

Predicting Hi-C contact matrices through coarse-grained simulations of the chromatin

Maurizio Gilioli

Contents

1	ABSTRACT	2
2	INTRODUCTION	3
2.1	Chromatin as the information center of a cell	3
2.2	ATAC-sequencing and CTCF sequencing	3
2.3	<i>ChromHMM</i> allows the sequential characterization of chromatin states	4
2.4	Hi-C matrices and their patterns	5
2.5	Chromatin Coarse-Graining, chromatin as a polymeric fluid	6
3	Aim of the project	8
4	METHODS	9
4.1	Data used during the project	9
4.2	The Model	9
4.3	Finding enriched states in 5 kbs long beads	11
4.4	Trajectories analysis	11
4.5	Algorithms used for comparison	12
4.6	the Stratum Adjusted Correlation Coefficient (SCC) metric	13
5	RESULTS AND DISCUSSION	14
5.1	<i>ChromHMM</i> results	14
5.2	Results obtained while defining the models	16
5.3	Trajectory analysis results	17
5.4	Results from model selection	20
6	CONCLUSIONS	23
7	Glossary	24
	References	24

1 ABSTRACT

The chromatin is one of the most important parts of a cell. In fact, it contains in its volume the largest part of the cell DNA, and a great number of proteins, such as the histones, which help in the functional compaction of the nuclear DNA. However, the direct study of this substance encounters significant difficulties, and the analysis of related data do not give straightforward results. All-atomistic simulation approaches to predict the conformations of the chromatin in time are completely unfeasible, due to the large amount of atoms to simulate. Because of that, it is necessary to adopt a justified coarse-grained approach, which allows for simpler and less complex simulations. To this scope, I had the great pleasure of working in collaboration and on the code written by professor Marco Di Stefano. The aim of the project (which is still on-going) is to predict Hi-C matrices of contact by using coarse-grained simulations for stretches of DNA 2,000,000 bp long. Those maps are generally particularly hard to obtain and of high cost, however, the information that they contain can unveil very interesting mechanisms, such as the promoter-enhancer interactions. At first, a decompaction and a relaxation trajectories were made for 100 replicates. Then, the matrices were produced by tuning the parameters used as weights for the potentials with an iterative-stepwise approach. The results were confronted during the procedure and after with experimentally obtained maps by computing the SCC metric. Despite the fact that it has still not been reached a final and complete result, I believe that this report and its presentation could help significantly in improving the procedure and the produced analysis. To expand this study, we are already thinking about using the method for other *loci* and cell-types, and refining some of the used processes.

2 INTRODUCTION

2.1 Chromatin as the information center of a cell

All the living organisms have, inside the nucleus, the largest portion of their DNA, which is the main molecule through which information is passed from the old generations to the daughter cells. Due to the extreme length of the chromosomes, an assembly of DNA, proteins and RNA, called chromatin, is built in a necessary ordered and functional manner¹. The most important proteins used to reach this scope are the histones, towards which DNA is wrapped around, forming the nucleosomes. To govern the functioning of the DNA, the histones and the nucleic acid itself are subjected to a variety of modifications, such as methylation. The latter alteration, in mammals, occurs in specific sites of the genome, called CpGs, where a cytosine is connected directly to a guanine. Methylations of regulatory elements have been implicated in determining cell identity and chromatin structure^{2,3}. Among the proteins that interact with the DNA, CTCF is one of the most important to be mentioned. It is a protein conserved in eukaryotes and is ubiquitous in mammals⁴, and contains a Zinc-finger that binds to the DNA. The act of binding is performed in cooperation with cohesins, and causes the folding of the chromatin.^{4,5}

When it comes to its structure, the DNA is thought to fold in a hierarchical manner, as depicted in figure 1. However, this hypothesis makes a simplification: the fibers could have a range of diameters which depend on their activity and their location.^{6,7} Importantly, the diameter of the nucleosomes was determined to be, on average, equal to 10 nm. This type of measure was taken into account when building the Fine Scale (FS) model (chapter 4.2).

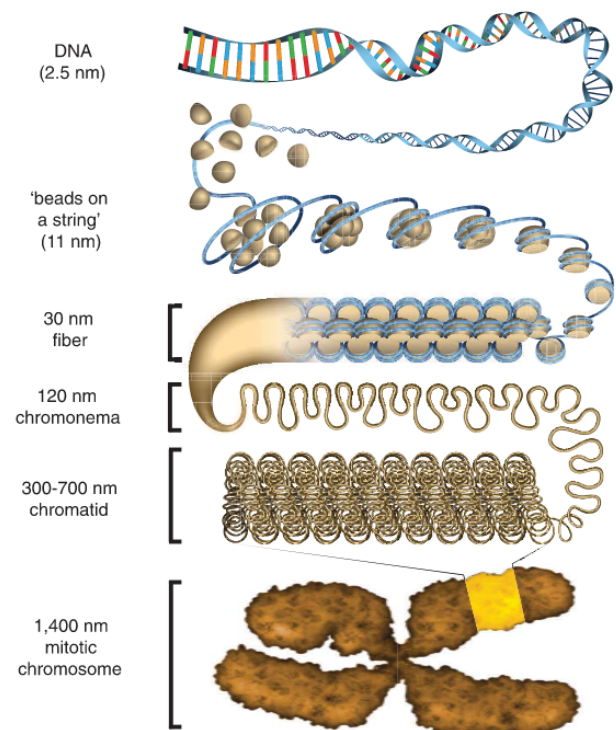


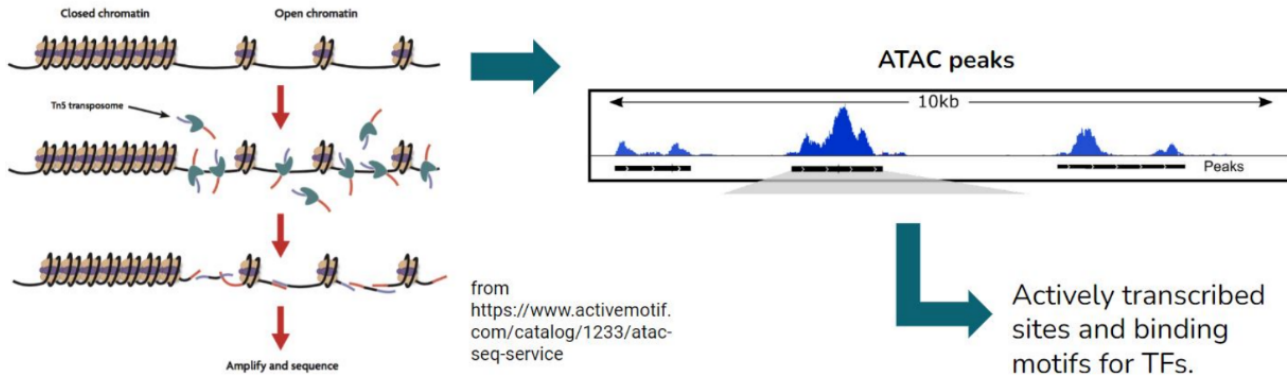
Figure 1. Image representing the hierarchical compaction of DNA in subsequently more compact and dense fibers.

2.2 ATAC-seq and CTCF sequencing

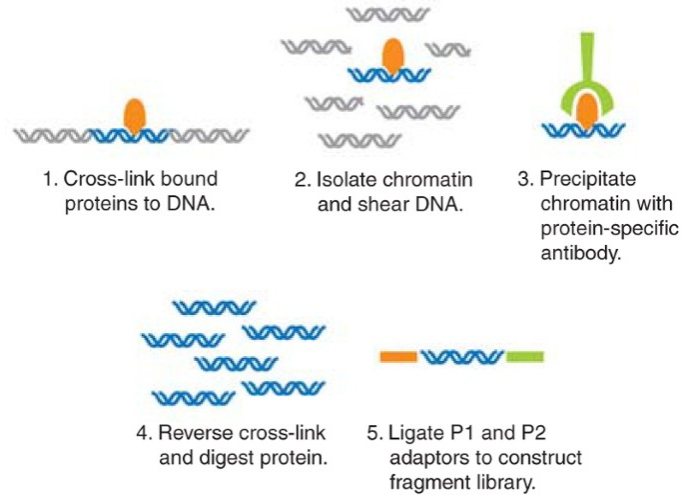
When studying chromatin, it is often very interesting to know where it is more open and accessible, and where the CTCF binds to the structure. To know this information, it is possible to perform ATAC-seq and CTCF ChIP-seq, respectively. Below is a brief description of the two methods:

- **ATAC-seq** is a technology that allows for the identification of open chromatin regions^{8,9}. In order to work, it requires the addition of Tn5, a hyper-active transposase. The latter is preloaded with sequencing adapters⁹ to induce a contemporaneous reaction of fragmentation and ligation of the pieces released, in a process called segmentation. The obtained adapted fragments are then amplified and sequenced. Once the reads are generated, a peak-calling algorithm (generally MACS-2¹⁰) is used to determine which portions of the genome present ATAC peaks, and areas where there are significant enrichments of aligned reads with respect to the background. A significant enrichment of reads is possible only in accessible regions, which are generally also the most active ones and with available sites for transcription factors binding.
- The **CTCF ChIP-seq** data used for this experiment, named in table 1, were obtained through a classical procedure: it consisted of combining a process of chromatin immunoprecipitation, made with antibodies specific for CTCF, and one of DNA sequencing¹¹. The scope of the technique is to infer the possible binding sites of the transcription factor on the DNA.

More details about the used data can be found by consulting the ENCODE¹² entries written in table 1.



(a) Basic concept behind the ATAC-seq technique. The images on the left and on top were taken from the following [link](#)¹³.



(b) Pipeline to perform a generic ChIP-seq analysis. The image was taken from the work of Anjali Shah¹¹.

Figure 2. The ATAC and CTCF procedures explanations.

Very interestingly, chromatin portions can be associated to a finite number of states, called chromatin states, on the base of the epigenetic data they produce¹⁴. Some machine learning approaches, like *ChromHMM*¹⁵ were built with the intention of predicting these configurations. The next chapter (2.3) will talk about that.

2.3 *ChromHMM* allows the sequential characterization of chromatin states

*ChromHMM*¹⁵ is a tool which helps in the annotation of genomic DNA by using epigenomic information¹⁴. It learns chromatin states signatures by using a multivariate hidden Markov model: in each genomic position (segment), it returns the most probable chromatin state and gives other useful information^{14,15}.

The package works through two functions in particular, which are the following¹⁴:

1. ***BinarizeBam***: it converts a set of *.bam* files of aligned reads into binarized data files in a specified output directory. The produced data can be used as input for the *LearnModel* function. When using this command, it has to be specified the segment size, which is set by default to be equal to 200 bps.
2. ***LearnModel***: it takes a set of binarized data files, learns chromatin state models, and by default produces data reporting the emission/transition parameters of the states, the abundance of the states at the TSSs (Transcriptional Starting Sites), at the TESs (Transcriptional Ending sites), and other relevant portions of the genome (CPG islands, exons, genes). Additionally, a webpage is generated with links to all the files and images created.

The results obtained from *ChromHMM* are shown in chapter 5.1.

2.4 Hi-C matrices and their patterns

Hi-C maps are useful tools to detect the interactions occurring inside a genome in analysis. Indeed, they allow to gain insights about the structural disposition of chromatin domains, loops and regions.¹⁶ The experimental procedure is described in figure 3. Each position (i, j) in a Hi-C matrix represents the number of contacts occurring between the coordinates i^{th} and j^{th} of the segment. The resolution of an Hi-C map will be dependent on the sequencing process, and have to be decided on the base of the type of information that has to be recovered from the data¹⁶.

The following list of patterns can be found by inspecting an Hi-C matrix^{16,17}.

1. **Cis/trans interaction ratio:** There are higher interaction frequencies on average between pairs of *loci* in the same chromosome (*cis*), with respect to those among *loci* which reside on different chromosomes (*trans*). The ratio between *cis/trans* interactions could be indicative of the quality of the obtained Hi-C data. The viewed specificity is related to the presence of genomic territories, which govern and establish the disposition of the chromosomes in the nucleus¹⁸. The *cis*-interactions can be easily seen in an Hi-C matrix along the diagonal and are depicted in panel *a* of figure 4.
2. **Distance-dependent interaction frequency:** From the visualization of an Hi-C matrix, it is possible to observe that the largest number of interactions are registered at small distances. On the other hand, only a few contacts can be observed with high spacial separations. Because of this recurring pattern, several studies tried to predict this interesting behavior. Importantly, it was found that in a number of situations it is possible to do that: for example, in yeast the probability of interaction could be described with the following equation¹⁶:

$$p_{\text{interaction}}(x, y) = Z * \text{dist}(x, y)^{-1,5}$$

Where $\text{dist}(x, y)$ represents the distance between a point x and a position y .

3. **Genomic compartments:** Genomic compartments (which can be seen in panel d of figure 4) have been found to be correlated with chromatin states, involving DNA accessibility, gene density, replication timing, GC content and histone marks¹⁹. The compartments can pertain to two categories, A and B, and are found by performing a principal components analysis with the matrix generated with Pearson Correlation coefficients (a formula for their calculation can be found in chapter 4.6). In general A-type compartments are defined as the euchromatic gene-dense regions, while B blocks are defined as gene-poor heterochromatic regions. The positions where they are found differ depending on the type and biological conditions of the analyzed cells¹⁶.
4. **Topological domains:** Also called TADs, they can be visually found in Hi-C matrices as larger squared boxes centered on the diagonal of the maps (panel *b* figure 4). They are contiguous portions in which *loci* tend to interact much more with each other than with *loci* outside the region¹⁶. It is hypothesized that TADs specify elementary regulatory micro-environments in which promoters interact with local enhancers. Finally, some proteins like cohesins and CTCF tend to interact with the genome at the boundaries of the Topological Domains¹⁶.
5. **Point interactions:** Those are connections occurring between small regions, and involve sequences of a few kb length. Biologically speaking, those points could indicate for example the interactions between enhancers and promoters. When considering a specific point connections, the observed value should be compared to the expected number of interactions for the distance in analysis, and the significance should be computed¹⁶.

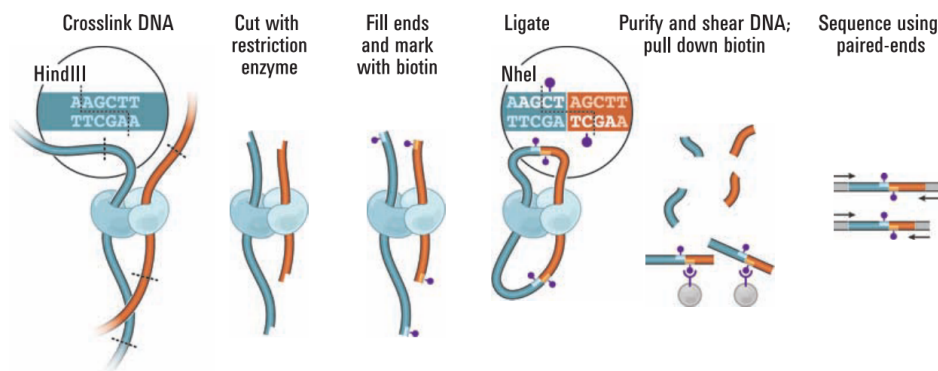


Figure 3. Image taken from the following [link¹⁹](#), representing the Hi-C sequencing technique in a schematized way.

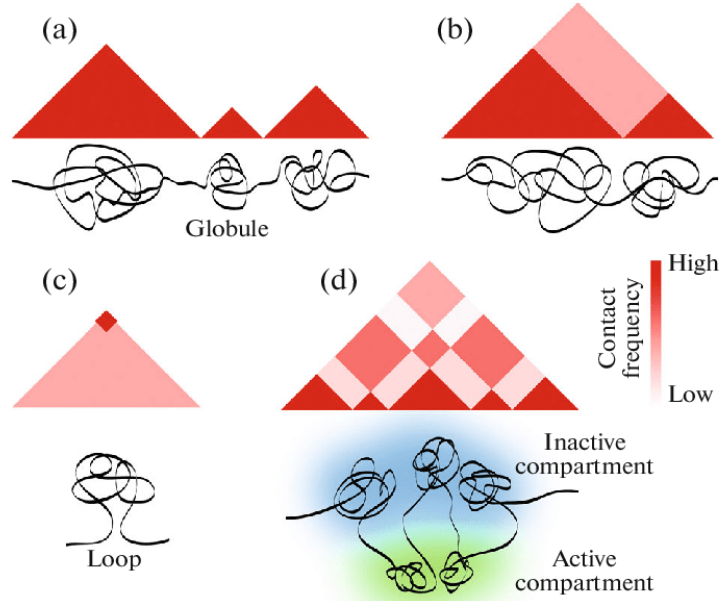


Figure 4. The image was taken from the work of Razin and colleagues²⁰. Figure representing the contact typologies described in this chapter. (a) Isolated triangles on a heat map are commonly interpreted as chromatin globules deriving from cis-interactions. (b) a TAD, represented as a combination of small triangles into larger ones. (c) An intense signal at the apex of a triangle suggests the interaction of TAD boundaries and the formation of a chromatin loop. (d) Representation of some chromatin compartments.

2.5 Chromatin Coarse-Graining, chromatin as a polymeric fluid

The coarse-graining procedure performed for this project was inspired by the article written by Kremer and Grest in 1990²¹.

In a very interesting chapter, which is summarized in this section, the book "*Giant molecules: here, there, and everywhere*"²² compares the DNA to a chain of beads, and more in general, to a polymeric fluid inside the cell nucleus. It affirms that the random motion of a DNA fiber could be compared to the stochastic Brownian dynamics of a particle: several small connected segments would move randomly, however maintaining their order and connectivity. It is easy to understand that, given this hypothesis, the most important quantity that should be computed to perform the coarse-graining of a polymer is the length of the rigid segments.

To perform the following steps and derive that value, also called as Kuhn length, it is made the assumption that a DNA polymer moves in a Brownian manner. Although this type of dynamics would largely reduce the volume of the chromatin, still the dimension of the latter would be too high to permit its entire confinement inside the nucleus of any cell. For this reason, it was suggested that some other mechanisms intervene to rule the condensation behavior of the filaments.²². For a DNA polymer, the Kuhn length is thought to be approximately equal to 100 nm.

To start, it is defined the difference between a Brownian motion and a straight movement. That discrepancy could be written as follows:

- For a straight motion $R = v(t_2 - t_1)$
- For a Brownian particle $R = l_{\text{eff}}^{1/2} [v(t_2 - t_1)]^{1/2}$

Where $R = |\vec{R}_1 - \vec{R}_2|$ is the difference between the initial position \vec{R}_1 and the final one \vec{R}_2 . The obtained equation can be rewritten as a square-root displacement in the following way:

$$R = l_{\text{eff}}^{1/2} [v(t_2 - t_1)]^{1/2} = \langle (\vec{R}_2 - \vec{R}_1)^2 \rangle^{1/2}$$

By easily substituting $v(t_2 - t_1)$ with L , it is possible to derive the equation for a polymer, which becomes

$$R = l_{\text{eff}}^{1/2} * L^{1/2}$$

Where L is the maximal possible length of the polymer, and is called contour length. By computing the squared value of the previous equation, it is possible to obtain

$$\begin{aligned}
 R^2 &= l_{\text{eff}} * L = \langle \vec{R}^2 \rangle \\
 \rightarrow l_{\text{eff}} &= \frac{R^2}{L}
 \end{aligned} \tag{1}$$

Importantly, the derived equation 1 contains the definition of the Kuhn length l_{eff} . This quantity allows to understand the degree of bendability of the chain. By using a more visual approach, this length could be considered as a memory that is maintained along a path on the polymer. Indeed, keeping in mind the idea of following a "journey" on the chain, the average angle that is obtained at a contour length s , obtained through the intersection of the tangents at the starting and the ending points of the segment (figure 5), is big or low depending on the ratio between the analyzed and the Kuhn lengths. By looking to the equation 2, in general, the lower is the contour segment inspected with respect to the Kuhn length, the higher is the probability of having a low degree angle ($\langle \cos \theta \sim 1 \rangle$). On the contrary, by analyzing larger lengths, it is possible to obtain a wider range of angles, with a calculated cosine that becomes $\langle \cos \theta \sim 0 \rangle$.

$$\langle \cos \theta(s) \rangle = \exp\left(-\frac{s}{l}\right) \tag{2}$$

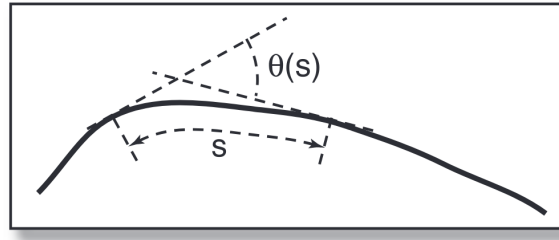


Figure 5. Image taken from Grosberg *et al.*²². Angle formed between the tangents at the extremes of a contour length.

To conclude, the Kuhn length is directly related the persistence length, which is another quantity, exploited in the polymer model of this work to compute the angle bending potentials, which are calculated through the equation 5. For this project, the relationship between the persistence and the Kuhn lengths was parametrized as in the formula 6.

Because of the fact that it is not really clear which and how many stages of compaction exist (chapter 2.1), it was decided to generate chains with beads including 5000 bp. The parameters for the most coarse-grained model (CG) were partially obtained by taking into account a Fine Scale (FS) configuration.

3 Aim of the project

The current project is a part of the thesis work, whose aim is to predict chromatin matrices of contact through the results obtained with molecular coarse-grained simulations of 2 Mb long chains. Our focus was in particular placed on a region including and surrounding the ANPEP *locus*, starting at position 89,000,000, and finish at the 91,000,000th base of the 15th. The type of system that was generated had beads incorporating 5000 bp (CG model). As a last step, which has still to be done, a comparison between this modelling approach and others currently available, such as *Hip-Hop* and *cOrigami*^{23,24}, would give us a better idea about the potential of the proposed method.

The scope of this work, and specifically of this report, is also to gather opinions and useful feedbacks to improve the project, which will continue during the next months.

4 METHODS

All the simulations were performed making use of the simulation software LAMMPS^{25,26}, and some already made codes of Marco di Stefano, researcher at IGH-CNRS (France).

4.1 Data used during the project

The data of CTCF and ATAC for the IMR90 cell line, included in the paper written by Jimin and colleagues in 2023²⁴, were used for the project (see table 1). It was decided to exploit the same data of *cOrigami* to allow a better comparison between its predictions and those produced by our modelling, which will be done as a last step. Both the two technical replicates, included in the listed ENCODE¹² entries, were considered.

Cell-Type	CTCF ChIP-seq	ATAC-seq
IMR90	ENCSR000EFI	ENCSR200OML

Table 1. Table referring to the data used for the analysis.

The ATAC-sequencing data were produced following the standard ENCODE procedure²⁷. In particular, the processes of read trimming, alignment, and filtering were performed making use of the *Bowtie 2*, *Samtools*, *Sambamba*, *Picard* and *cutadapt* softwares²⁸. An explanation of the named processes could be found in the work made by Feng Y. and colleagues in 2020²⁹.

When it comes to the ChIP-sequencing data, again, the standard procedure of ENCODE was used to produce the online available datasets. An overview about the method can be found at the following link³⁰. To sum up, at first the reference genome was indexed with the *BWA*, then, the alignments between the reads and the reference genome (*hg38*³¹) were produced and filtered with the *BWA*, *Samtools*, *Picard*, *BEDTools*, *Phantompeakqualtools* and *SPP* softwares.

In the case of our study, the analyzed region included ANPEP, and it was taken from the 89,000,000th base to the 91,000,000th position of the 15th chromosome.

4.2 The Model

The model creation was done by using and modifying scripts and code written by Marco Di Stefano. The code is included in the *TADphys* unpublished package. The objective of the modelling process was to create a CG polymer model with a resolution of 5000 bp. To start the analysis, some information about the IMR90 cell type had to be collected. For example, the total genome length had to be determined, and was obtained by consulting the UCSC genome browser³². Secondly, the dimensions of the nucleus and the nucleolus of the IMR90 cell had to be specified. By consulting some accessible literature^{33–35}, the volumes of the two structures were respectively equated to 520 μm and 100 μm .

With the aim of being able to make statistics out of the generated data, and have more chains to analyze, it was decided to produce simulations for 100 replicates. Additionally, three chains were simulated at the same time within each repetition. Everything was included inside a box with a volume written in table 3. It is important to specify that all the simulations were performed making use of periodic boundaries, and were run with the *run_lammps* function included in the *TADphys* package. Three types of potential energies were set and always present in all the simulations that will be presented in this chapter. Those are described in the list below:

- **FENE potential:** The FENE potential is a finite extensible nonlinear elastic potential energy and is usually used for polymer models^{25,26}. The first term in the equation is attractive, whilst the second one is repulsive and is a Lennard-Jones (LJ) potential. The first term extends until R_0 , the maximum extent of the bond. The second term has a cutoff set at $2^{\frac{1}{6}}\sigma$, where the value of the LJ potential found is the minimum²⁶. Indeed, in that position, $V_{\text{LJ}} = -\epsilon$. As usual, the σ in the LJ formulation is the distance at which the intermolecular repulsive potential between two particles is zeroed.

$$E = -0.5KR_0^2 \ln \left[1 - \left(\frac{r}{R_0} \right)^2 \right] + 4\epsilon \left[\left(\frac{\sigma}{r} \right)^{12} - \left(\frac{\sigma}{r} \right)^6 \right] + \epsilon \quad (3)$$

The following specifications were made to include the FENE interactions:

1. K (energy/distance²) = 30.0: It weights the contributions deriving from the attractive part of the equation, and represents the stiffness or strength of the bond.

2. R_0 (distance unit) = 1.5: Maximum extension of the bond. This is the maximum distance at which the bond can be stretched.
 3. ϵ (energy unit) = 1.0: This term scales the LJ potential contribution.
 4. σ (distance unit) = 1.0: It depicts the equilibrium bond length for the LJ potential. It represents the ideal or equilibrium distance between the bonded particles.
- **Excluded-volume interactions:** To allow for the excluded-volume interactions, a simple Lennard-Jones potential was included, set as depicted in equation 4.

$$E_{LJ} = 4\epsilon \left[\left(\frac{\sigma}{r} \right)^{12} - \left(\frac{\sigma}{r} \right)^6 \right] \quad r < r_c \quad (4)$$

Three parameters are set:

1. ϵ (energy unit) = 1.0
 2. σ (distance unit) = 1.0
 3. LJ cutoff r_0 (distance unit) = $\sigma * 2^{\frac{1}{6}} = 1.12246152962189$
- **Angle bending potentials:** The angle bending associated potentials were set through the *angle_style* and *angle_value* functions of LAMMPS²⁶, and were calculated with the formula 5.

$$E = K[1 + \cos(\theta)] \quad (5)$$

In particular, K was equated to the persistence length, which was calculated as in equation 6, and its value is presented in table 3.

$$PL = lk_{CG}/b_{CG}/2 \quad (6)$$

The larger is the value of K , the bigger is the potential generated after the bending of the chain with a specific angle θ .

When it comes to the effective modelling process, the performed steps are described in the following paragraphs:

The computation of the parameters for Coarse Graining: Both parameters for the fine-scaled (FS) model (table 2) and the CG model (table 3) were calculated. The values for some of the variables obtained in the first case were used to compute and derive some of the latter system (see table 3). The number of bp wrapping around a bead in the FS method was considered to be 150, while instead the linker portion was set to be of length 50 bp (table 2). The thickness of a bead, which corresponds to a nucleosome was taken as equal to 10 nm, while instead the default Kuhn length was set to 50 nm. The genome densities ($\rho_{FS} = \rho_{CG} = 0.012 \text{ bp/nm}^3$) were imposed to be the same for both the FS and the CG model; this value has been found experimentally and represents the density of bp inside the human nuclei³⁶. Because of that, before computing the parameters for the coarse-grained model, the DNA content of the Kuhn segments in the CG system (Dlk_{CG}) was tuned to match the desired value of DNA content in CG beads (v_{CG}). The spheres and bonds had all the same length in the FS and the CG models, consequently, the contour lengths (which represent the maximal lengths of the polymer chains) were exactly calculated as the products between the number of beads and their size.

Generation of the initial conformation rosettes: Once the coarse graining parameters were calculated, rosettes for all the replicates were built, with *radii* of 12.0 nm inside a cubic box whose edge length was equal to 300 nm. The particle *radius* of the CG model was set to 0.5 nm. In each replicate, three equal chains were built, by setting a different random seed each time. The total number of particles in each chain was of 400 beads. Indeed, if the calculation is made, $2,000,000/5,000 = 400$.

Finding the optimal pressure: When the rosettes were successfully made, a decompaction was performed taking as inputs the compacted configurations. A range of values of pressure was tested from 0.1 to 1 with steps every 0.1. The obtained sizes were compared to the target calculated size (table 3). The precision of the estimation for the pressure was improved by taking more decimals and restricting the length of the range of tested values. For each replicate, a new random seed was generated and stored, again. The simulations done in order to find the optimal pressure parameter were long 1000 steps with a duration of 0.001 ps. Before attempting the decompression, the minimal energy structure was found by taking into consideration a stopping energy tolerance of $1 * 10^{-4}$, a stopping tolerance force of $1 * 10^{-6}$, and a maximal number of iterations and evaluations of 100,000 steps. At the end, the optimal pressure was attested to be at 0.192.

Decompaction and relaxation: Once the optimal value for the pressure was found (0.192), other two simulations, respectively 5,000,000 and 25,000,000 steps long, were performed for each replicate. This time, the step was set to have a temporal length of 0.0012 ps. In both the cases MSD values were collected every 100 steps. A frame was dumped (saved) every 5,000 steps. At the end, the trajectories were collected and analyzed by computing the *RMSD*, the *Rg* and the autocorrelation function as described in chapter 4.4. For the sake of simplicity, to save memory and computational costs, the collection was accomplished by capturing one frame every 50,000 of them.

Computing matrices of contact: Once defined the step at which the simulations were considered to be at the equilibrium (which, in our cases, seemed to be reached at the $50 * 50000 = 2,500,000$ step, as reported in chapter 5.3), some dictionaries produced through the output of *ChromHMM*^{14,15}, whose input were CTCF and ATAC-sequencing datasets (chapter 4.1), were used to define the identity of the 5,000 base pairs beads. Then, the best attraction parameters were selected in the iterative manner described by the 3 procedure. To add the attraction potentials between the beads, new Lennard-Jones energies (equation 4) were added to the systems. In particular, the cutoff distance of the potential r_0 was set to be equal to: $r_0 = r_{\text{cutoff}} = \sigma * 2.5$, where σ corresponds to the sum of the *radii* of the interacting particles (set in all the cases to 0.5 nm).

Once the interaction parameters were set, other 1 second long simulations were performed with the intention of generating contact maps. Those simulations were firstly preceded by an energy minimization process, very similar to the one already explained in the "Finding the optimal pressure" paragraph. After the small simulation time, the contact matrices were generated, by taking into account just the interactions occurring in the same chain (intra-chain). A contact was established whenever two coarse-grained particles were found at a distance lower than 100 nm.

4.3 Finding enriched states in 5 kbs long beads

To find the enriched states in 5 kb long bins, the procedure described in process 1 was followed. The fold change of a state within a bin was determined by dividing the proportion in the bin by the corresponding proportion in chromosome 15. Once done that, the state with the highest fold-change was assigned to the bin. A visual investigation was performed afterwards to check the quality of the "binarization" procedure by using the IGV visualization software³⁷.

Algorithm 1: Finding enriched states in 5-kb long bins

```

Result: Enriched states
forall chromosomes do
  | Find proportions of the states in the chromosome
foreach bin do
  | Calculate bin proportions for each state
  | Compute the fold changes
  | Assign the state with the highest fold change
Generate .bed files
Visualization in IGV of regions of interests
/* each bin is 5 kb long

```

*/

4.4 Trajectories analysis

The analysis of the trajectories was done by taking together and considering as independent the results obtained from each single chain of each replicate. For this reason, it was also decided to put together all the outputs and produce the collective plots represented in figures 10b, 11b and 12b.

Three types of analysis were performed:

1. **Root Mean Square Deviation (RMSD):** The RMSD is evaluated by using the *MDanalysis* package³⁸ (*rmsd* in *MDAnalysis.analysis.rms*) and is calculated as follows:

$$RMSD = \rho(t) = \sqrt{\frac{1}{N} \sum_{i=1}^N w_i (\vec{x}_i(t) - \vec{x}_i^{REF})^2} \quad (7)$$

Before performing this type of calculation, the structures were aligned to the first frame (each frame of each replicate was aligned to the first frame of the replicate). This type of alignment was done making use of the *AlignTraj* function³⁸. When interpreting the results obtained from RMSD calculation, it is generally considerable true the concept that the smaller is the difference between two structures, the lower is the value of RMSD. The results are written in graph 10a and 10b.

2. R_g : The Radius of Gyration was calculated as written in equation 8 through the *MDanalysis* package (*radius_of_gyration* function). This quantity is a measure of how the mass of an object is spread out relative to a particular axis of rotation. In general, it tells "how spherical" is an object^{38,39}; the higher is the value of the radius of gyration, the lower is the sphericity of the substance. The results are written in graph 11a and 11b.

$$R_g = \sqrt{\frac{\sum_i m_i \vec{r}_i^2}{\sum_i m_i}} \quad (8)$$

3. **Autocorrelation function:** The autocorrelation function represents the degree of similarity between a given time series and a lagged version of itself over successive time intervals. It can be written as shown in equation 9⁴⁰. The results are depicted in graphs 12a and 12b.

$$r_k = \frac{C_k}{C_0} = \frac{\frac{1}{M} \sum_{t=1}^{M-k} (A_t - \bar{A})(A_{t+k} - \bar{A})}{\frac{1}{M} \sum_{t=1}^{M-k} (A_t - \bar{A})^2} \quad (9)$$

Where $C_k = \frac{1}{M} \sum_{t=1}^{M-k} (A_t - \bar{A})(A_{t+k} - \bar{A})$ is the autocovariance function at lag k and $C_0 = \frac{1}{M} \sum_{t=1}^{M-k} (A_t - \bar{A})^2$ is the variance function.

4.5 Algorithms used for comparison

As stated in the methods section 4.2 about the simulations, the attraction potentials are sequentially added to the model in order to improve the predictions. The SCC metric (equation 10) was used to compute the difference between the control contact matrix and the CG derived one (chapter 4.6).

Two methods were then considered to add the variables, which are expressed in algorithm 2 and 3. Since the states are differentially populated, as stated in chapter 5.1, it is possible to argue that the largest contribution to the result would be given, in a hierarchical manner, by the most populated states. As a consequence of this consideration, it should be possible to fix the values related to the most prevalent states, before considering those that are less present. This type of mechanism is described in the greedy process of algorithm 2. If that assumption is not accepted, then either all the possible configurations have to be considered, either a better way to test the generated models should be thought.

Because of the fact that beads of different states could act coordinately in defining the quality of the results, the algorithm 3 was devised. It represents a stepwise solution which allows to solve partially the problem at the cost of more simulations to perform. In fact, any choices of parameters are deleted and rethought during the process, in a manner that tries to avoid local *optima*.

Algorithm 2: Greedy matrix comparison

Result: Best performing greedy model

forall attraction parameter **do**

Construct models by adding the most present attraction parameter to the previous step configuration ;
 Compute SCC of the model with respect to the reference ranging among a list of possible values;
 Select the value of the parameter which gives the best results;
 Add that parameter (attraction) with that value to the model;

return Best greedy model;

Algorithm 3: Step-wise process for matrix comparison**Result:** Best model $loop = 0;$

Continue = False;

while Continue == True **do** $loop = loop + 1;$

Continue = False;

 Construct models by adding the most present attraction parameter to the $(loop - 1)^{th}$ step configuration ;

Compute SCC of the model with respect to the reference ranging among a list of possible values;

Select the value of the parameter which gives the best results;

if addition gives better results **then**

Add that parameter (attraction) with that value to the model;

Continue = True

 $loop = loop + 1;$ **forall** state in model **do**

Remove that state from the model;

Vary the value associated to the state with the highest frequency among those remaining;

Compute SCC of each model with respect to the reference ranging among a list of possible values

Select the reduced model which gives the best results;

if removal gives better results **then**

Perform the reduction;

Continue = True

return Best model**4.6 the Stratum Adjusted Correlation Coefficient (SCC) metric**

The SCC metric is described in the paper written by Yang and colleagues in 2017^{41,42}. It can quantify the similarity between an Hi-C matrix and another. In general, the most common techniques to use in these situations are either to analyze the matrices by eye, or, in a certainly more precise way, to calculate a Pearson/Spearman correlation coefficient. However Hi-C data have certain unique characteristics, including domain structures, such as topological association domain (TAD), A/B compartments and distance dependencies, that require a more precise approach. Indeed, the chromatin interaction frequencies between two genomic loci, on average, decrease substantially as their genomic distance increases. Standard correlation approaches do not take into consideration these structures and may lead to incorrect conclusions^{41,42}.

The SCC metric could be seen as a weighted Pearson coefficient, as written in equation 10.

Variables

N_k	$k \in K$	Number of observations in stratum k ;
X_k	$k \in K$	Observations in stratum k in matrix X ;
Y_k	$k \in K$	Observations in stratum k in matrix Y ;
$r_{1k} = \frac{\sum_{i=1}^{N_k} x_{ik} y_{ik}}{N_k} - \frac{\sum_{i=1}^{N_k} x_{ik} \sum_{j=1}^{N_k} y_{jk}}{N_k^2} = E(X_k Y_k) - E(X_k)E(Y_k)$	$k \in K$	Correlation between X_k and Y_k ;
$r_{2k} = \sqrt{\text{var}(X_k) \cdot \text{var}(Y_k)}$	$k \in K$	Square root of the product between the variances of X_k and Y_k ;
$\rho_k = r_{1k} / r_{2k}$	$k \in K$	Pearson coefficient related to bin k ;

Formula

$$\rho_s = \sum_{k=1}^K \left(\frac{N_k r_{2k}}{\sum_{k=1}^K N_k r_{2k}} \right) \rho_k \quad (10)$$

5 RESULTS AND DISCUSSION

5.1 ChromHMM results

In total, 4 states were considered to be present. Two functions in particular were used: *BinarizeBam* and *LearnModel*. The data shown in 4.1 were aligned to *hg38* reference genome³¹. Results are shown in image 6; by taking a look to its subfigures, the following considerations could be done:

1. Clear absence or presence of ATAC and/or CTCF signals could be detected in figure 6a. for this reason, the following states are defined:
 - **State 1:** State without the presence of ATAC and CTCF signal.
 - **State 2:** State with ATAC but not CTCF peaks.
 - **State 3:** State with the presence of both ATAC and CTCF signal.
 - **State 4:** State with CTCF but not ATAC peaks.
2. The states 1 and 2, in particular, tend to perform transitions towards themselves instead of different states (figure 6b)
3. State 2 (with ATAC) and 3 (with ATAC and CTCF) tend to localize in CpG islands, exons and Transcriptional Starting Sites (TES) (figures 6c, 6d, 6e) as well as, although to a lesser extent, in Transcriptional Ending Sites (TES).

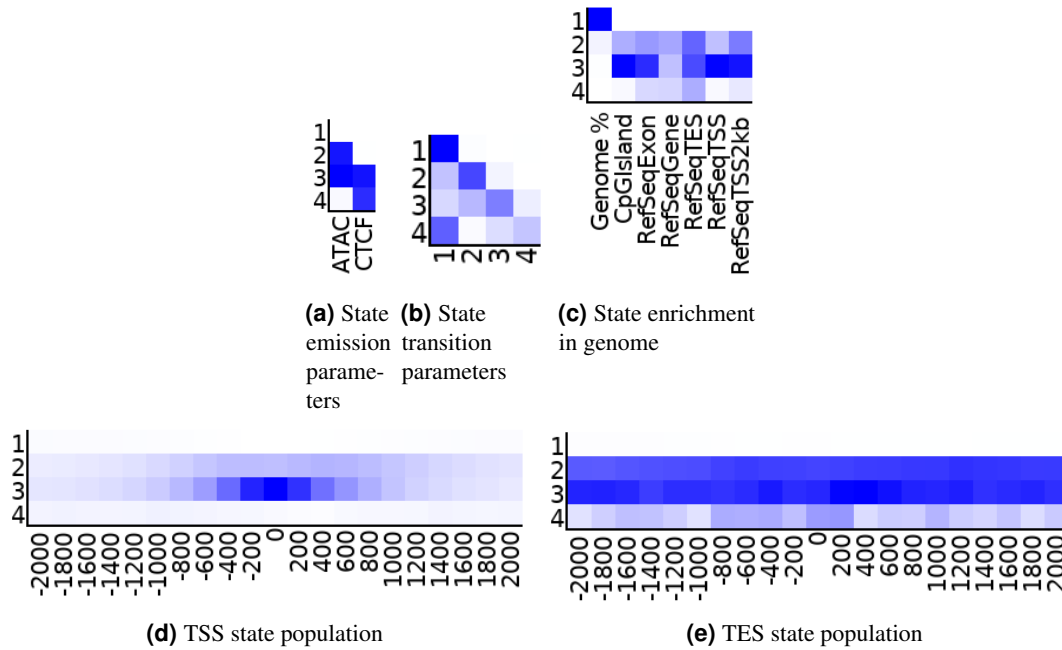
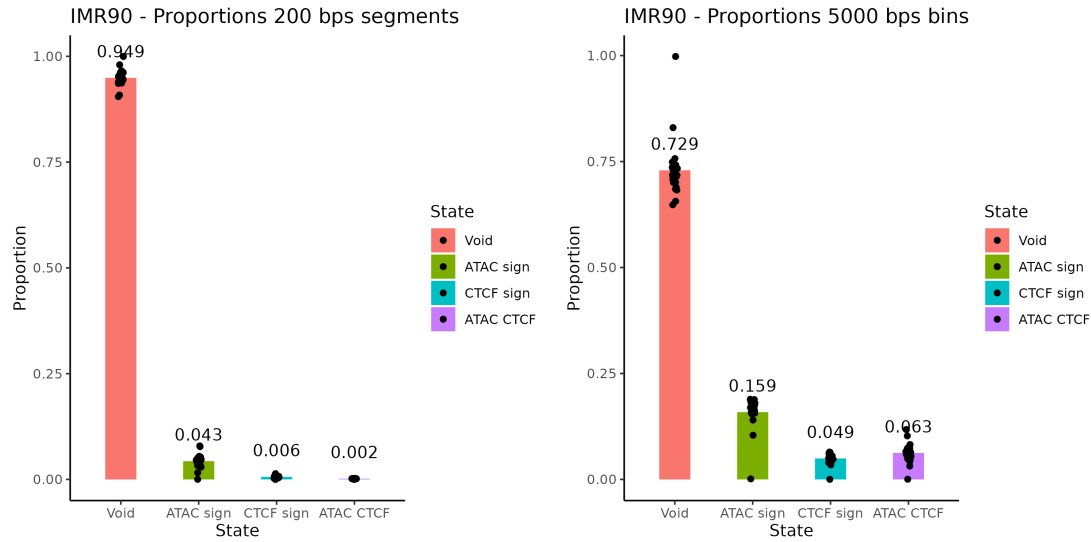


Figure 6. Results from *ChromHMM* for the IMR90 replicates

The proportions represented in image 7b were obtained for the 5 kbp long bins. In general, the quantities calculated taking into consideration the 200 bp segments were much lower with respect to those calculated with the 5kb bins. The reason is that the first state is in general much more present than the other three, and whenever there are a few occurrences of the rarely present states, the long bins are assigned with a great probability to those states, and as a consequence the proportions tend to smooth out their differences.



(a) Proportions of states in 200 bps segments. The segments were directly found by *ChromHMM*. Each dot represents the proportion relative to a chromosome.

(b) Proportion of states in 5000 bps long bins, obtained through the method described in process 1. Each dot represents the proportion relative to a chromosome.

Figure 7. Proportions found in 200 bp segments and in 5000 bins.

The following image (8) was created by using IGV, a visualization tool^{37,43}. After a visual inspection of the results, it was decided to trust the assignment performed. However, some defects become evident while viewing the results: whenever the *ChromHMM* signaled the presence of the 4th state, the relative bin was assigned to it. What happened was that, when the fourth state was found, if the 3th (with both ATAC and CTCF) was not signaled enough, the information about the presence of ATAC peaks was lost. Problems about the precision of the state assignment process couldn't be easily solved and are a direct consequence of the coarse-graining process. Indeed, it's impossible to retain the same amount of information before and after the bin simplification.



Figure 8. IGV snapshot of a portion of the ANPEP region analyzed. The first track reports the alignment results obtained from the ATAC data, the second the CTCF data. The straight traces in the bottom represent the state of the bins in that genomic region. An orange trace represents the state 1 (no ATAC and not CTCF), while instead the purple highlighting signals the third state (with ATAC and CTCF).

5.2 Results obtained while defining the models

The results in this section will be inserted by following the paragraph order written in section 4.2.

- **The computation of the parameters for Coarse Graining:** The values in table 2 and 3 were obtained.

Property	Formula	Value
c	<i>const.</i>	19
v_{FS} (DNA content of a monomer in b.)	<i>const.</i>	150+50 bp = 200 bp
b_{FS} (Diameter of a bead in nm)	<i>const.</i>	10 nm
lk_{FS} (Kuhn length of the chain in FS)	<i>const.</i>	50 nm
ρ_{FS} (Genome density)	<i>const.</i>	0.012 bp/nm ³³⁶
N_{FS} (Number of monomers to represent the chromosome)	$\frac{DNA_{content}}{v_{FS}} * ncopies$	30000 mon.
N_{FS}^k (Number of Kuhn lengths of the chain)	$\frac{N_{FS} * b_{FS}}{lk_{FS}}$	6000 k. l.
ρ_{FS}^k (Genome density in Kuhn lengths)	$\frac{\rho_{FS} * b_{FS}}{v_{FS} * lk_{FS}}$	1.2e – 05 1/nm ³
L_{FS} (Polymer contour length)	$N_{FS} * b_{FS}$	300000 nm
Le_{FS} (Entanglement length of the chain in nm)	$lk_{FS} * \left(\frac{c}{\rho_{FS}^k * lk_{FS}^3} \right)^2$	8022.22 nm
Number of monomers in a Kuhn length FS	lk_{FS} / b_{FS}	5 mon.
Blk_{FS} (Bead content of a Kuhn length FS)	$(lk_{FS} * b_{FS}) / v_{FS}$	2.5 nm ² /bp
Dlk_{FS} (DNA content of a Kuhn length FS)	$(lk_{FS} * v_{FS}) / b_{FS}$	1000 bp

Table 2. Parameters calculated for the Fine Scale (FS) model

Property	Formula	Value
c	<i>const.</i>	19
v_{CG} (DNA content of a monomer in b.)	<i>const.</i>	5000 bp
Dlk_{CG} (DNA content of a Kuhn length CG)	<i>tuned const.</i>	33791 bp
ϕ_{CG} (Volumetric density of the chain in the CG model for IMR90 cell-type)	<i>const.</i>	0.1
ρ_{CG} (Genome density in bp/nm ³)	<i>const.</i>	0.012 bp/nm ³
b_{CG} (Diameter of a bead in nm)	$\sqrt{\left(\sqrt{\frac{Dlk_{CG}}{Blk_{FS}}}\right) / \rho_{CG} \cdot \frac{6}{\pi} \cdot \phi_{CG}}$	43.0155 nm
lk_{CG} (Kuhn length of the chain in CG)	$\sqrt{Dlk_{CG} * Blk_{FS}}$	290.65 nm
Number of monomers in a Kuhn length CG	lk_{CG} / b_{CG}	6.75687 mon.
N_{CG} (Number of monomers to represent the chromosome)	$\frac{DNAcontent}{v_{CG}} * ncopies$	1200 mon.
side _{CG} (size of the cubic simulation box)	$\frac{(N_{CG} * v_{CG} / \rho_{CG})^{1/3}}{b_{CG}}$	18.4515 nm
N_{CG}^k (Number of Kuhn lengths of the chain)	$(N_{CG} * b_{CG}) / lk_{CG}$	177.597 k. l.
ρ_{CG}^k (Genome density in Kuhn lengths bp/nm)	$\frac{\rho_{CG} * b_{CG}}{v_{CG} * lk_{CG}}$	3.55194e-07 bp/nm
L_{CG} (Polymer contour length)	$N_{CG} * b_{CG}$	51618 nm
Le_{CG} (Entanglement length of the chain in nm)	$lk_{CG} * \left(\frac{c}{\rho_{CG}^k * lk_{CG}^3}\right)^2$	1379.51 nm

Table 3. Parameters calculated for the coarse-grained (CG) model

- **Finding the optimal pressure:** The values of pressure and the respective sizes are plotted in figure 9. To perform this step, just 5 replicates of the 100 total replicates were used for simplicity.

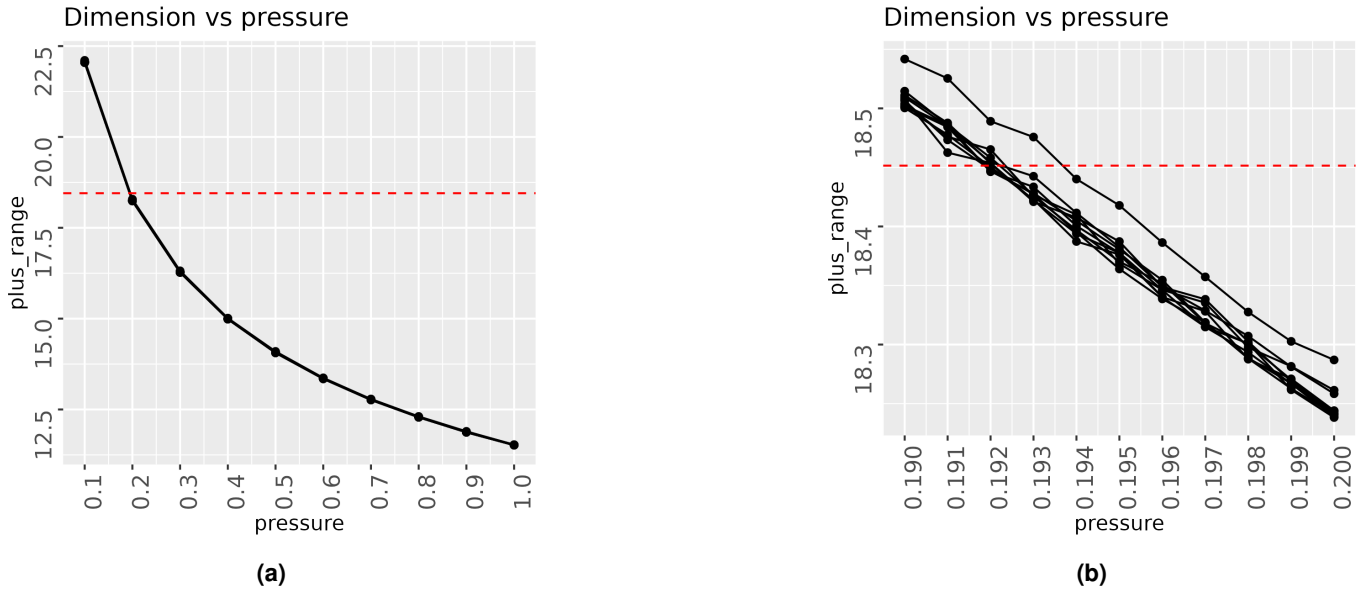
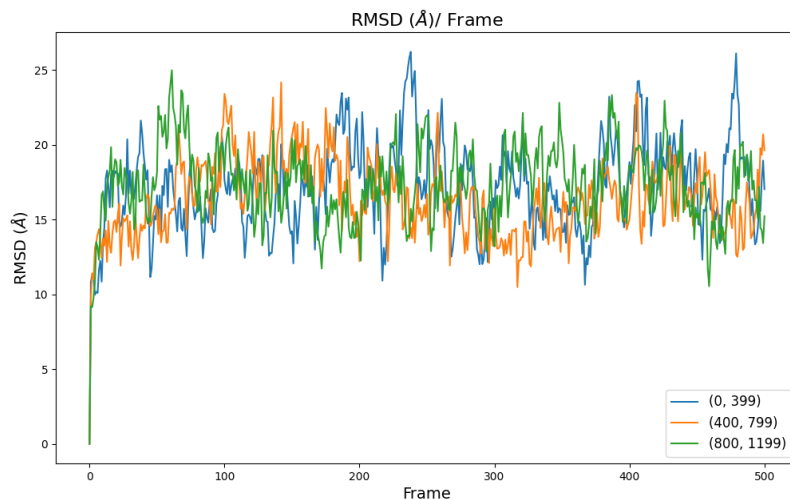


Figure 9. The side estimates obtained for different pressure values are represented in the two graphs. As said in the Methods chapter, the range of trial values for the pressure was shortened by adding more decimals to the quantity. For example, at first the values between 0.1 and 0.9, with a difference of 0.1, were tested (a); after only the region between 0.19 and 0.2 was investigated (b).

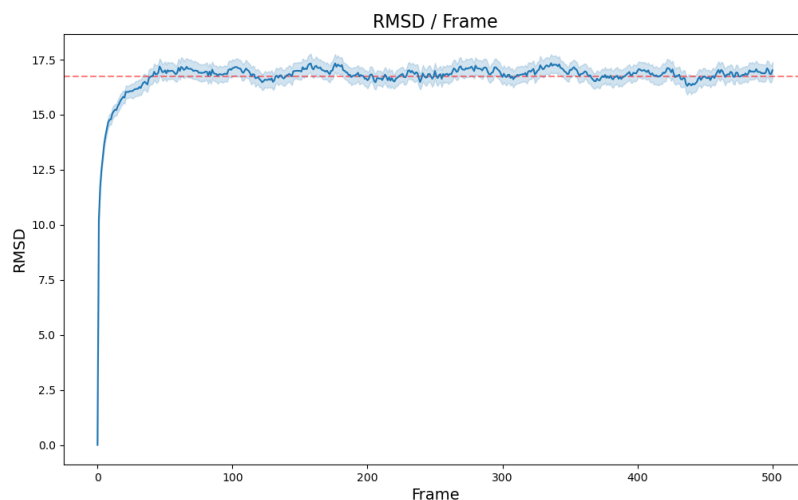
5.3 Trajectory analysis results

Results from the RMSD analysis are reported in figure 10a and 10b. The R_g results are instead plotted in images 11a and 11b. Finally the autocorrelation function graphs are reported in 12a and 12b. Although the RMSD seems to be less variable then the R_g profile, in reality the interval is approximately large the same. By looking to the RMSD profile, it can be argued that the final distension is obtained at the 50th frame, which corresponds to the $50 * 50000 = 2,500,000$ step. As stated in 4.2, these evaluations

were done in absence of specific interactions between the beads of the polymer. As a consequence, the equilibrium that was thought to be obtained was due to the reaching of the maximal and most stable extension of the chain.

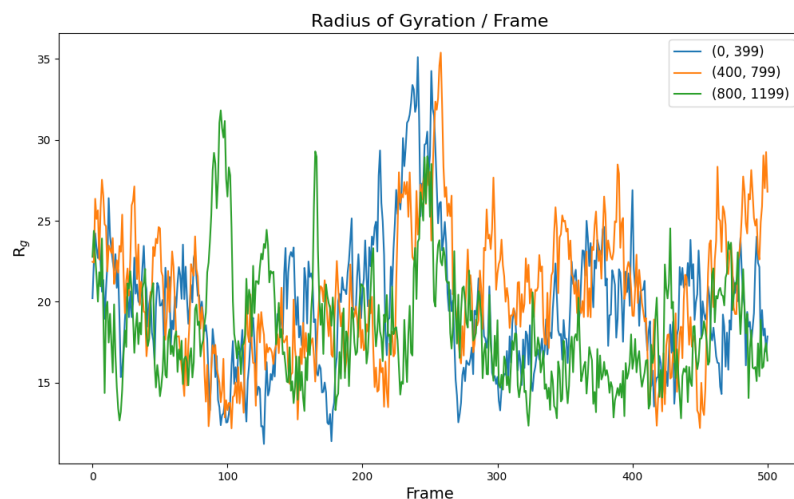


(a) Graph representing the **RMSD** of the chains pertaining to the first replicate. As it is possible to see from the legend, the first, the second and the third chain are represented respectively in blue orange and green.

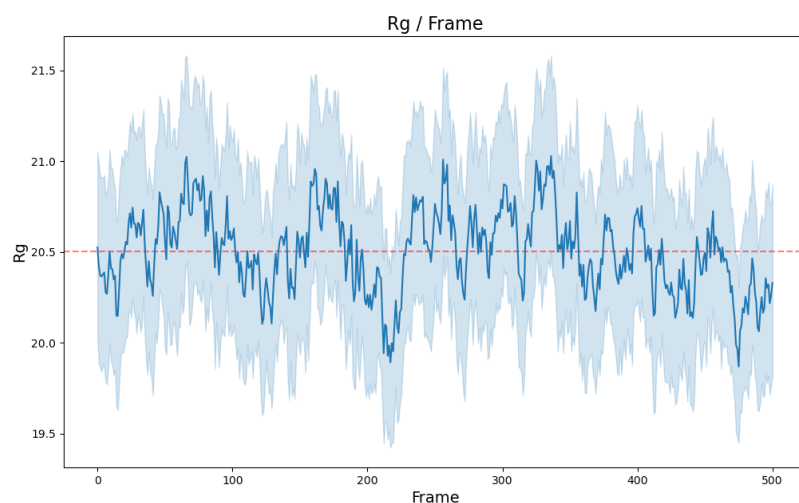


(b) Figure representing the collective behaviour of all the chains of all the 100 replicates. It is possible to observe a plateau at approximately 50×50000 steps. The red dashed line represents the average value.

Figure 10. RMSD profiles

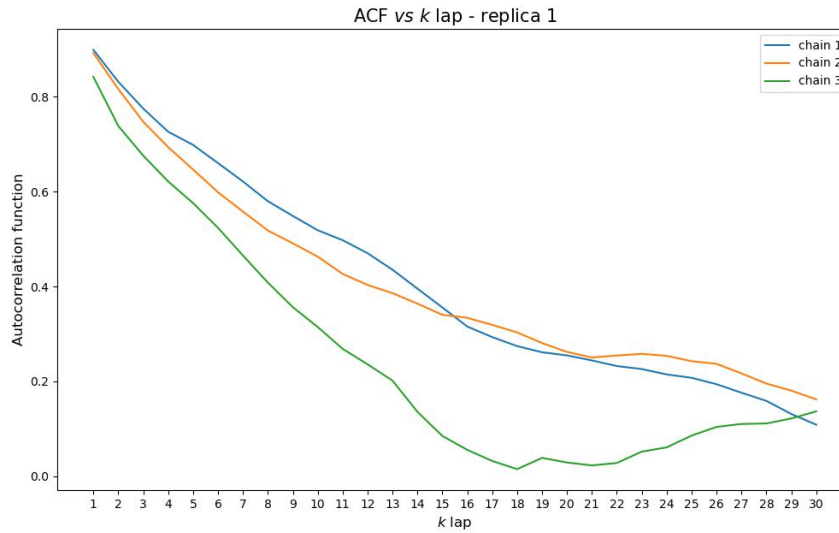


(a) Graph representing the R_g of the chains pertaining to the first replicate. As it is possible to see from the legend, the first, the second and the third chain are represented respectively in blue orange and green.

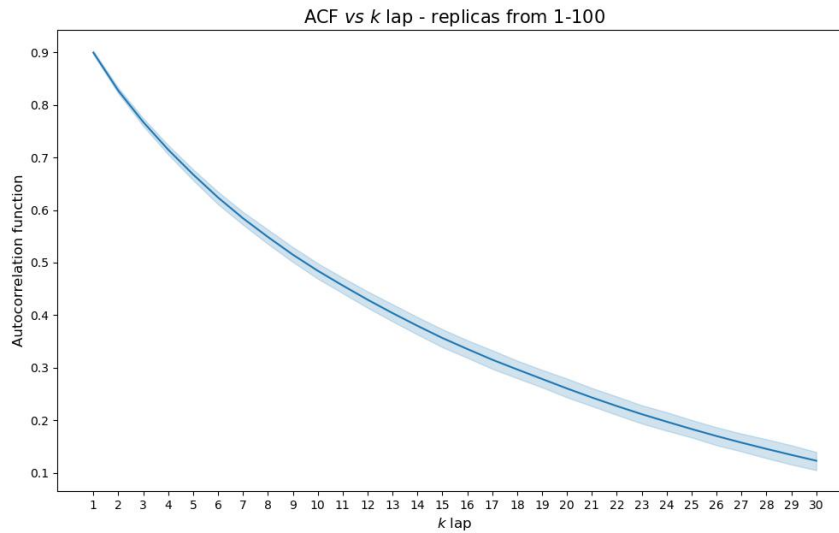


(b) Figure representing the collective behavior of all the chains of all the 100 replicates. The red dashed line represents the average value.

Figure 11. R_g profiles



(a) Graph representing the **autocorrelation function** of the chains pertaining to the first replicate. As it is possible to see from the legend, the first, the second and the third chain are represented respectively in blue orange and green.



(b) Figure representing the collective behavior of all the chains of all the 100 replicates. All the chains of all the replicates were considered independent from each other and taken as singular examples.

Figure 12. Autocorrelation function results

5.4 Results from model selection

The steps described in algorithm 3 were partially followed. A number of configurations for the model were tested, one after the other (loops). To assess each time the best version of a loop, the variations, obtained by tuning the value associated to the same variable, were confronted in terms of the produced SCC quantities. In table 4, the performed loops are listed, and only the characterizations giving the best results are shown.

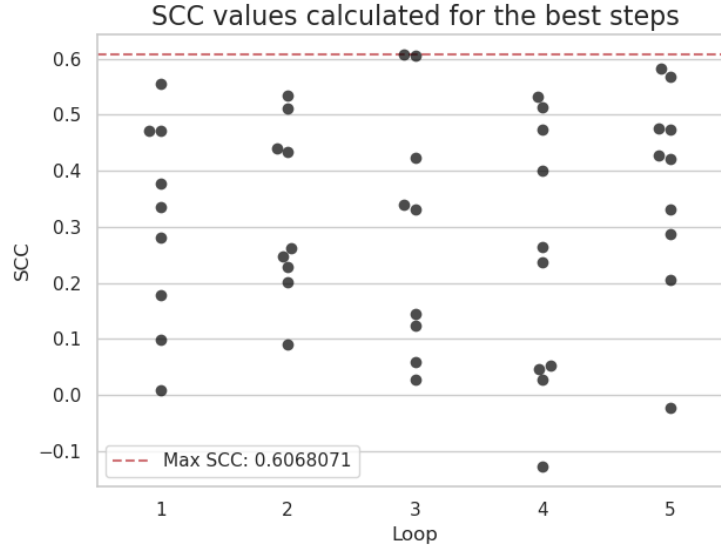


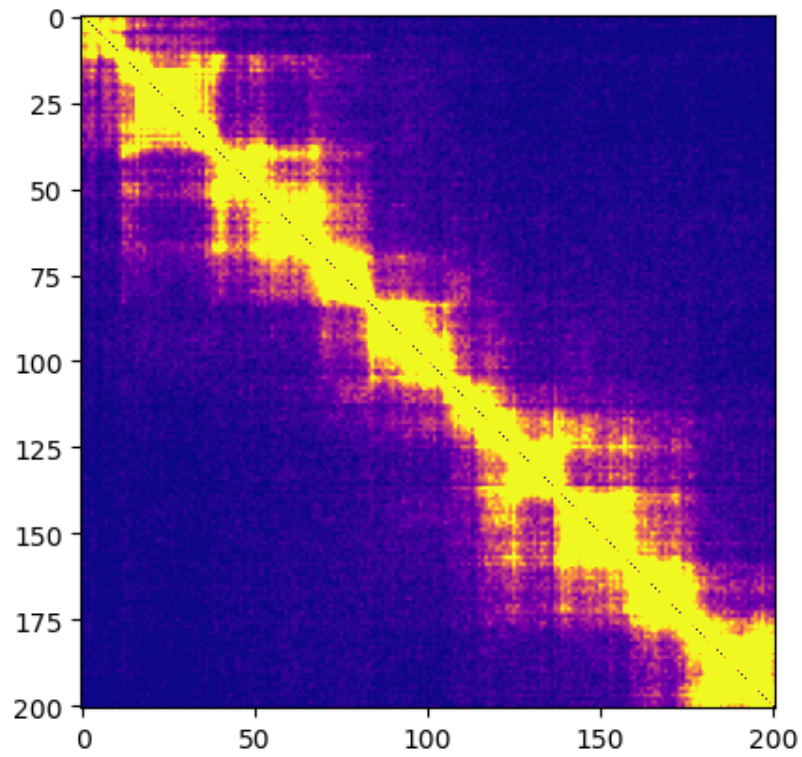
Figure 13. The results obtained from the best steps of the loops. The information reported are also in table 4.

Loop	Model	SCC
Loop 1	$E_{33} = 0.7$	0.5549067
Loop 2	$E_{33} = 0.7, E_{22} = 0.8$	0.5352417
Loop 3	$E_{33} = 0.7, E_{22} = 0.8, E_{44} = 0.6$	0.6068071
Loop 4 (removal)	$E_{22} = 0.8, E_{44} = 0.6$	0.5324294
Loop 5	$E_{33} = 0.7, E_{22} = 0.8, E_{44} = 0.6, E_{11} = 0.6$	0.5822261

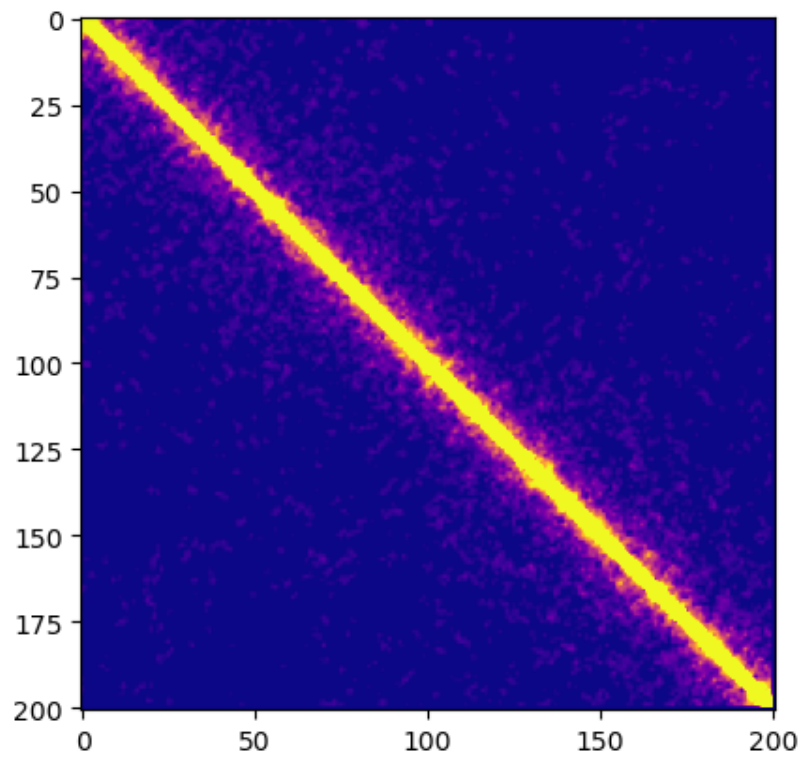
Table 4. Results from the first tested loops. Interactions between beads of the same type are written in the form " $E_{nm} = value$ "

In contrast with the initial guesses, state 1 does not appear in the best model found until now. This meant that the "greedy" approach, described in 4.5, should be discarded, and more combinations should be considered in order to find the best combination of parameters. The data in the table could also be visualized in figure 13.

Although this process has not been finished yet, I decided to produce the contact map of the best model tested ($E_{33} = 0.7$, $E_{22} = 0.8$, $E_{44} = 0.6$), to make the first considerations that would aid in the next future (figure 14). Despite the evident difference existing among the two matrices, it is possible to see larger densities of points where the original matrix signals the presence of TADs, and, in general on the boundaries of the squared shapes. The value obtained of SCC, which is a Pearson Correlation coefficient weighted with factors that depend on the population of the distance bends, is quite high, and indicate a moderate correlation. Expectedly, the contacts measured by the simulated model were very lower with respect to those of the original one.



(a) Original matrix



(b) Simulated matrix

Figure 14. Representation of the real and the simulated matrices.

6 CONCLUSIONS

To conclude, several simulations of 2 Mb chromatin regions containing the ANPEP locus were performed. The length was considered as composed by beads with fixed dimension (chapter 4.2), which were assigned to one of the state presented in the *ChromHMM* results (chapters and 2.3 and 4.3) whose input were ATAC-seq and CTCF Chip-Seq data (chapter 4.1). After the tuning of the parameters associated to the interaction potentials generated between beads of the same type, very interesting correlation coefficients between the simulated matrices and the true experimental matrices were found. The maps were compared making use of the SCC coefficient, however, another possible way to see the differences would be to compute the Spearman correlation coefficient. Ideally, it would be interesting to see if there are cases where the two metrics produce different results, and to understand which of them is better in what situations. As a future perspective, it could be considered the extension of the analysis towards new cell-types and/or new *loci*. In particular, we would be interested in investigating the MYC, SOX9, ITG45, MSX2, NT5E genes, and the GM12878 cell-type. Also, the tuning process for the parameters could be improved and automated better. To allow a better comparison with the already present models, the results obtained with the model could be compared to those resulting from other very interesting simulation softwares, such as *Origami* and *Hip-Hop*^{23,24}.

7 Glossary

TAD	Topological domains
bp	Base Pairs
Bead	The complex formed by the DNA and the histone proteins
TSS	Transcriptional Starting Site
TES	Transcriptional Ending site
FS	Fine Scale
CG	Coarse Grained. It is used to refer to the model with 5,000 bp beads
k. l.	Kuhn length
mon.	Monomer
PL	persistence length
LJ	Lennard-Jones
FENE potential	Finite Extensible Nonlinear Elastic potential
RMSD	Root Mean Square Deviation
Rg	Radius of Gyration
Map	The term is used as a synonym for the term matrix

References

1. Paro, P. D. R., Grossniklaus, P. D. U., Santoro, D. R. & Wutz, P. D. A. Biology of Chromatin. In *Introduction to Epigenetics [Internet]*, DOI: [10.1007/978-3-030-68670-3_1](https://doi.org/10.1007/978-3-030-68670-3_1) (Springer, 2021).
2. Liao, B. B. *et al.* Adaptive Chromatin Remodeling Drives Glioblastoma Stem Cell Plasticity and Drug Tolerance. *Cell Stem Cell* **20**, 233–246.e7, DOI: [10.1016/j.stem.2016.11.003](https://doi.org/10.1016/j.stem.2016.11.003) (2017).
3. Shareef, S. J. *et al.* Extended-representation bisulfite sequencing of gene regulatory elements in multiplexed samples and single cells. *Nat Biotechnol* **39**, 1086–1094, DOI: [10.1038/s41587-021-00910-x](https://doi.org/10.1038/s41587-021-00910-x) (2021).
4. Kim, S., Yu, N.-K. & Kaang, B.-K. CTCF as a multifunctional protein in genome regulation and gene expression. *Exp Mol Med* **47**, e166–e166, DOI: [10.1038/emmm.2015.33](https://doi.org/10.1038/emmm.2015.33) (2015).
5. Hsieh, T.-H. S. *et al.* Enhancer–promoter interactions and transcription are largely maintained upon acute loss of CTCF, cohesin, WAPL or YY1. *Nat Genet* **54**, 1919–1932, DOI: [10.1038/s41588-022-01223-8](https://doi.org/10.1038/s41588-022-01223-8) (2022).
6. Ou, H. D. *et al.* ChromEMT: Visualizing 3D chromatin structure and compaction in interphase and mitotic cells. *Science* **357**, eaag0025, DOI: [10.1126/science.aag0025](https://doi.org/10.1126/science.aag0025) (2017).
7. Robinson, P. J. J., Fairall, L., Huynh, V. A. T. & Rhodes, D. EM measurements define the dimensions of the “30-nm” chromatin fiber: Evidence for a compact, interdigitated structure. *Proc. Natl. Acad. Sci.* **103**, 6506–6511, DOI: [10.1073/pnas.0601212103](https://doi.org/10.1073/pnas.0601212103) (2006).
8. Buenrostro, J. D., Giresi, P. G., Zaba, L. C., Chang, H. Y. & Greenleaf, W. J. Transposition of native chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-binding proteins and nucleosome position. *Nat Methods* **10**, 1213–1218, DOI: [10.1038/nmeth.2688](https://doi.org/10.1038/nmeth.2688) (2013).
9. Grandi, F. C., Modi, H., Kampman, L. & Corces, M. R. Chromatin accessibility profiling by ATAC-seq. *Nat Protoc* **17**, 1518–1552, DOI: [10.1038/s41596-022-00692-9](https://doi.org/10.1038/s41596-022-00692-9) (2022).
10. Zhang, Y. *et al.* Model-based Analysis of ChIP-Seq (MACS). *Genome Biol.* **9**, R137, DOI: [10.1186/gb-2008-9-9-r137](https://doi.org/10.1186/gb-2008-9-9-r137) (2008).
11. Shah, A. Chromatin immunoprecipitation sequencing (ChIP-Seq) on the SOLiD™ system. *Nat Methods* **6**, ii–iii, DOI: [10.1038/nmeth.f.247](https://doi.org/10.1038/nmeth.f.247) (2009).
12. ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57–74, DOI: [10.1038/nature11247](https://doi.org/10.1038/nature11247) (2012).
13. ATAC-Seq Services - End-to-End Open Chromatin Analysis Service. <https://www.activemotif.com/catalog/1233/atac-seq-service>.
14. Ernst, J. & Kellis, M. Chromatin-state discovery and genome annotation with ChromHMM. *Nat Protoc* **12**, 2478–2492, DOI: [10.1038/nprot.2017.124](https://doi.org/10.1038/nprot.2017.124) (2017).
15. Chilled House Vibes. Learning Chromatin States from ChIP-seq Data: ChromHMM - Luca Pinello (2015).
16. Lajoie, B. R., Dekker, J. & Kaplan, N. The Hitchhiker’s Guide to Hi-C Analysis: Practical guidelines. *Methods* **72**, 65–75, DOI: [10.1016/j.ymeth.2014.10.031](https://doi.org/10.1016/j.ymeth.2014.10.031) (2015).
17. Di Stefano, M., Paulsen, J., Lien, T. G., Hovig, E. & Micheletti, C. Hi-C-constrained physical models of human chromosomes recover functionally-related properties of genome organization. *Sci Rep* **6**, 35985, DOI: [10.1038/srep35985](https://doi.org/10.1038/srep35985) (2016).
18. Halverson, J. D., Smrek, J., Kremer, K. & Grosberg, A. Y. From a melt of rings to chromosome territories: The role of topological constraints in genome folding. *Rep Prog Phys* **77**, 022601, DOI: [10.1088/0034-4885/77/2/022601](https://doi.org/10.1088/0034-4885/77/2/022601) (2014).
19. Lieberman-Aiden, E. *et al.* Comprehensive mapping of long range interactions reveals folding principles of the human genome. *Science* **326**, 289–293, DOI: [10.1126/science.1181369](https://doi.org/10.1126/science.1181369) (2009).
20. Razin, S., Ulianov, S. & Gavrilov, A. 3D Genomics. *Mol. Biol.* **53**, 802–812, DOI: [10.1134/S0026893319060153](https://doi.org/10.1134/S0026893319060153) (2019).
21. Kremer, K. & Grest, G. S. Dynamics of entangled linear polymer melts: A molecular-dynamics simulation. *J. Chem. Phys.* **92**, 5057–5086, DOI: [10.1063/1.458541](https://doi.org/10.1063/1.458541) (1990).
22. Grosberg, A. J., Chochlov, A. R. & de Gennes, P.-G. *Giant Molecules: Here, There, and Everywhere* (World Scientific, New Jersey, 2011), 2. ed edn.
23. Buckle, A., Brackley, C. A., Boyle, S., Marenduzzo, D. & Gilbert, N. Polymer Simulations of Heteromorphic Chromatin Predict the 3D Folding of Complex Genomic Loci. *Mol. Cell* **72**, 786–797.e11, DOI: [10.1016/j.molcel.2018.09.016](https://doi.org/10.1016/j.molcel.2018.09.016) (2018).
24. Tan, J. *et al.* Cell-type-specific prediction of 3D chromatin organization enables high-throughput in silico genetic screening. *Nat Biotechnol* **41**, 1140–1150, DOI: [10.1038/s41587-022-01612-8](https://doi.org/10.1038/s41587-022-01612-8) (2023).

25. Steve Plimpton, Aidan Thompson, Stan Moore, Axel Kohlmeyer & Richard Berger. LAMMPS. <https://www.lammps.org/#gsc.tab=0>.
26. Thompson, A. P. *et al.* LAMMPS - a flexible simulation tool for particle-based materials modeling at the atomic, meso, and continuum scales. *Comput. Phys. Commun.* **271**, 108171, DOI: [10.1016/j.cpc.2021.108171](https://doi.org/10.1016/j.cpc.2021.108171) (2022).
27. ATAC-seq (unreplicated) – ENCODE. <https://www.encodeproject.org/pipelines/ENCPL344QWT/>.
28. Michael Cherry, Jason Buenrostro, Alicia Schep & Will Greenleaf. ATACSeq Pipeline. https://docs.google.com/document/d/1f0Cm4vRyDQDu0bMehHD7P7KOMxTOP-HiNoIvL1VcBt8/edit?usp=sharing&usp=embed_facebook.
29. Yan, F., Powell, D. R., Curtis, D. J. & Wong, N. C. From reads to insight: A hitchhiker's guide to ATAC-seq data analysis. *Genome Biol.* **21**, 22, DOI: [10.1186/s13059-020-1929-3](https://doi.org/10.1186/s13059-020-1929-3) (2020).
30. Transcription Factor ChIP-seq Data Standards and Processing Pipeline – ENCODE. https://www.encodeproject.org/chip-seq/transcription_factor/.
31. Homo sapiens genome assembly GRCh38 - NCBI - NLM. https://www.ncbi.nlm.nih.gov/datasets/genome/GCF_000001405.26/.
32. UCSC Genome Browser Home. <https://genome.ucsc.edu/>.
33. Ehler, E., Babiychuk, E. & Draeger, A. Human foetal lung (IMR-90) cells: Myofibroblasts with smooth muscle-like contractile properties. *Cell Motil.* **34**, 288–298, DOI: [10.1002/\(SICI\)1097-0169\(1996\)34:4<288::AID-CM4>3.0.CO;2-4](https://doi.org/10.1002/(SICI)1097-0169(1996)34:4<288::AID-CM4>3.0.CO;2-4) (1996).
34. Ingram, S. P. *et al.* Hi-C implementation of genome structure for in silico models of radiation-induced DNA damage. *PLOS Comput. Biol.* **16**, e1008476, DOI: [10.1371/journal.pcbi.1008476](https://doi.org/10.1371/journal.pcbi.1008476) (2020).
35. Maiser, A. *et al.* Super-resolution in situ analysis of active ribosomal DNA chromatin organization in the nucleolus. *Sci Rep* **10**, 7462, DOI: [10.1038/s41598-020-64589-x](https://doi.org/10.1038/s41598-020-64589-x) (2020).
36. Golkaram, M., Jang, J., Hellander, S., Kosik, K. S. & Petzold, L. R. The Role of Chromatin Density in Cell Population Heterogeneity during Stem Cell Differentiation. *Sci Rep* **7**, 13307, DOI: [10.1038/s41598-017-13731-3](https://doi.org/10.1038/s41598-017-13731-3) (2017).
37. Robinson, J. T., Thorvaldsdóttir, H., Turner, D. & Mesirov, J. P. Igv.js: An embeddable JavaScript implementation of the Integrative Genomics Viewer (IGV). Preprint, Bioinformatics (2020). DOI: [10.1101/2020.05.03.075499](https://doi.org/10.1101/2020.05.03.075499).
38. Gowers, R. *et al.* MDAnalysis: A Python Package for the Rapid Analysis of Molecular Dynamics Simulations. In *Python in Science Conference*, 98–105, DOI: [10.25080/Majora-629e541a-00e](https://doi.org/10.25080/Majora-629e541a-00e) (Austin, Texas, 2016).
39. Tuckerman, M. E. *Statistical Mechanics: Theory and Molecular Simulation*. Oxford Graduate Texts (Oxford Univ. Press, Oxford, 2015), reprinted (with corr.) edn.
40. Suma, A., Di Stefano, M. & Micheletti, C. Electric-Field-Driven Trapping of Polyelectrolytes in Needle-like Backfolded States. *Macromolecules* **51**, 4462–4470, DOI: [10.1021/acs.macromol.8b00019](https://doi.org/10.1021/acs.macromol.8b00019) (2018).
41. Lin, D., Sanders, J. & Noble, W. S. HiCRep.py: Fast comparison of Hi-C contact matrices in Python. *Bioinformatics* **37**, 2996–2997, DOI: [10.1093/bioinformatics/btab097](https://doi.org/10.1093/bioinformatics/btab097) (2021).
42. Yang, T. *et al.* HiCRep: Assessing the reproducibility of Hi-C data using a stratum-adjusted correlation coefficient. *Genome Res* **27**, 1939–1949, DOI: [10.1101/gr.220640.117](https://doi.org/10.1101/gr.220640.117) (2017).
43. Robinson, J. T. Integrative genomics viewer. *C O Rresp O N N Ce* **29**, 3 (2011).