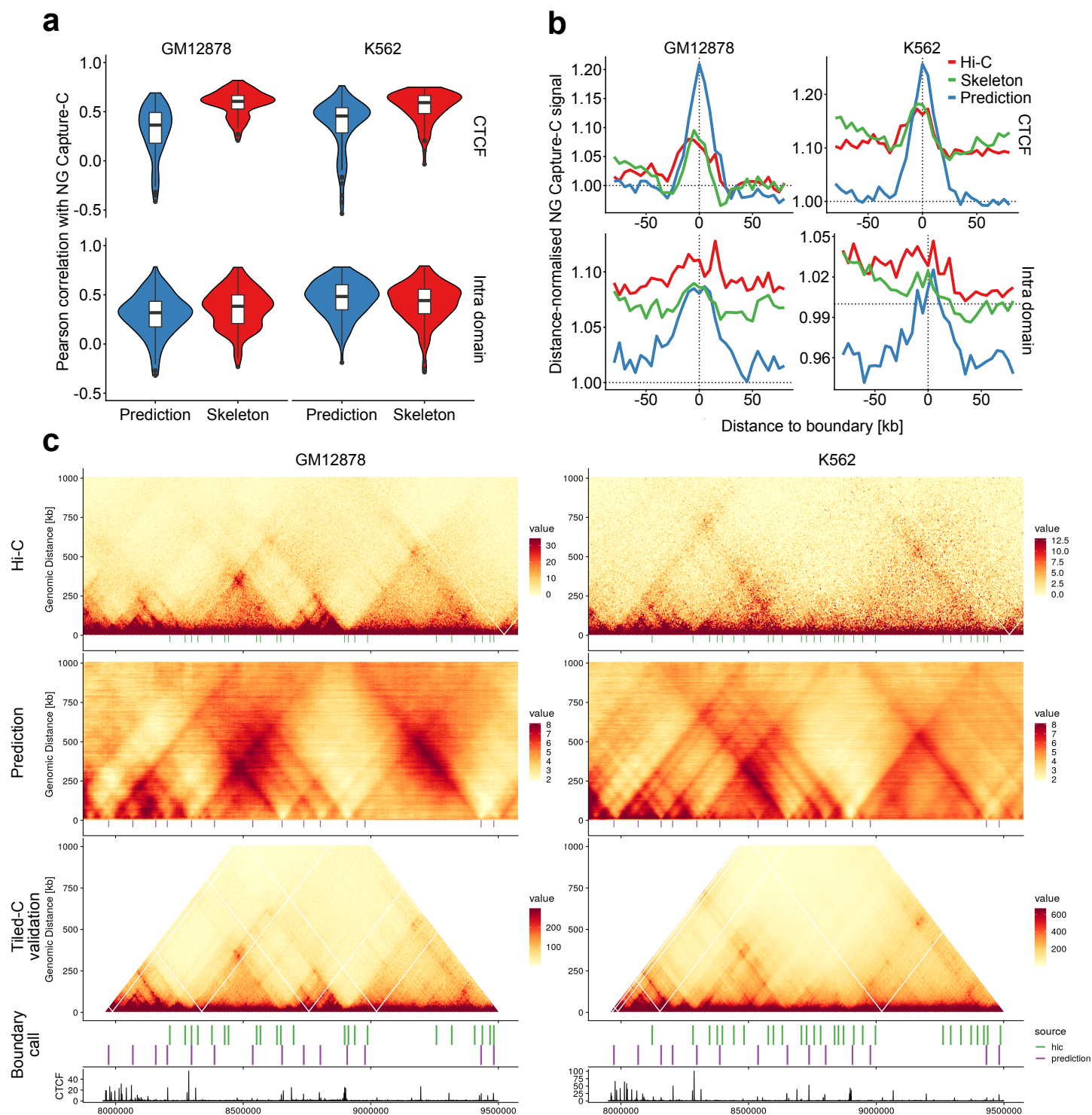
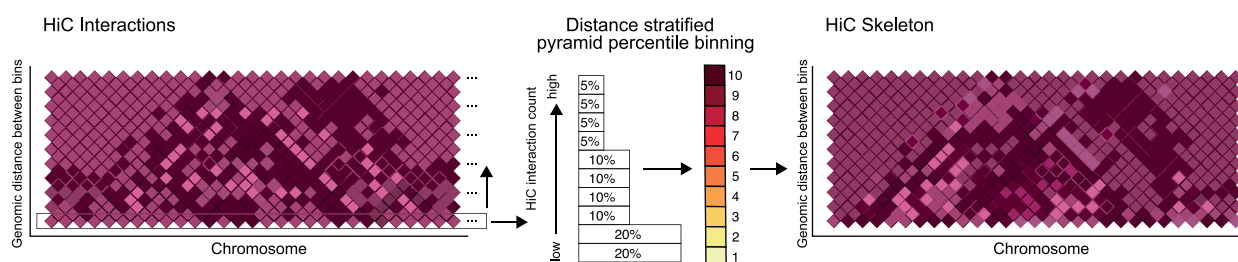
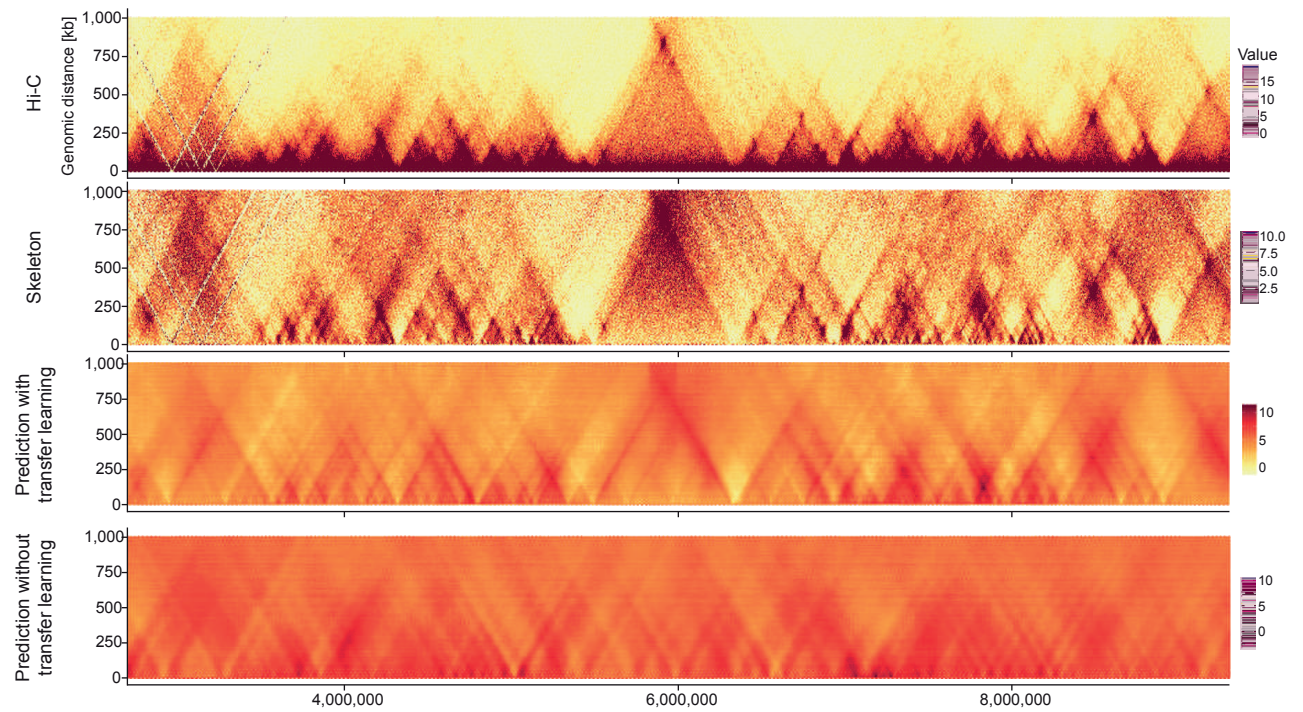


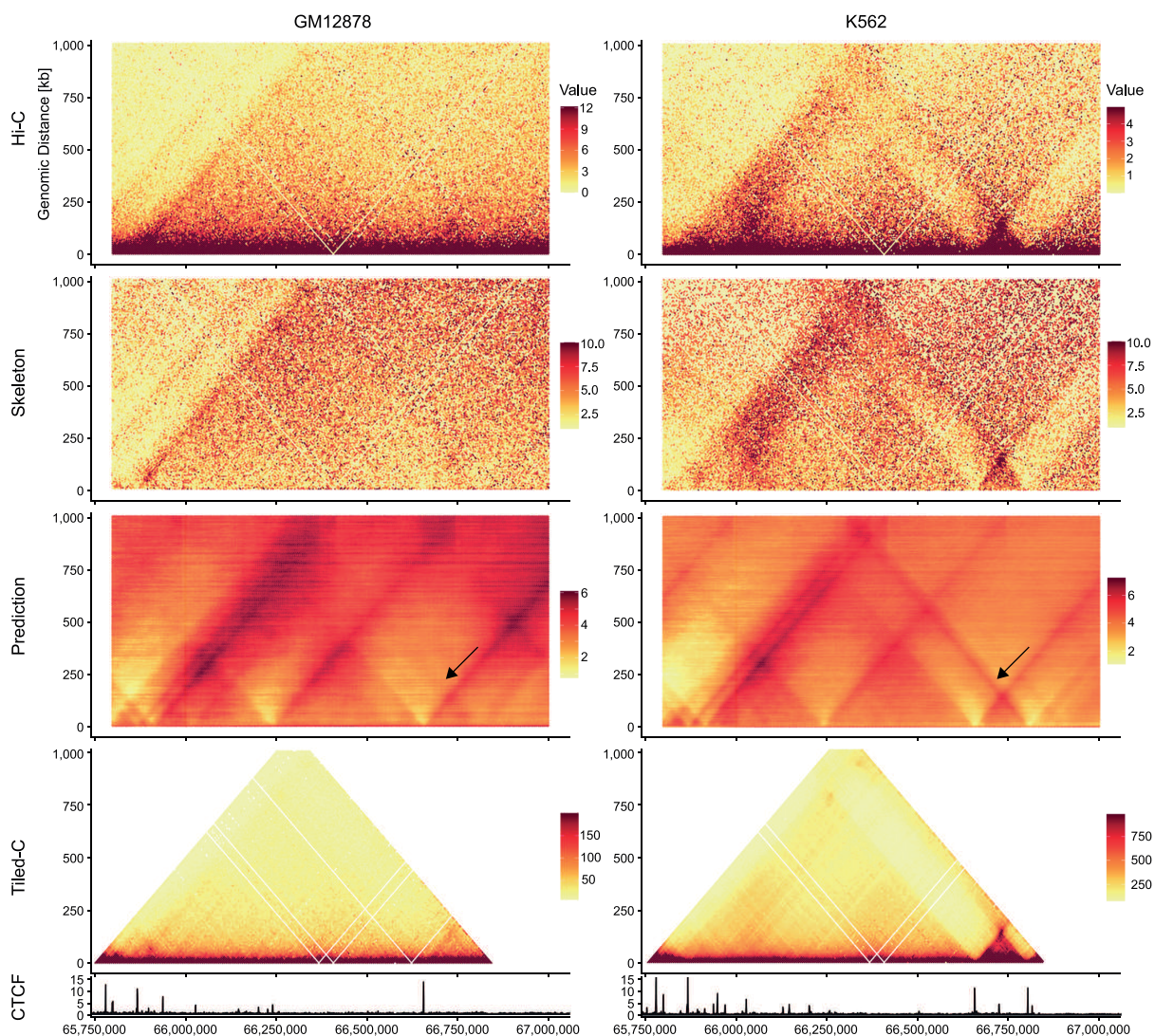
Figure 2 | Validation of deepC predictions. **a)** Comparing the correlation between the validation NG Capture-C profiles and the virtual 4C profiles derived from the Hi-C skeleton (red); and the deepC prediction map (blue) from all viewpoints in two cell types. Boxplots: median middle thick line, 25th and 75th percentile left and right hinge respectively, whiskers stretch up 1.5 times the IQR (inter quartile range). **b)** Meta-profiles of the average NG Capture-C signal over domain boundaries called at high resolution from the Hi-C data, the skeleton and the deepC predicted interaction map respectively. Shown is the mean distance-normalised NG Capture-C signal relative to the boundary center. The labels "CTCF" and "Intra domain" refer to the NG Capture-C viewpoint fragments. These were designed to overlap either CTCF sites or to lie within insulated domains but not overlap with regulatory genomic elements, so as to capture the domain structure and not the interactions of specific genomic elements. **c)** Shown are Hi-C data, the deepC predicted interaction map and the Tiled-C high sensitivity map over a locus on chr17, a hold out chromosome. Boundaries called at high resolution from Hi-C (green) and deepC predictions (purple) are aligned under the respective map and the Tiled-C map. Cell-type-specific CTCF ChIP-seq tracks are visualized below. For contrast, Hi-C and Tiled-C data were bounded between 5 and 95 % of coverage and deepC predictions were bounded between 2 and 8 predicted regression score.



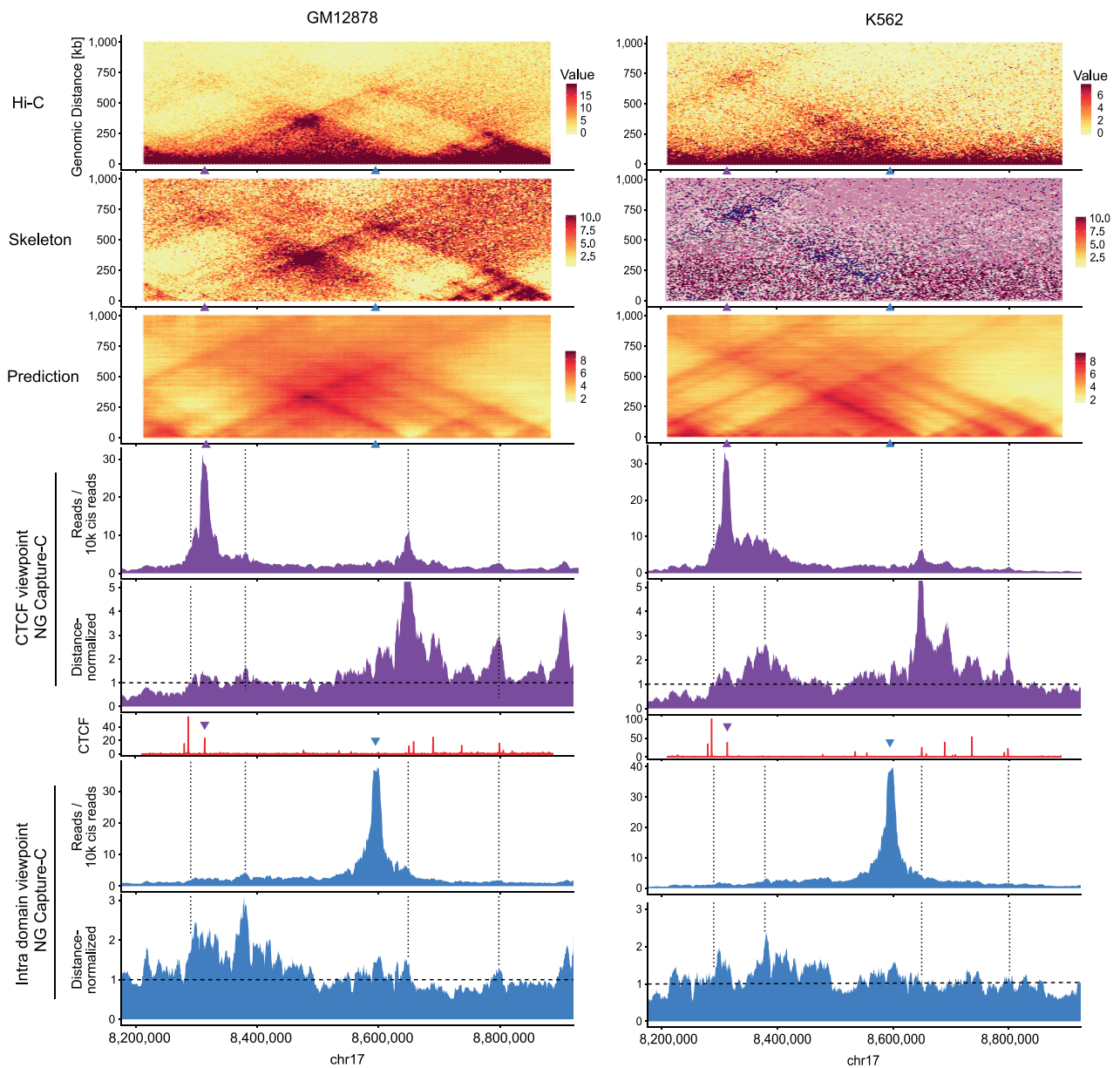
Extended Data Figure 1 | Percentile normalizing Hi-C data for deep learning. The Hi-C interactions are percentile-binned in a distance-stratified manner. For every genomic distance, in steps equal to the bin size, the Hi-C signal is split into unequal percentiles ranging from 20 % bottom to 5 % top. The percentiles are attributed the values 1 to 10 yielding the Hi-C skeleton. The unequal percentile sizes ensure a finer distinction of the differences at the high Hi-C interaction value range, while minor differences in the low interaction value range are squished. Effectively, this procedure reduces the proximity signal and enhances domains and domain boundaries.



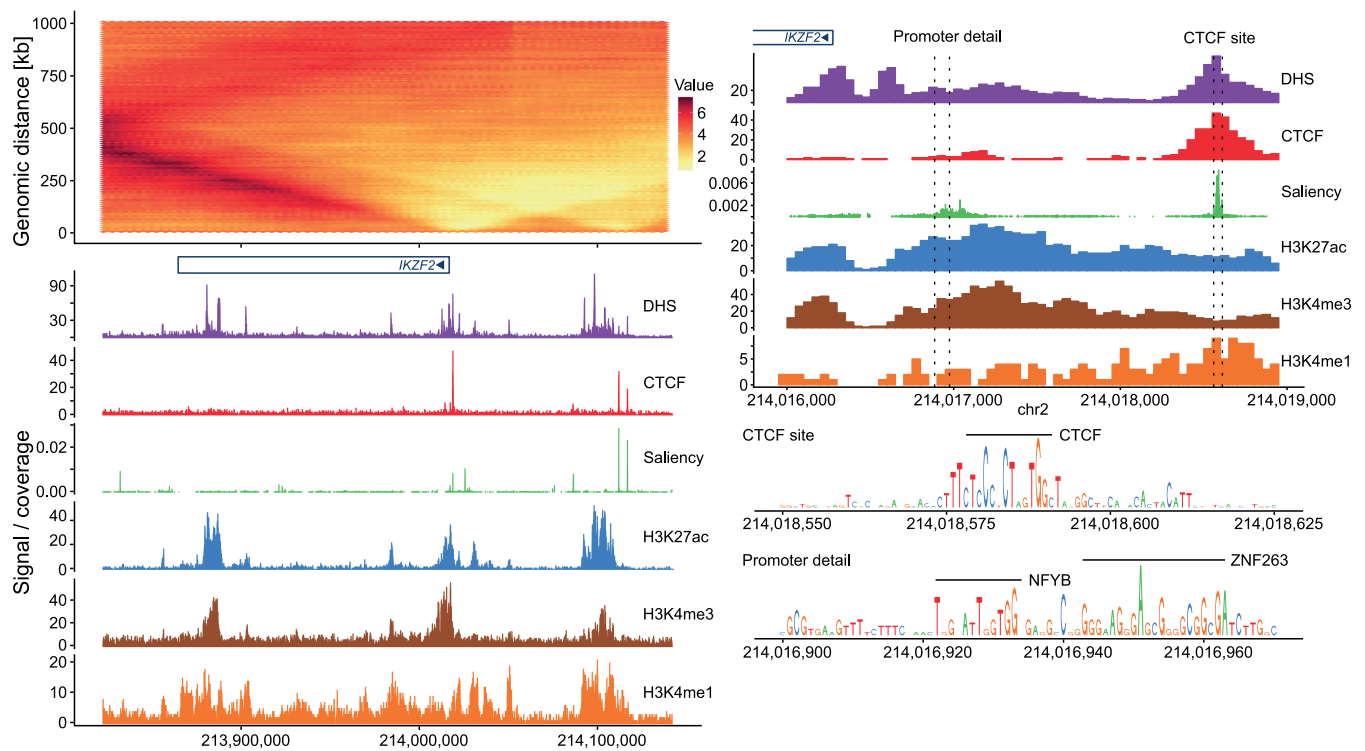
Extended Data Figure 2 | Comparison of deepC training with and without transfer learning. Training a deepC model with the same architecture but without pre-seeding the lower convolutional layers with the chromatin feature model weights results in the emergence of triangular structures. Their positioning however does not match with the Hi-C structures. In contrast, with pre-seeding the predicted domains overlap well with the Hi-C skeleton.



Extended Data Figure 3 | Tissue-specific deepC predictions. Shown is a region on chromosome 2 around the MEIS1 locus. DeepC predicts a small domain with insulation to the upstream regions (black arrow) in a tissue specific manner. The domain is only visible in K562 Hi-C data and matches with tissue-specific CTCF binding. Tiled-C confirms the tissue-specific domain. For contrast, Tiled-C data were bounded between the 5 and 95 percentiles.



Extended Data Figure 4 | NG Capture-C validation of deepC predictions. a) Example region with overlap of GM12878: Hi-C, skeleton and deepC prediction; NG Capture-C tracks, distance-normalized NG Capture-C tracks and CTCF ChIP-seq track (red). Shown is a CTCF viewpoint (purple triangle) and an intra domain viewpoint (blue triangle) not overlapping with any active elements. Dashed lines in the distance-normalized NG Capture-C tracks indicate the expected interaction value. Dotted black lines highlight deepC prediction details that correspond to boundaries in the NG Capture-C tracks. b) K562 data of the same region.



Extended Data Figure 5 | Mapping important features for genome folding. Shown are GM12878 deepC predictions over the IKZF2 locus (a) on chromosome 2 and focused on the IKZF2 promoter (b). Aligned are DHS as well as ChIP-seq tracks for CTCF and histone modifications. Shown in green is the saliency score which is a proxy for the importance every base has in predicting the chromatin interactions of that region. The saliency score shows sharp peaks overlapping CTCF binding sites and broader peaks overlapping active gene promoters. Resolving the saliency score at base pair resolution (b) highlights CTCF and general transcription factor binding motifs.