

METHODOLOGY ARTICLE

Open Access



Inferring chromosome radial organization from Hi-C data

Priyojit Das¹, Tongye Shen² and Rachel Patton McCord^{2*} 

*Correspondence:
rmccord@utk.edu

² Department of Biochemistry and Cellular and Molecular Biology, University of Tennessee, Knoxville, TN 37996, USA
Full list of author information is available at the end of the article

Abstract

Background: The nonrandom radial organization of eukaryotic chromosome territories (CTs) inside the nucleus plays an important role in nuclear functional compartmentalization. Increasingly, chromosome conformation capture (Hi-C) based approaches are being used to characterize the genome structure of many cell types and conditions. Computational methods to extract 3D arrangements of CTs from this type of pairwise contact data will thus increase our ability to analyze CT organization in a wider variety of biological situations.

Results: A number of full-scale polymer models have successfully reconstructed the 3D structure of chromosome territories from Hi-C. To supplement such methods, we explore alternative, direct, and less computationally intensive approaches to capture radial CT organization from Hi-C data. We show that we can infer relative chromosome ordering using PCA on a thresholded inter-chromosomal contact matrix. We simulate an ensemble of possible CT arrangements using a force-directed network layout algorithm and propose an approach to integrate additional chromosome properties into our predictions. Our CT radial organization predictions have a high correlation with microscopy imaging data for various cell nucleus geometries (lymphoblastoid, skin fibroblast, and breast epithelial cells), and we can capture previously documented changes in senescent and progeria cells.

Conclusions: Our analysis approaches provide rapid and modular approaches to screen for alterations in CT organization across widely available Hi-C data. We demonstrate which stages of the approach can extract meaningful information, and also describe limitations of pairwise contacts alone to predict absolute 3D positions.

Keywords: Hi-C, Chromosome territory, Chromosome radial organization, Principal component analysis, Network modeling, Nucleus shape, Gene density chromosome length

Background

The three dimensional (3D) structure of the human genome is composed of different structures at different length scales. At smaller length scales, nucleosome positions, loops, and topologically associating domains are the most salient features, followed by compartmentalization at a longer length scale [1]. At the largest scale of this genome organization, the 3D bodies of individual chromosomes arrange mostly



as discrete entities, known as chromosome territories (CTs) [2]. The arrangement of all the CTs inside the nucleus with respect to nuclear center and periphery forms the higher-order genome architecture. This CT organization is nonrandom with respect to the nucleus periphery and can play important roles in different nuclear mechanisms ranging from DNA replication and gene expression to the processing of RNA [3]. Recently, it has been shown that chromosome territorial organization can also protect genome from deleterious rearrangements during DNA damage [4]. Alterations in CT organization can also be important during cell differentiation [5] and in different disease conditions. For example, in Hutchinson–Gilford progeria syndrome cells, chr18 shifts toward the nucleus interior as compared to its position in normal proliferating fibroblast cells [6], while certain gene-rich chromosomes localize near the periphery in blebs in progeria cells [7]. Alterations in CTs can influence the likelihood of chromosomal translocations. For example, during adipogenesis, chr12 and chr16 become spatially proximal, increasing the chance of translocation between those chromosomes, which is the driving event of liposarcoma tumorigenesis [5]. Recently, researchers have shown that in breast cancer, a gain in inter-chromosomal interactions for chrX correlates with its gene expression changes [8]. Approaches to characterize CT positions can thus further our understanding of the implications of CT arrangements in health and disease.

From careful microscopic measurements over the past several decades, largely using sequence specific probes in fluorescence in situ hybridization (FISH), principles of CT organization in certain cell types have been identified. It has been observed that CTs are often organized according to one of two different distributions: either a gene density or chromosome length based pattern [9, 10]. Opposing forces of gene activity and lamina associated domain (LAD) density on different chromosomes likely contribute to these different distributions [11, 12]. In general, LADs are likely to be repressive to gene activity and occur more frequently on gene poor chromosomes, which also tend to be longer than gene rich chromosomes [13]. The proliferation rate of cells and nuclear shape have also been implicated as factors influencing CT organization. For example, proliferating cells tend to follow a gene density-based organization compared to the length-based distribution in quiescent or senescent cells [14]. Further, the spherical human lymphocyte nucleus follows a gene density based organization, which is conserved across several related species [15, 16], while the chromosome size based distribution is more prevalent in ellipsoidal fibroblast nuclei [17]. Our understanding of the relationships between factors influencing CT distribution are limited, however, by the fact that relatively few different cell types have been characterized in depth by this type of microscopic analysis.

Thorough analysis of CT positions requires not only observing their average, but also the distribution of possible positions across the cell population. The position of each chromosome varies between cells in the population, even while following certain tendencies [18]. Throughout interphase the CT positions remain stable, but change from one generation to another during mitosis as the nuclear envelope is broken down and re-established [19–22]. Improvements to microscopy experiments [23, 24] and image analysis methods make sampling this variation in CT position across the population increasingly feasible [25–27], yet such analyses of all CT positions in a large number of cells exist for only few cell types [28].

In contrast, genome wide chromosome conformation capture (Hi-C) approaches are being applied to characterize the 3D genome structure of rapidly increasing numbers of cell types and conditions for the past decade [1, 29–31]. This Hi-C technique and its variants capture a snapshot of pairwise chromosomal interactions ranging from specific enhancer-promoter interactions [32] to large scale inter-chromosomal interactions. Due to the effect of noise and high variability, the inter-chromosomal interactions have received less attention compared to their intra-chromosomal counterparts. In a past few years, with the improvement of experimental protocols, the effect of noise on inter-chromosomal contacts has been reduced and studies based on relevant inter-chromosomal interactions have started to emerge [33–36]. For example, genes corresponding to olfactory receptors from different chromosomes form specific inter-chromosomal contacts in mouse olfactory sensory neurons which strengthen upon differentiation from a progenitor cell [37]. In another set of studies, by analyzing Hi-C inter-chromosomal contacts obtained from different malignant diseases, researchers have identified several novel chromosomal rearrangements [38, 39]. Though Hi-C does not directly capture the radial position of the CTs, approaches that infer 3D chromosome positioning information from Hi-C contact data provide a valuable supplement to microscopic data, greatly increasing the number of cell types for which CT positions can be analyzed. In this study, we explore a set of Hi-C analysis approaches focused on rapid and efficient prediction of CT radial organization, ranging from very simple direct calculations on the Hi-C contact matrix to a network model tuned by additional chromosome property information.

Many approaches have been developed to reconstruct the 3D folding of individual chromosomes and the genome from Hi-C data at different resolutions using restraint and polymer physics based approaches [40, 41]. Some models focus on detailed structures of local regions of chromatin rather than the whole genome, and for others, the primary focus is often to understand the mechanistic principles underlying the organization of the interphase and metaphase genome rather than a prediction of CT arrangement [42–48]. But, some of the approaches have also explored the radial arrangement of CTs in their 3D models. An early restraint-based model of the 3D genome based on tethered chromosome conformation capture (TCC) data predicted CT positions that agreed with the major principles of lymphoblast genome organization characterized by microscopy [42]. In another work, Hi-C maps were probabilistically deconvoluted into a population of single cell structures, and the averaged radial position of the CTs predicted from those structures matched fairly well with microscopic measurements from a single cell type [49]. A series of coarse-grained polymer simulation based studies have been performed to characterize non-random organization of the chromosomes using gene activity and random and biological looping constraints [50–52]. In a recent polymer modeling based study, the researchers used Hi-C derived properties and a chromatin state based energy function to study the principles of the spatial as well as radial genome organization [53].

Because numerous arrangements can potentially be consistent with a set of Hi-C contacts, in many cases, whole genome 3D models generated from Hi-C also have to take into account external information in order to increase the accuracy of CT positioning predictions. For example, Stevens et al. combine imaging and Hi-C contacts on single cells to orient their 3D chromosome models [54]. Similarly, the Chrom3D algorithm

combines LAD data along with Hi-C data to capture spatial and radial organization of the chromosomes [55]. Di Stefano et al. used a steered molecular dynamics based approach to reconstruct diploid genome organization from Hi-C data [56]. Using significant Hi-C interactions as the constraints, the modeling technique was able to capture the preferential nuclear position of different genomic regions based on their gene density, lamina association and epigenetic marks.

To supplement this landscape of approaches to predict 3D genome conformations from Hi-C data, we have several specific goals in this study. We describe and test direct, computationally non-intensive analysis approaches that have the focused aim of inferring CT radial positions from Hi-C data rather than a full model of 3D chromosome folding. Specifically, we demonstrate that radial organization patterns can be inferred from PCA analysis of thresholded inter-chromosomal contact matrices and show the utility of a force-directed graph layout algorithm to infer the average and variation around the average CT positions. These approaches thus do not require the computational resources necessary to calculate high resolution polymer models, and could be used to screen for potentially important differences in CT organization across a wide range of cell types and conditions without creating detailed 3D models in each case. We further evaluate the strengths and limitations of Hi-C contact data when it is used toward the goal of inferring large scale 3D positions of chromosomes, and where additional reference information needs to be added to the contacts to generate reliable radial positioning information. We finally evaluate the different stages of our approach on a variety of cell types and conditions, comparing to a variety of published microscopic data and predicting additional details of CT organization where limited microscopy data exists.

Results

Thresholding contacts extracts meaningful chromosomal interaction patterns from Hi-C data

Hi-C experiments capture spatial genome organization by measuring the interaction frequency between different genomic fragments, yielding information about both intra-chromosomal and inter-chromosomal interactions. Inter-chromosomal interactions generally occur much less frequently than intra-chromosomal interactions, and true interactions are mixed in with noise that arises from random background ligation [30, 57]. Despite this inherent noise and sometimes low signal, the interactions between chromosomes also contain information that reflects the radial organization of the chromosome territories inside the nucleus. But, excluding contacts that may primarily reflect the background is important to prevent those contacts from masking the true signal. In order to extract strong interactions that are more likely to distinguish radial chromosome positions between cell types, we applied a thresholding technique to the genome-wide contact matrix (see “Methods” section). The thresholding cutoff h_{cut} value was calculated by taking a certain percentile of all the genome-wide Hi-C interactions and then the interactions greater than this cutoff limit are considered as strong interactions. To determine the desired cutoff value, we compared a Hi-C matrix from a specific cell type with a corresponding simulated random ligation matrix (see “Methods” section) for different values of h_{cut} ranging from the 5th to 95th percentile. We examined two

different cell types with different nuclear shapes (Additional file 1: Table 1), since the nuclear shape has been observed to correlate to some extent with the non-random radial organization of the CTs. The human blood lymphoblastoid cell, GM12878, has a spherical nucleus and follows a gene density based radial CT organization [9]. On the other hand, BJ1-hTERT human skin fibroblast cell has an ellipsoidal nucleus, which shows a chromosome length based CT organization [17]. The whole genome Hi-C contact matrices were obtained from Sanders et al. [58].

For each h_{cut} value, the number of chromosomal contact bin pairs passing the threshold were summed between each pair of chromosomes to a single bin (Eqn. 4).

$$CHR_{ij}^{STRONG} = \text{total number of strong interaction bins between chr } i \text{ and } j \quad (1)$$

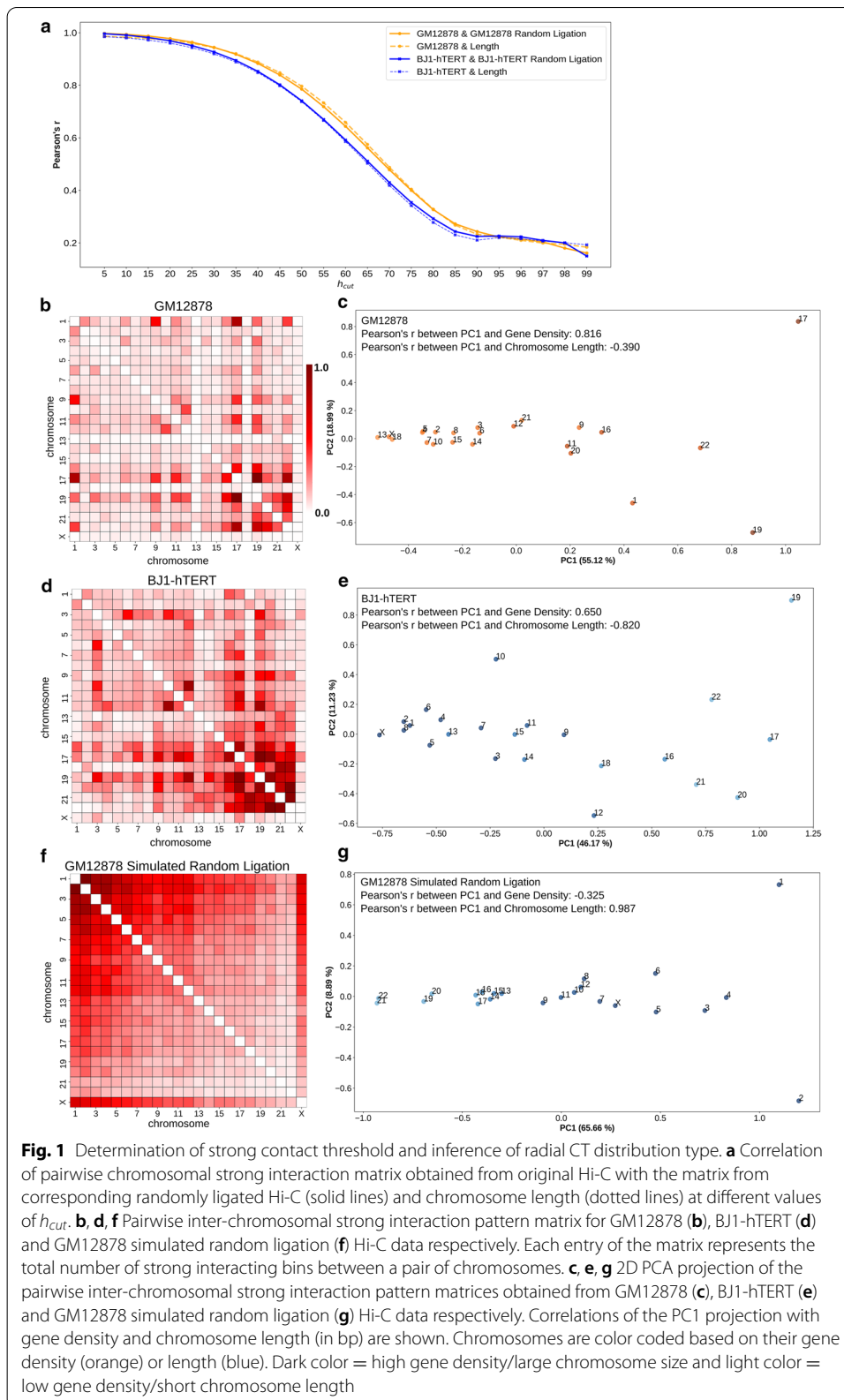
Then, Pearson's correlation was applied to compare the whole chromosome pair interaction sums obtained from the original and corresponding random ligation Hi-C data. From Fig. 1a, it can be seen that with the increase in h_{cut} values from the 45th to the 85th percentile, the pairwise strong chromosomal interaction similarity between random and real data decreased rapidly and reached a stable value around the 90th percentile. Similarly, we compared the strong chromosomal interaction sums with the pairwise product of the chromosome lengths to measure how much the number of interactions was primarily driven by chromosome length (i.e. two large chromosomes will have more interactions at random overall than two small chromosomes). This analysis produced a similar correlation trend as the random ligation effect comparison: strong interaction sums are no longer primarily explained by chromosome length at 90th–95th percentile h_{cut} (Fig. 1a). Based on these two comparison results for both GM12878 and BJ1-hTERT, we chose the 95th percentile as the final value of h_{cut} that leads to a minimized effect of chromosome length and random ligation for both cell types. While this optimal value was similar for two different datasets we considered, we note that Hi-C library complexity, read depth, and cis/trans ratio could affect the most appropriate h_{cut} value. As an alternative to this thresholding approach, we also explored the FitHiC [59] algorithm to extract significant chromosomal interactions from the Hi-C data. The analysis results obtained using those interactions are discussed in the following subsection.

Radial chromosome ordering can be inferred from PCA on inter-chromosomal strong interaction pattern matrix

The pairwise chromosomal strong interaction pattern is not only able to distinguish the true biological interaction pattern from the random ligation, but also reveals distinct patterns specific to each cell type. By looking at the inter-chromosomal component of the pairwise strong interaction patterns (Eq. 5) for GM12878 and

$$CHR_{ij}^{TRANS} = \begin{cases} CHR_{ij}^{STRONG} & \text{if } i \neq j \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

BJ1-hTERT represented in Fig. 1b, d respectively, we can clearly see that these two different cell types have distinct patterns—in BJ1-hTERT, the smaller chromosomes have higher strong inter-chromosomal interactions among themselves, whereas in GM12878 that pattern is dispersed. In order to capture the major interaction trends between



chromosomes from those distinct patterns, we applied principal component analysis (PCA) to these pairwise strong interaction matrices. PCA has been frequently applied to contacts within or between chromosomes to detect the spatial segregation of A/B compartments [29, 60–62]. Further, the A/B compartment status of bins within a chromosome has been found to correlate with their lamina association and radial positioning, and therefore has been used in models predicting the lamin associations of domains within a chromosome [63]. Here, by contrast, we are not detecting the compartment status of regions within a chromosome but instead applying PCA to the pattern of thresholded pairwise interactions between whole chromosome territories and then examining the projection of chromosomes onto the first two principal components. We find that PCA on the pattern of strong contacts between all pairs of chromosomes can detect the spatial segregation and relative ordering of chromosome territories. (Here we note that since the Hi-C matrix represents the average of the two copies of each chromosome, we use chromosome (chr) and chromosome territory (CT) interchangeably, noting that the CT positions will represent the average of the chromosome locations of each homolog across the cell population). Figure 1c shows the 2D PCA projection of the pairwise inter-chromosomal strong interaction pattern matrix for the GM12878 cell. From this figure, it can be seen that chromosome 17, 19 and 22 and chromosome 13, 18 and X are on the opposite ends along the PC1 axis due to the high dissimilarity in their inter-chromosomal interaction patterns. In addition to that, this separation also correlates highly with the gene density based distribution of the chromosomes, as gene-rich chromosomes (e.g., chr19, chr17) are on the right extreme and gene-poor chromosomes (e.g., chr18, chr13) on the other extreme. This ordering thus corresponds to previous reports showing that lymphoblast cells with spherical nuclei tend to position their chromosomes in the nucleus in a gene density associated pattern [9]. Indeed, the PC1 values of the chromosomes show a higher absolute value correlation with their gene density than with their chromosome length. Next, we analyzed the pairwise inter-chromosomal significant interaction pattern matrix obtained using FitHiC for the GM12878 Hi-C data. We again see a higher PC1 correlation with gene density than chromosome length, but the correlation is much weaker (Additional file 1: Fig. 1b). This is not surprising given that the interaction pattern in the FitHiC chromosome pair interaction map is less distinct and more uniform across chromosomes (Additional file 1: Fig. 1a). Therefore, we choose to proceed with the strong interaction thresholding approach for our remaining analyses to detect the inherent radial arrangement of the CTs.

When we applied PCA to the BJ1-hTERT pairwise inter-chromosomal strong interaction pattern matrix, PC1 showed a different separation, where chromosomes are ordered from left to right roughly according to decreasing length (Fig. 1e). The strong correlation between this ordering and chromosome length matches previous observations from fibroblast nuclei [17]. There is still some gene density correlation with this pattern, likely reflecting the complex picture of the radial CT organization in fibroblast nuclei. It has indeed been reported that while chromosomes are generally radially positioned by length in fibroblast nuclei, CT18 (gene poor and short) is still nearer to the nuclear envelope than CT19 (gene rich and short). When we applied PCA to the simulated random-ligation matrices generated from the GM12878 and BJ1-hTERT, we find that chromosomes are ordered along PC1

in a strong length-based distribution (Fig. 1g and Additional file 1: Fig. 2d). This is expected, since by chance longer chromosomes will have more random strong interactions than shorter chromosomes, as is evident in Fig. 1f.

The above results suggest that a simple application of PCA to the pairwise inter-chromosomal strong pattern matrix can reveal chromosome spatial distribution type, while also showing the limitations of this direct use of Hi-C contacts alone. Though PC1 can infer the radial CT distribution type, the direction of the ordering (whether inwards to outwards or outwards to inwards) is arbitrary and cannot be inferred based only this PCA result. Further, the random ligation result provides a caution that a length-based radial distribution would be detected even when no distinct interaction patterns are present.

To further explore the utility of such PCA ordering of CT positions, we next applied this analysis technique to Hi-C data from conditions which have been previously shown to exhibit large-scale rearrangement of chromosome territories. Specifically, in the premature aging disease Hutchinson–Gilford Progeria syndrome, it has been shown that chr18 moves to the nuclear interior compared to its peripheral location in normal proliferating fibroblasts, while chr10 shows the opposite trend [6, 64]. To check whether our analysis technique can also detect these CT reorganizations, we analyzed WI38-hTERT proliferating fibroblast Hi-C data [65] and Hi-C data from progeria patient fibroblasts at passage 19, when the cells are approaching senescence [66]. Figure 2a, b show the 2D PCA projections of the pairwise inter-chromosomal strong interaction pattern matrices obtained from the proliferating fibroblast and the progeria cells respectively. As mentioned above, while these projections reflect the underlying CT organization, inferring the directionality requires additional experimental information. There is evidence that chrX does not change its peripheral position in progeria cells compared to normal proliferating fibroblasts [6]. Indeed, in our analyses, chrX is positioned at a far extreme of PC1 in both progeria and proliferating fibroblasts, so we assigned the end of the PC1 axis near chrX as the periphery and the opposite end as the nucleus center, as seen in Fig. 2a, b. Now, if we examine the positions of chr18 in these projections, we can see that chr18 occupies a more internal position in the progeria cell (Fig. 2a, b). We can visualize the relative changes in chromosome positions by ordering the chromosomes based on the PC1 value in an increasing fashion from center to periphery in a radar plot (Fig. 2d). Here, we can see the internal shift of chr18 in progeria compared to the proliferating fibroblast. We did not find any significant change in the chr10 position in our analysis. However, we noticed that chr13 in progeria also moves to interior similar to chr18, which can be found in proliferating laminopathy fibroblasts [64]. Progeria cells approaching senescence have a CT arrangement that to some extent matches other quiescent and senescent cells [6]. Therefore, next, we performed PCA ordering analysis on WI38-hTERT oncogene induced senescent cell Hi-C data [65]. When we compared the result with normal proliferating cells, we found chr10, chr13, chr18 and chrX radial rearrangements similar to the progeria cells (Fig. 2c, d bottom). Overall, this analysis shows that the strong inter-chromosomal interaction pattern from Hi-C has the potential to infer changes in the underlying CT distribution within the nucleus.

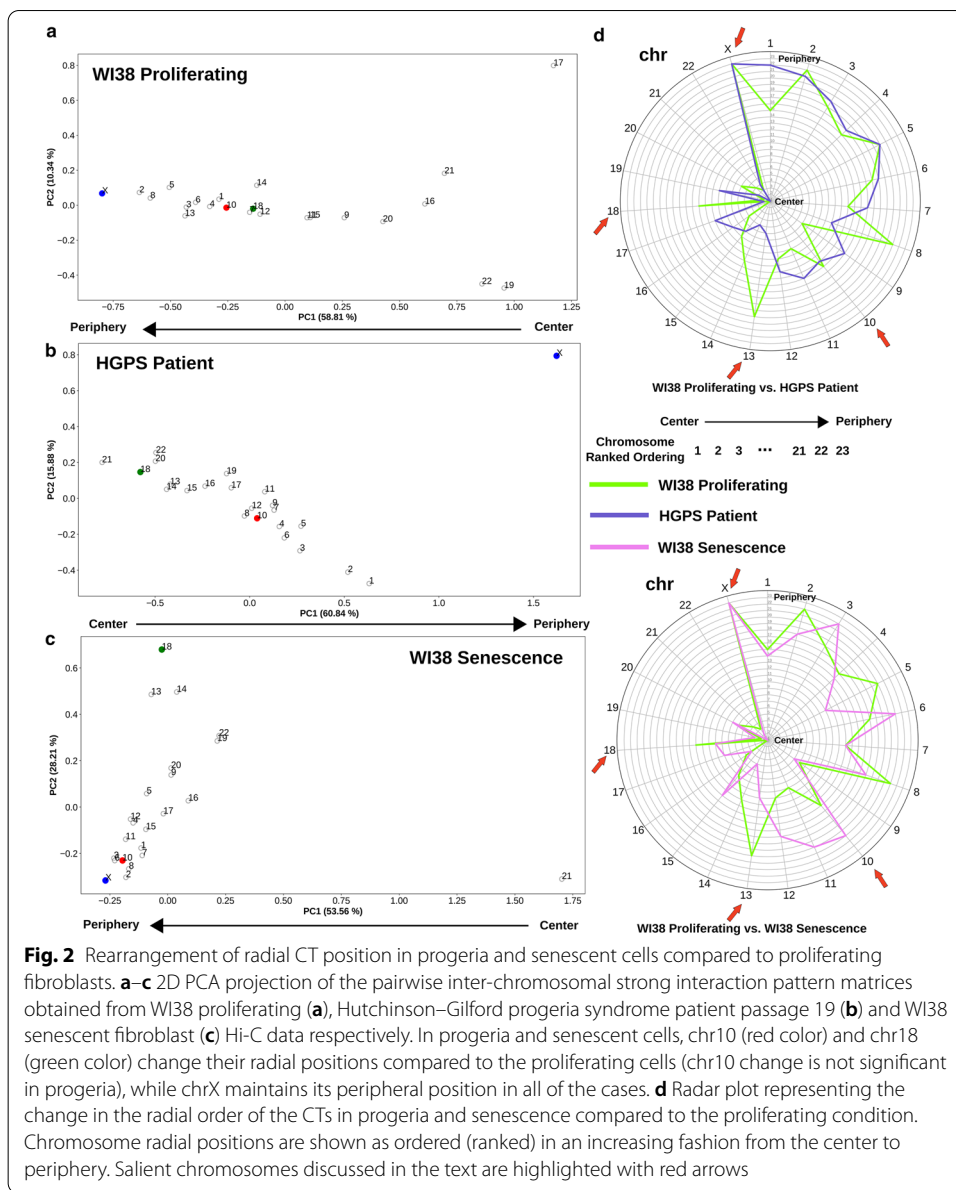


Fig. 2 Rearrangement of radial CT position in progeria and senescent cells compared to proliferating fibroblasts. **a–c** 2D PCA projection of the pairwise inter-chromosomal strong interaction pattern matrices obtained from WI38 proliferating (**a**), Hutchinson–Gilford progeria syndrome patient passage 19 (**b**) and WI38 senescent fibroblast (**c**) Hi-C data respectively. In progeria and senescent cells, chr10 (red color) and chr18 (green color) change their radial positions compared to the proliferating cells (chr10 change is not significant in progeria), while chrX maintains its peripheral position in all of the cases. **d** Radar plot representing the change in the radial order of the CTs in progeria and senescence compared to the proliferating condition. Chromosome radial positions are shown as ordered (ranked) in an increasing fashion from the center to periphery. Salient chromosomes discussed in the text are highlighted with red arrows

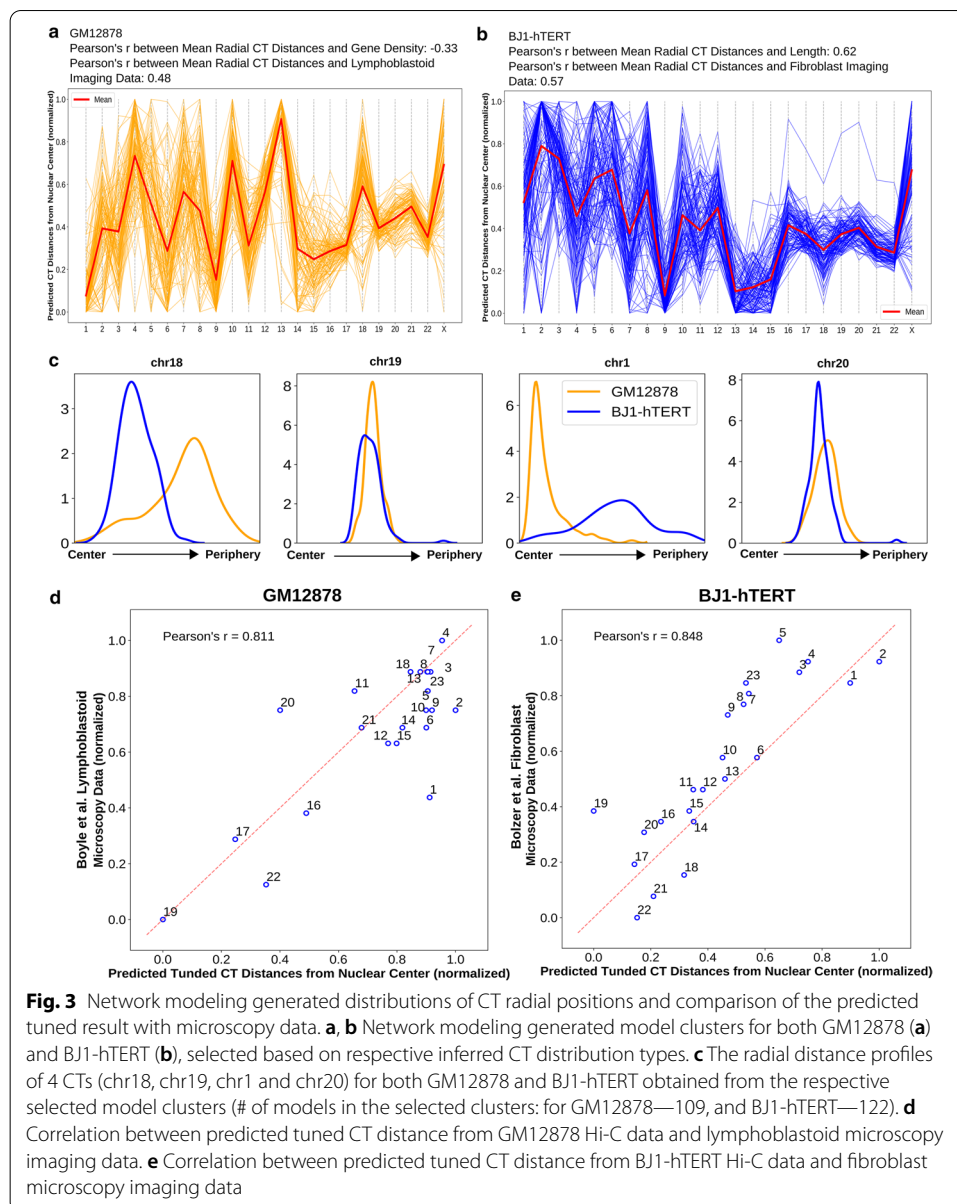
Network modeling predicts probabilistic radial CT organization

Though applying a simple PCA approach to a pairwise inter-chromosomal strong interaction pattern matrix is surprisingly powerful at detecting relative chromosome ordering in different cell types, as shown above, it is limited to predicting a single position for each chromosome, and will predict ordering patterns even for data that is actually random noise. In reality, chromosome territory positions vary substantially within individual cells, even while following certain general trends. The positioning of CTs in individual cells changes from mother cell to daughter cell during mitosis [19, 20]. Thus, we sought an approach to model not only the average position, but the variability of chromosome positions in an ensemble of arrangements derived from Hi-C data. For this, we tested a network modeling approach. We used the strong

interaction matrix obtained from the thresholding step to construct an undirected weighted graph in three-dimensional space using the 3D Fruchterman-Reingold (FR) force-based layout [67] algorithm. Each node of the graph represents a genomic region of fixed size and the weights are the contact frequencies. In such a graph layout, the nodes within each chromosome become clustered together due to high intra-chromosomal contacts while chromosomes that have higher inter-chromosomal contacts among themselves will end up close together in 3D space compared to the other chromosomes with less inter-chromosomal contacts. To model the variable CT arrangements possible within the average Hi-C data, 1000 independent runs of the FR algorithm were performed with different random initial configurations, which led to the formation of 1000 different 3D network graphs (see “Methods” section). Next, a minimum volume ellipsoid algorithm was used to fit a geometrical object around the models to serve as the nuclear periphery. From each of those models, we calculated the distance of the center of mass of each chromosome from the nucleus center (center of the geometrical fitted object).

When we examined the models, we found that they cluster into several different possible chromosome organization patterns (Additional file 1: Fig. 3). As we discovered with the PCA analysis, the Hi-C data alone cannot determine the absolute ordering from interior to periphery, so some clusters of models represent inverted patterns of organization. We also find some clusters of models that do not coincide with the length or gene density based distributions inferred from PCA CT ordering. These may represent local minima reached from certain initial conditions of the network model which do not reflect the true chromosome arrangements. Given these observations, we combine information from the previous PCA analysis with these network models to select the cluster of models for further analysis. After observing that the PCA analysis reports whether the cell type in general follows a length or gene density-based radial distribution, we use this information to select the cluster of models which has the highest absolute correlation between its mean CT distances and gene density or chromosome length (based on the inferred CT distribution type from the pairwise inter-chromosomal PCA transformation).

With the selected cluster of models (Fig. 3a, b), we next examined the predicted heterogeneity in chromosome territory positioning within each cell type (Fig. 3 and Additional file 1: Fig. 4). Figure 3c shows the radial distance profiles of four example CTs - CT1, CT18, CT19 and CT20 obtained from GM1878 and BJ1-hTERT network model clusters. Among these CTs, CT18 (length 78.1 Mbp and gene density 3.4 genes/sequenced Mbp) and CT19 (length 59.1 Mbp and gene density 23.9 genes / sequenced Mbp) have comparable DNA content but have drastically different gene densities. The network modeling result shows a trend in which, in GM12878, CT18 has a peak near the periphery and CT19 has a more internal peak. On the other hand, in BJ1-hTERT, both CT18 and CT19 peaks are located toward the nuclear center with a similar distribution. The next contrasting pair is CT1 (length 249.3 Mb and gene density 8.70 genes/sequenced Mb) and CT20 (length 63.0 Mb and gene density 8.71 genes/sequenced Mb). They both have comparable gene densities but strikingly different lengths. Again based on modeling results, we can see that CT1 and CT20 have internal locations in GM12878, but in BJ1-hTERT CT1 occupies a much peripheral



location. As expected given the overall gene density or length correlations of the model cluster, these mean positions are consistent with microscopy data regarding the different positioning of chr18, chr19, and chr1 in different cell types. Beyond mean position shifts, we also observe that the distributions of possible chromosome positions also match microscopic evidence for chromosome position variation in some respects. For example, in the models, CT4 in lymphoblasts is highly skewed toward peripheral positions and almost never observed near the center of the nucleus. This matches the distribution of positions of a gene located on chr4 measured by FISH [27] in which the gene was almost never observed in the interior 30 percent of the nuclear radius in any individual cell (Additional file 1: Fig. 5d). In contrast, the models predict that CT10 in fibroblasts can be found throughout the middle of the nucleus,

but almost never at the extreme center or periphery (Additional file 1: Fig. 5c). This matches the measured chr10 distribution by FISH [68]. In some cases, however, the model-inferred CT distribution matches the mean position observed in microscopy, but not the distribution of positions. For example, the models locate chr19 near the nucleus center on average in both cell types we considered, as in microscopy data (Additional file 1: Fig. 5a, b). However, in single cell images, chr19 is found over a broader range of positions [69], while our models predict a more tightly focused consistent positioning of chr19. This narrow predicted distribution likely stems from the highly specific interaction pattern for chr19, visible in the strong contact matrix and distinct PCA position of this chromosome in Fig. 1. Additionally, while chr13 has been observed by microscopy to be strongly skewed toward the nuclear periphery [68, 70], our model sometimes predicts a strongly peripheral distribution (GM12878) and sometimes a strongly internal distribution (BJ1-hTERT) (Additional file 1: Fig. 4). This occurs because there are relatively few strong inter-chromosomal interactions detected by chr13 in the Hi-C data, and the model cannot distinguish whether this means the chromosome is located on its own far to the center or far to the periphery. So, while our network model approach shows a strong ability to predict chromosome position variability for some chromosomes, in other cases, there are inherent limitations of what can be predicted from Hi-C contacts alone.

In contrast to the original Hi-C contact maps, simulated random ligation maps produced network models that were much more variable and did not form tight clusters (Additional file 1: Fig. 2a, e). This primarily happens due to non-specific pairwise inter-chromosomal interaction pattern driven by length which ultimately leads to the generation of wide variety of random network model configurations. Although these simulated random ligation datasets show a strong length based distribution of CTs from the PCA ordering, the pairwise inter-chromosomal strong interaction pattern is quite different from true length based distribution as shown by ellipsoidal cells. When we examine the distributions of chromosome positions predicted for these random models, most chromosomes showed broad indistinct distributions that did not vary based on the cell type the random matrix was derived from (Additional file 1: Fig. 6). This demonstrates that real Hi-C contacts contribute important information to our network modeling predicted radial distance distributions.

Chromosome property-based tuning improves predicted consensus radial arrangement of CTs

Although the network modeling generated radial distribution of the CTs can reveal several interesting features of the CT organization that match with prior observations, the direct correlation with microscopy measured mean positions is only modest (Fig. 3a, b). Meanwhile, the direct correlation between gene density or chromosome length and microscopy measured positions may be quite high (Additional file 1: Table 2), but these values are cell type invariant and can never be used alone to infer different CT positions between different cell types. Thus, we next tested an approach of using chromosomal properties, informed by PC1 ordering, to further tune the network model inferred mean radial positions. We add the effects of gene density and chromosome length to the averaged positions using weighted averaging and loess (locally estimated scatterplot

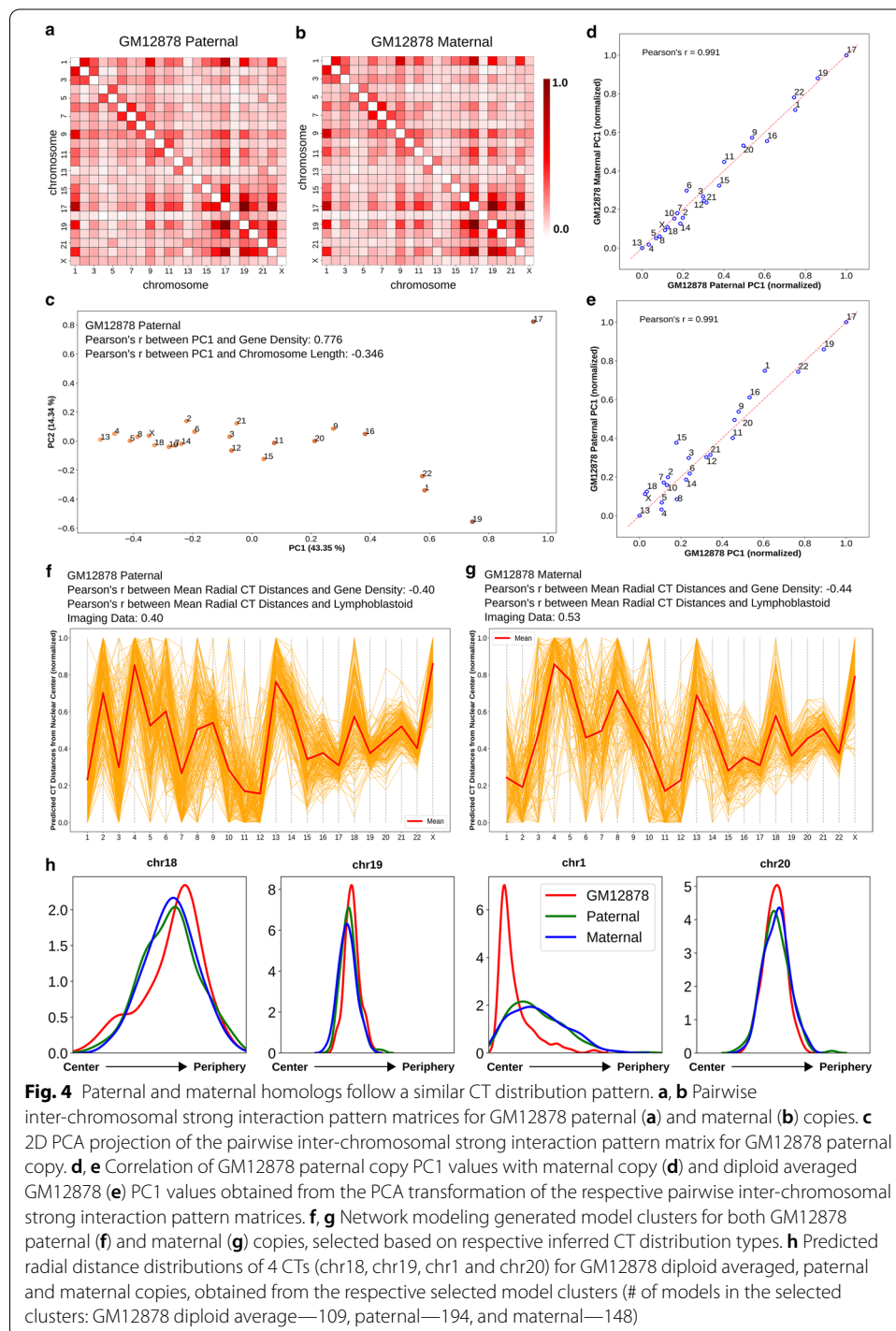
smoothing) [71] techniques. This tuning procedure is described in detail in the “[Methods](#)” section, but the basic idea behind it involves several concepts. First, we incorporate information about the length and gene density of each chromosome explicitly. Second, going back to our earlier observation that the PC1 of the pairwise inter-chromosomal strong interaction pattern matrix can provide meaningful ordering information, we combine the network model output with PC1 ordering result obtained from selected model cluster through weighted averaging. Finally, in calculating weights for our averaging calculation, we take into account how much of the variance in the chromosome positioning pattern in the selected model cluster is captured by PC1 and PC2. We find this metric captures a major distinction between random and real Hi-C data. In the previous steps, random data will give some pattern, often highly chromosome length related, but we find that PC1 and PC2 explain a much lower percentage of the overall variance in the chromosome contact pattern for random data, while for real data these first two PCs capture most of the variance, indicating that the orderings along these PCs are highly meaningful. Thus, weighting by the variance explained by PC1 and PC2 will emphasize meaningful patterns over random patterns.

In order to evaluate the accuracy of the tuning procedure, we compared the tuned radial CT positions of GM12878 and BJ1-hTERT with corresponding microscopy imaging data, as shown in Fig. 3d, e. From Fig. 3d, it can be seen that, in GM12878, most of the CT positions correlate well with the experimentally obtained position with slight displacements, apart from the CT1 position, which moves outwards compared to the experimental data. Similarly, predicted CT distances for the BJ1-hTERT cell shows high similarity with the corresponding imaging data. For this particular cell, out of all the CTs, CT19 showed a higher amount of displacement towards the center in the predicted result. When the tuning technique was tested on simulated random ligation Hi-C data generated from GM12878 and BJ1-hTERT, it produced a far weaker correlation with the fibroblast imaging result in both cases (length based inferred CT distribution type) compared to the original Hi-C analysis results (Additional file 1: Fig. 2b, f). Furthermore, we checked the consistency of the predicted tuned results between two Hi-C replicates of the BJ1-hTERT cell and found very little difference in chromosome positioning (Additional file 1: Fig. 7f). In case of GM12878, variation between replicates was higher (Additional file 1: Figs. 8f and 9), perhaps reflecting the larger variation in quality metrics between these two datasets. However, a strength of employing the PCA ordering is revealed in comparing these replicates: the PCA ordering of the strong interaction matrix is robust to such variations in dataset quality (Additional file 1: Fig. 8d).

Paternal and maternal homologs show mostly similar radial arrangement of CTs

Our modeling approach places only one copy of each chromosome in the radial organization network, and we assume that this represents the average of the two homolog positions. In any given cell, it is known from microscopy that the CT positions of homologs can be quite different [69], but we assume that both homologs would follow the same overall distribution of positions, and thus averaging their positions is reasonable. To test this assumption, we applied our approach to the deeply sequenced GM12878 Hi-C data from Rao et al. that can be mapped to maternal and paternal homologs based on allele-specific single nucleotide polymorphisms (SNPs) [72]. Once we had the genome wide

Hi-C contact matrices for the paternal and maternal copies, thresholding was applied and pairwise inter-chromosomal strong interaction pattern matrix was analyzed using PCA transformation for each of them. Figure 4c shows the 2D PCA projection of the pairwise inter-chromosomal strong interaction pattern matrix obtained from the paternal copy. By looking at this figure, we can clearly see that the paternal homologs are



following a gene density driven distribution along the PC1 axis, which we also found for the the maternal homologs (Fig. 4d y-axis). In addition, the PC1 values from the paternal and maternal copies are highly correlated to each other and also with the PC1 values obtained from the previous GM12878 data (Fig. 4d, e). This suggests that the paternal and maternal homologs follow a similar radial distribution inside the nucleus, and that it is fair to infer a single average position for both homologs. We also applied the network modeling approach to the paternal and maternal copies individually and selected the model cluster for each of them whose mean CT distances has highest absolute correlation with the inferred CT distribution type, as above. By looking at the mean CT distances from Fig. 4f, g, we can see that the mean radial CT distribution pattern between two copies are highly comparable with a single drastic exception in case of chr2, which shows an internal position in maternal copy and occupies a peripheral position in the paternal one. Again this can be explained by distinct patterns of inter-chromosomal strong interaction of chr2 paternal and maternal homologs. Finally, when we compared the radial distance profiles of the 4 CTs - CT1, CT18, CT19 and CT20 among the paternal, maternal and the standard GM12878, we found the density peaks in similar radial positions (Fig. 4h and Additional file 1: Fig. 10). However, in case of chr1, we found a strong narrow peak from the standard GM12878 data compared to the short broad peak obtained from both the paternal and maternal copies, which might arise due to the less specific inter-chromosomal interaction pattern of chr1 in both maternal and paternal copies compared to the interaction patterns of rest of the chromosomes in those copies. Overall, we observe that representing both homologs with an average predicted position is valid.

Radial CT organization in epithelial cells changes depending on their nuclear shape

After testing the performance of the analysis technique on the GM12878 and BJ1-hTERT cells, we applied our analysis to MCF10A non-tumorigenic breast epithelial cells. The Hi-C contact data of this cell was downloaded from Barutcu et al. [73] and binned at 2.5 Mb resolution. We obtained corresponding imaging data for 8 CTs (CT1, CT4, CT11, CT12, CT15, CT16, CT18 and CT21) of the MCF10A from Fritz et al. [8]. When we applied PCA to the pairwise inter-chromosomal strong interaction matrix, PC1 showed higher absolute correlation with chromosome length compared to gene density, as expected from experimental data (Fig. 5b), and corroborating the association of the length based distribution with ellipsoidal nucleus shape. Next, we used the network modeling approach to generate radial distance profiles of the 8 CTs and found the distributions match the ordering of CT positions reported in imaging data (Fig. 5c). We note that we predict some chromosomes (chr4, chr15) have a much broader distribution than others (chr16, chr21), but this distribution is not reported for comparison in microscopy data. Finally, when the tuned CT positions were compared with the experimental data, a high correlation was obtained (Fig. 5d).

Since epithelial cells can take different nuclear shapes depending on their properties in culture, we next applied our analysis to Hi-C data from another human normal mammary epithelial cell (HMEC) [72]. In this cell type, the PCA analysis of the pairwise inter-chromosomal strong interaction pattern inferred a gene density based distribution as represented in Fig. 5e. We hypothesize that this is related to the fact

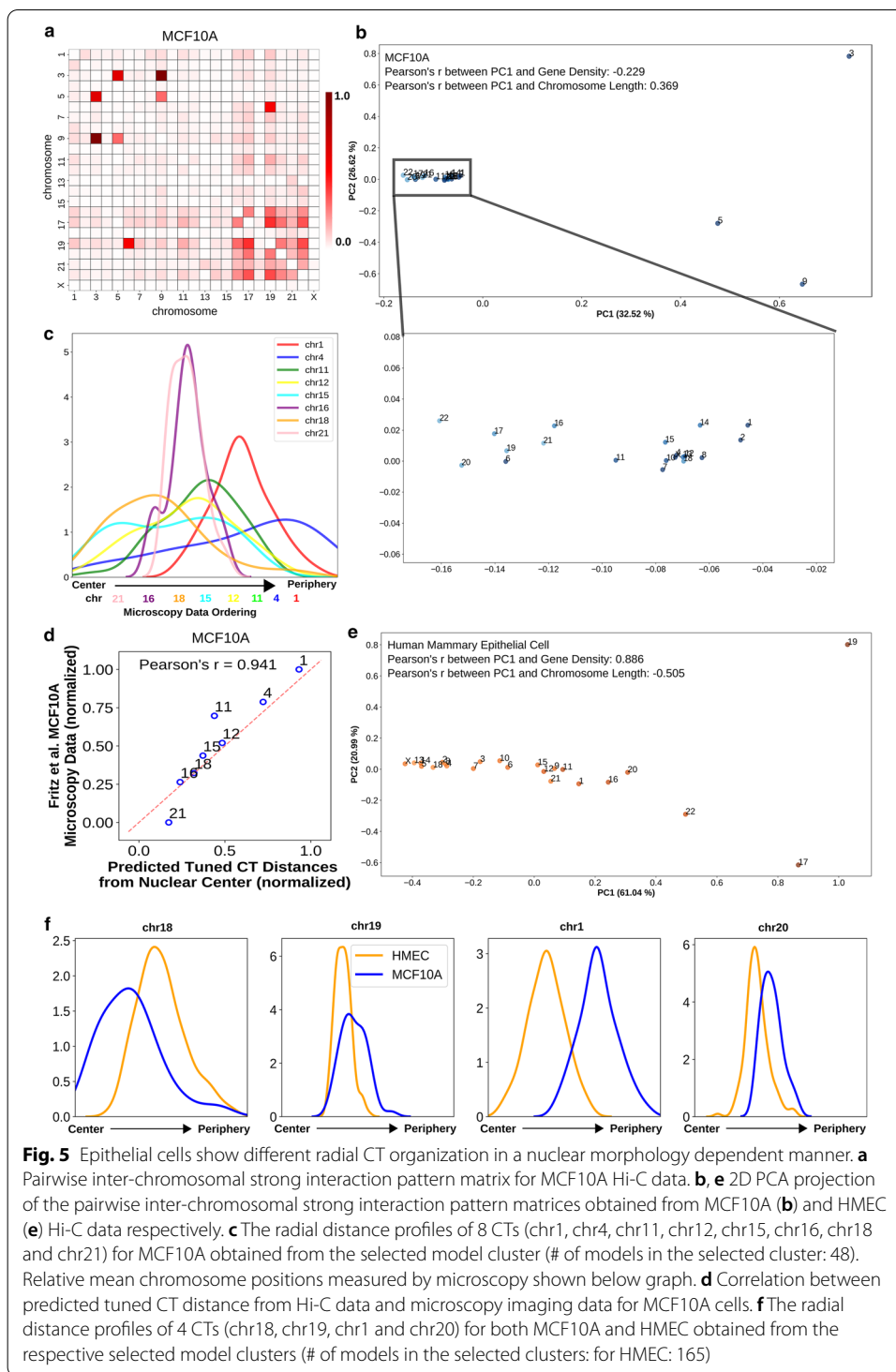


Fig. 5 Epithelial cells show different radial CT organization in a nuclear morphology dependent manner. **a** Pairwise inter-chromosomal strong interaction pattern matrix for MCF10A Hi-C data. **b, e** 2D PCA projection of the pairwise inter-chromosomal strong interaction pattern matrices obtained from MCF10A (**b**) and HMEC (**e**) Hi-C data respectively. **c** The radial distance profiles of 8 CTs (chr1, chr4, chr11, chr12, chr15, chr16, chr18 and chr21) for MCF10A obtained from the selected model cluster (# of models in the selected cluster: 48). Relative mean chromosome positions measured by microscopy shown below graph. **d** Correlation between predicted tuned CT distance from Hi-C data and microscopy imaging data for MCF10A cells. **f** The radial distance profiles of 4 CTs (chr18, chr19, chr1 and chr20) for both MCF10A and HMEC obtained from the respective selected model clusters (# of models in the selected clusters: for HMEC: 165)

that these cells are closer to normal epithelium than MCF10A, which harbor chromosomal translocations, and that classical epithelial patterns would lead to a more spherical nucleus shape even in 2D culture, rather the flat, spreading growth pattern of MCF10A cells [74]. When the network modeling inferred radial distance profiles

of four CTs - CT1, CT18, CT19 and CT20 were compared between MCF10A and HMEC (Fig. 5f), CT19 and CT20 showed density peaks at a similar radial position in both cell types. On the other hand, CT18 showed a preference for interior radial positions in MCF10A (length based) and for peripheral positions in HMEC (gene density based) and in case of CT1 that trend was the opposite. These results suggest that similar cell types can show different CT organization that correlates with their cell morphology, and Hi-C contact patterns can capture these differences.

Variation between the radial arrangement of lymphoblastoid cells and neutrophils

To test the utility of our analyses on a cell having an irregular nuclear shape, we focused on neutrophils and obtained the corresponding Hi-C data from Javierre et al. [75]. Neutrophils have a multi-lobed nucleus with a toroid shaped genome where lobes are connected by thin filaments [76]. When we analyzed the pairwise inter-chromosomal strong interaction pattern matrix using PCA, we found that the PC1 values have a higher absolute correlation with gene density (Fig. 6b), though the effect is much weaker compared to GM12878 case, which is consistent with the divergence of a neutrophil from its precursor cell’s round nucleus shape. This is another indication that the Hi-C contact data provides more information about chromosome positioning than just the underlying length based or gene density based distribution. Based on microscopy imaging data, it

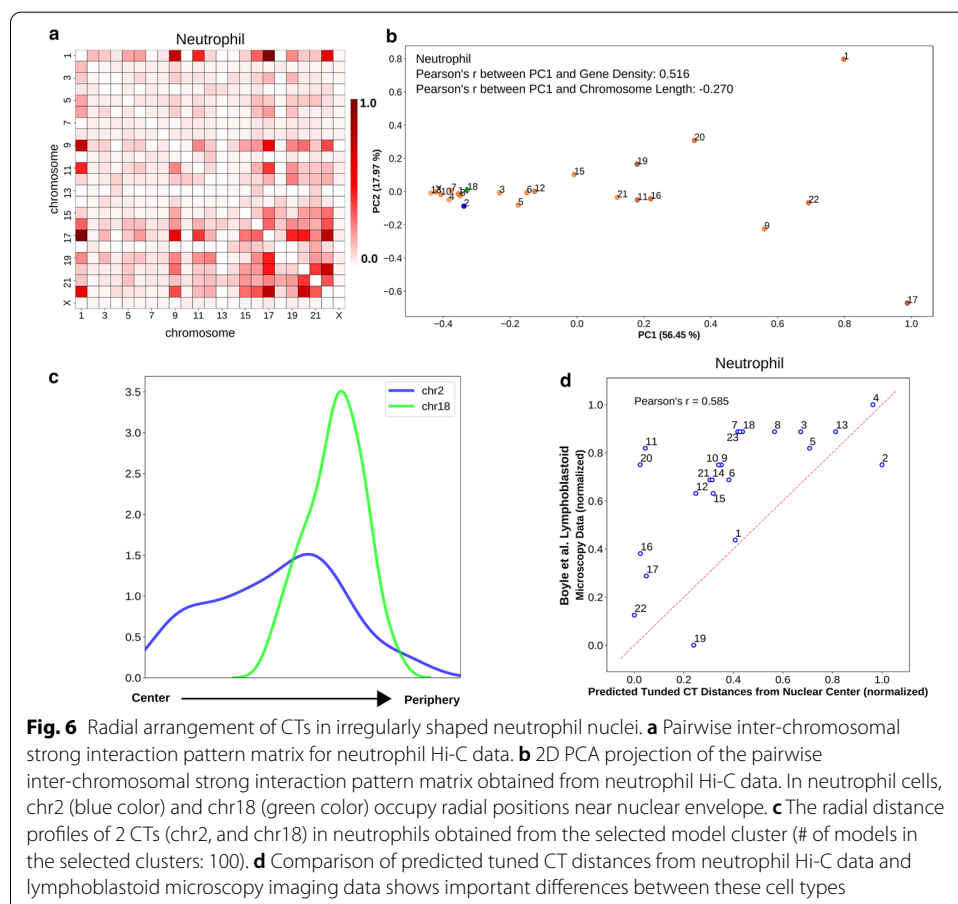


Fig. 6 Radial arrangement of CTs in irregularly shaped neutrophil nuclei. **a** Pairwise inter-chromosomal strong interaction pattern matrix for neutrophil Hi-C data. **b** 2D PCA projection of the pairwise inter-chromosomal strong interaction pattern matrix obtained from neutrophil Hi-C data. In neutrophil cells, chr2 (blue color) and chr18 (green color) occupy radial positions near nuclear envelope. **c** The radial distance profiles of 2 CTs (chr2, and chr18) in neutrophils obtained from the selected model cluster (# of models in the selected clusters: 100). **d** Comparison of predicted tuned CT distances from neutrophil Hi-C data and lymphoblastoid microscopy imaging data shows important differences between these cell types

has been reported that chr2 and chr18, which have drastic differences in their lengths, both occupy a position near nuclear envelope in neutrophil cells [77]. Upon inspecting the positions of those chromosomes along the PC1 axis in Fig. 6b, we found the two chromosomes in very close proximity near the left extreme of the PC1 axis. Although projection data does not infer the directionality of the ordering, based on its higher similarity with gene density, we can assume the left extreme as the periphery which contains mostly gene poor chromosomes. Furthermore, for those two chromosomes, we also observed similar trend in the network model-derived radial CT distance profiles as represented in Fig. 6c. The network model predicts a highly variable positioning of chr2, likely related to the overall dearth of strong contacts between chr2 and other chromosomes in the initial thresholded contact map. We observe that the predicted tuned CT positions for neutrophils are only weakly correlated with lymphoblastoid imaging data (Fig. 6d), revealing the ability of the model to predict different chromosome positions from different initial Hi-C contact data.

Discussion

In this study, our goal has not been only to predict one final set of CT positions from Hi-C data. Instead, we have explored what aspects of 3D chromosome radial positioning can and cannot be inferred with a series of direct and non-computationally intensive Hi-C contact analysis approaches that have not been previously used for this purpose.

Our results demonstrate that a straightforward statistical calculation (PCA) on the pattern of strong inter-chromosomal contacts can capture important biological features of CT radial positioning. With this approach, we not only can capture important patterns of gene density or chromosome length based ordering of chromosomes previously observed for very different cell types, but also can detect meaningful shifts of individual chromosomes in related cell types. We were able to infer changes in CT ordering in premature aging and senescence directly from contact data, and these changes are supported by previous microscopy results. We also note that this PCA ordering approach is actually quite robust to differences in Hi-C data quality and depth. While our subsequent network graph layout approach was somewhat sensitive to different quality metrics of different Hi-C replicates, the PCA ordering of chromosomes was robust to different levels of noise in these datasets.

We have also demonstrated that a network graph layout algorithm approach can generate an ensemble of models that capture the experimentally validated mean and variation of CT positions in a cell type. We demonstrate that these approaches can be applied to a variety of cell types and can even detect differences in underlying CT radial distributions between highly related cell types (two different breast epithelial cell types). Interestingly, for these cell types, the difference in CT distribution corresponds to documented differences in nucleus geometry of these cell types when grown in culture. This suggests that not only do spherical and elliptical nuclei exhibit different CT organization in completely different cell lineages (lymphoblast vs. fibroblast), as previously documented, but that even cells within the same overall type may have different CT positioning associated with their nucleus shape. This result emphasizes the importance of approaches to detect changes in CT ordering directly from a given cell type in its particular condition, rather than assuming that a measurement in one circumstance

will define CT organization principles across, for example, all epithelial cells. Our rapid approaches to inferring CT distribution makes screening across such varieties of data-sets more feasible, as a complement to more intensive computational models or microscopic measurements.

In addition to the validation of these approaches that makes them useful for future applications, our results also show that while Hi-C contacts are useful for inferring a relative ordering of radial chromosome positions, inferring absolute ordering often requires an external reference point. We observe this in that the PCA based ordering of chromosomes can represent either direction (interior-exterior or vice versa) and that some clusters of network models display the same chromosome relative ordering in reverse orientation. This is an important factor to consider in any model that attempts to use Hi-C contacts in isolation to generate 3D structure models.

The analysis approach described in this paper is highly flexible due to its modular nature and can be integrated with different kinds of genome analysis applications. For example, both PCA ordering and inferred tuned positions can be used to characterize the changes in the radial CT organization in a perturbed (e.g. relocation of CTs during DNA damage response [78]) or diseased cell (e.g. mislocalization of CT18 and CT19 in lamin B2 depleted colorectal cancer cells [79]) compared to a control cell, which in turn allows us to study how these changes affect higher-order chromatin structure. The method also generates a population of 3D network structures, which can further be used to characterize inter-chromosomal dynamics and can be compared to single-cell Hi-C results.

Conclusions

We have described a set of approaches that can be used sequentially or as separate modules to predict chromosome territory radial positioning in the nucleus from Hi-C data. We find that analyzing only the strongest interchromosomal contacts emphasizes differences in CT arrangement between cell types, and that PCA on this strong contact matrix can be used as a simple, fast approach to detect relative CT radial positions. We describe a network modeling approach that builds on this overall pattern detected by PCA to simulate the variability in CT positioning across the cell population. Finally, we demonstrate both the strengths and limitations of what Hi-C data alone can predict about radial CT organization and show that intrinsic chromosome properties can be added to tune CT organization predictions. These methods provide researchers with additional tools to infer the properties of radial CT organization for the growing numbers of cell types that have available Hi-C data but not detailed microscopic measurements of chromosome positions.

Methods

The schematic representation of our whole analysis approach is given in Fig. 7.

Pre-processing of Hi-C data

All original Hi-C data sources for this study are listed in Additional file 1: Table 1. We mapped the data to hg19 to obtain a set of unique valid interacting fragment pairs and

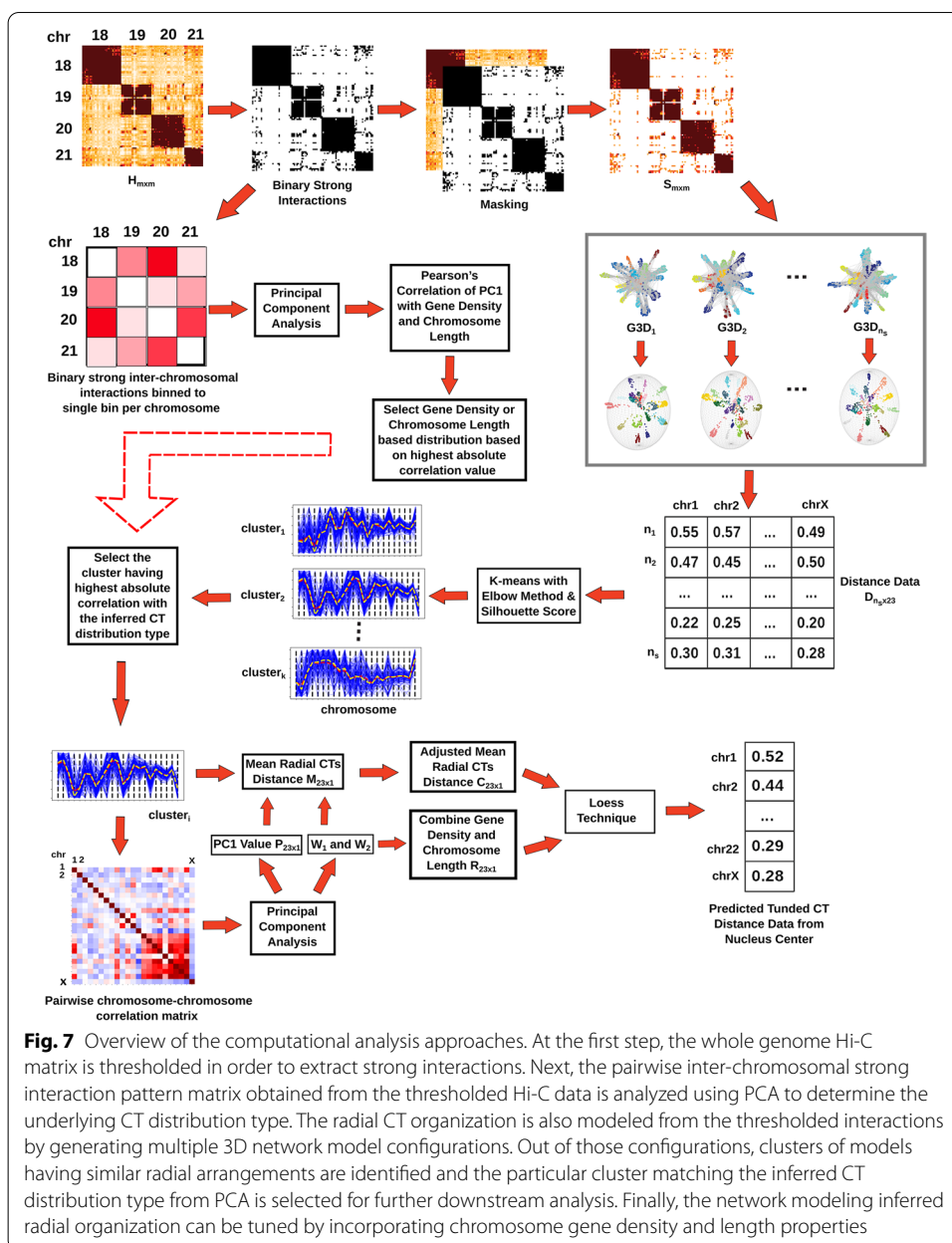


Fig. 7 Overview of the computational analysis approaches. At the first step, the whole genome Hi-C matrix is thresholded in order to extract strong interactions. Next, the pairwise inter-chromosomal strong interaction pattern matrix obtained from the thresholded Hi-C data is analyzed using PCA to determine the underlying CT distribution type. The radial CT organization is also modeled from the thresholded interactions by generating multiple 3D network model configurations. Out of those configurations, clusters of models having similar radial arrangements are identified and the particular cluster matching the inferred CT distribution type from PCA is selected for further downstream analysis. Finally, the network modeling inferred radial organization can be tuned by incorporating chromosome gene density and length properties

then binned the data at a resolution of 2.5 Mb, where each bin of the contact matrix represents the interaction frequency between a pair of 2.5 Mb genomic regions. We then perform two pre-processing steps on the Hi-C matrix. First, to remove biases related to GC content and cut site frequency, the raw contact matrix is normalized using the ICE technique [62]. Next, the unmappable (repetitive) genomic regions are removed from the normalized Hi-C matrix. Also, we do not consider chrY in our analysis since this chromosome is not present in both male and female cell types, and we remove those corresponding bins from the Hi-C matrix. This modified Hi-C matrix is used for further downstream processing.

Let $H_{m \times m}$ represent a square symmetric Hi-C contact matrix where each row and column correspond to genomic regions (bins) of specific size, with each element h_{ij} representing the normalized number of contacts between i th and j th bins.

Determining a threshold to capture strong interactions in Hi-C matrix

After preprocessing the Hi-C contact data, we proceed to identify the strong interactions which have the best potential to infer cell-type specific radial organization of the CTs, distinct from the levels of average background noise. To do this, we define a cutoff limit and apply that to $H_{m \times m}$, to obtain a matrix of only filtered strong interactions $S_{m \times m}$.

$$\begin{cases} s_{ij} = h_{ij} & \text{for } h_{ij} > h_{cut} \\ s_{ij} = 0 & \text{otherwise} \end{cases} \quad (3)$$

where s_{ij} is the element of matrix $S_{m \times m}$, which represents the normalized number of contacts between i th and j th bins and h_{cut} is the cutoff limit. For 2.5 Mb genomic resolution, the value of cutoff limit is set to: $h_{cut} = 95$ th percentile of all the genome-wide interactions from $H_{m \times m}$. The significance of and reason for choosing this 95th percentile cutoff is discussed in the “Results” section. Due to the application of this cutoff limit h_{cut} , the resultant matrix $S_{m \times m}$ contains mostly the intra-chromosomal interactions and a few inter-chromosomal interactions. In addition, to detect significant chromosomal interactions, we apply the FitHiC tool [59] with default parameters to the genome-wide Hi-C contact data, binned at 2.5 Mb resolution. From those significant intra- and inter-chromosomal interactions detected, we further select highly significant interactions for analysis by applying a q-value cutoff - 10^{-2} (for intra-chromosomal) and 10^{-12} (for inter-chromosomal). The reason behind selecting a very stringent q-value cutoff for inter-chromosomal interactions is to make the number of significant inter-chromosomal interactions from FitHiC comparable to the strong interactions from our method.

Simulating random ligation Hi-C from original Hi-C data

Random ligation Hi-C data are simulated by taking an original raw/non-normalized Hi-C contact map, binned at 2.5 Mb, and shuffling the bins (including the diagonal) of this matrix five times. Then, the shuffled matrix is passed through the ICE normalization step [62]. The reason for generating the random ligation matrix from the original matrix by random shuffling is to ensure that both the contact matrices have an equal number of total reads and a similar dynamic range of values, ensuring a matched comparison of the results from real and random Hi-C data.

Identifying CT radial distribution patterns with PCA transformation on the pairwise strong inter-chromosomal interaction pattern matrix

For our most direct approach to infer a radial organization pattern from Hi-C data, we begin by creating a pairwise strong inter-chromosomal interaction pattern matrix. For the chosen h_{cut} value, the number of bin pairs passing the threshold are summed between each pair of chromosomes to a single bin (Eqn. 4). This sum thus includes a count of the number of bin pairs that passed a threshold, rather than the total number of interactions in each included bin.

$$CHR_{ij}^{STRONG} = \text{total number of strong interaction bins between chr } i \text{ and } j \quad (4)$$

Then, we exclude cis contacts and look only at the inter-chromosomal component of this pairwise strong interaction pattern (Eqn. 5).

$$CHR_{ij}^{TRANS} = \begin{cases} CHR_{ij}^{STRONG} & \text{if } i \neq j \\ 0 & \text{otherwise} \end{cases} \quad (5)$$

To capture the major interaction trends from this matrix, we apply principal component analysis (PCA) to this pairwise strong inter-chromosomal interaction pattern. We then calculate the projection of all chromosomes onto PC1, and we find that the ordering of chromosomes along this PC can capture the radial organization ordering of the chromosomes.

Constructing 3D network graphs from Hi-C matrix

The thresholded strong interaction matrix $S_{m \times m}$ is treated as a weighted adjacency matrix to generate an undirected weighted graph G , where nodes represent genomic bins and the number of interactions between each pair of bins is used as the edge weight. Next, we apply a 3D Fruchterman-Reingold (FR) force based layout to draw the undirected Hi-C graph in 3D space. This graph drawing layout adds an attractive force between the connected nodes and creates repulsion between the nodes that are not connected. Along with this, gravitational force is used in this layout to pull the nodes towards the center. The rationale behind using the FR layout is that it uncovers the intrinsic structure of the network, as the strong interactions only matrix $S_{m \times m}$ has a large number of intra-chromosomal interactions and fewer inter-chromosomal interactions. The resulting 3D network graph infers the radial organization of the CTs inside the nucleus. Next, n_s number of 3D network graphs - $G3D_1, G3D_2, \dots, G3D_{n_s}$ are generated by performing independent runs of the FR algorithm with different random initial configurations to model the variability in the radial CT organization.

Fitting a geometrical structure to the 3D network graphs

For each 3D network graph, the next objective is to find the distance of each CT from the center of that network graph. This step is performed in two parts.

In the first part, given a 3D network graph, the 3D Cartesian coordinates of the nodes are extracted and then Khachiyan's algorithm is used to find a minimum volume ellipsoid enclosing the set of nodes [80]. The center of the fitted object is calculated and assigned as the nucleus center. Along with this, the center of mass of each of the individual CTs is calculated from the coordinates of the nodes of the network graph.

$$COM_j^i = \frac{\sum \text{coordinates of nodes of CT } j \text{ in } G3D_i}{\# \text{of nodes of CT } j \text{ in } G3D_i} \quad (6)$$

where COM_j^i represents the center of mass of a particular CT j in the structure $G3D_i$ and $j \in \{1, 2, \dots, 22, X\}$.

After obtaining the center of mass of each CT and nucleus center, the Euclidean distance is calculated between the center of mass of each CT and the nucleus center. Also, to remove the heterogeneity that arises from different minimum volume ellipsoid fits of

different 3D network structures, for each structure the distance values are normalized in the [0,1] range using Min-max normalization. In this way, by iterating over all of the 3D network graphs, for each CT, n_s distances from the center of the nucleus are obtained. This distance information is represented in a matrix $D_{n_s \times 23}$ of size $n_s \times 23$, where rows represent n_s 3D network graphs and columns correspond to 23 CTs (22 autosomes and one X chromosome), and each matrix entry d_{ij} is the distance of the CT j from the nucleus center of network graph $G3D_i$.

The parameter n_s represents the number of configurations required to mimic the heterogeneity in CT arrangement originating during cell division within the same cell type. To estimate this parameter, for different values of n_s we compare the fitted distance profiles of all CTs between GM12878 and BJ1-hTERT cells using non-parametric hypothesis testing - a Two-sided Mann-Whitney U test. From the statistical test results, it can be observed that with increasing n_s value, more chromosomes show significant differences in their CT distance profiles between lymphoblastoid and fibroblast cells, reaching a maximum at $n_s = 1000$ (Additional file 1: Fig. 11). Hence, we set the n_s parameter to 1000 in our analysis.

Identifying the specific cluster of network models having meaningful CT organization

The network modeling approach produces 3D graph models with heterogeneous CT organization, as intended, but we find that some groups of models may capture an inverted ordering of some groups of CTs or all CTs (reversing central to peripheral distances). Thus, rather than blindly averaging all models together, we first identify these different clusters of organization patterns and then choose for further consideration the cluster of models that follows the radial CT organization distribution captured by the PCA analysis described above. We perform K-means clustering [81] with deterministic initialization on the distance matrix $D_{n_s \times 23}$ by treating the models (rows) as samples and chromosomes (columns) as features. The initial centroid positions for the clustering technique are calculated using the algorithm from Nazeer et al. [82]. In addition to that, to estimate optimal number clusters for the K-means, we use a modified elbow method approach. Here, first we calculate inertia [83] which represents within cluster sum of squares for different increasing number of clusters and detect the elbow of the inertia curve with the help of the algorithm from Satopaa et al. [84]. After detecting the elbow, again we calculate the average silhouette score [83] for each of the different number of clusters and select the point as optimal number of clusters which will have the highest average silhouette score in the vicinity of the elbow point (two points upstream of elbow, two points downstream and the elbow itself). As our set of predictions for further analysis, we take the cluster whose mean radial CT positions shows the highest absolute correlation with the CT distribution type obtained by the PCA transformation of the pairwise inter-chromosomal strong interaction pattern matrix.

Gene density and chromosome length based tuning of consensus radial organization of CTs

The last aspect of our approach considers how the predicted averaged radial CT distances can be further tuned using two chromosomal properties - gene density and chromosome length. Human chromosome gene density (genes/ sequenced Mb) is

obtained from “Short guide to the human genome” [85] and length information from UCSC Genome Browser hg19 [86]. But before discussing the steps involved with this tuning procedure, let us assume, $D'_{t \times 23}$ represents the selected model cluster based on inferred CT distribution type. This cluster contains t number of network models and the average CT distance of each chromosome across all of the t models is denoted by the column vector $M_{23 \times 1}$. After having the selected model cluster $D'_{t \times 23}$, we calculate a pairwise chromosome correlation matrix of size 23×23 from that and perform PCA transformation on that correlation matrix. The PC1 value obtained from this transformation represents the major separation of the chromosomes based on pairwise interactions and is denoted by the column vector $P_{23 \times 1}$. Once we have the two column vectors - $M_{23 \times 1}$ and $P_{23 \times 1}$ representing the average radial position of the CTs and their separation respectively, we combine them using a weighted averaging technique as per Eqn. 7.

$$C_{23 \times 1} = W_1 * M_{23 \times 1} + W_2 * P_{23 \times 1}$$

$$W_1 = \frac{\% \text{ of variance explained by PC1}}{\% \text{ of variance explained by PC1} + \% \text{ of variance explained by PC2}} \quad (7)$$

$$W_2 = \frac{\% \text{ of variance explained by PC2}}{\% \text{ of variance explained by PC1} + \% \text{ of variance explained by PC2}}$$

where $C_{23 \times 1}$ represents the resultant combined column vector and W_1 and W_2 denote the weights calculated from the percentage of the the variance explained by PC1 and PC2.

Next, in order to model the effect of both chromosome length and gene density, we combine these two properties using weighted averaging as described in the Eqn. 8. The first and second components on the right hand side of each of the equations are related to the normalized gene density and normalized chromosome length respectively. Also, we use the same set of weights W_1 and W_2 in this equation but the weights are ordered based on the inferred CT distribution type. For example, if the inferred CT distribution follows a gene density based radial positioning, in that case the component having gene density will have the higher weight W_1 and vice versa.

$$R_{23 \times 1} = \begin{cases} W_1 * \exp(1 - GD_{23 \times 1}) + W_2 * LN_{23 \times 1} & \text{if inferred - gene density} \\ W_2 * \exp(1 - GD_{23 \times 1}) + W_1 * LN_{23 \times 1} & \text{if inferred - length} \end{cases} \quad (8)$$

Here, GD and LN are column vectors of size 23×1 and contain the gene density and length of each chromosomes respectively.

Following the chromosomal properties modeling procedure, in the final step of the tuning, the modified consensus radial distance of the CTs $C_{23 \times 1}$ is locally smoothed based on $R_{23 \times 1}$ using the Loess technique. The radial distance of the CTs obtained from this step represents the tuned arrangements of the CTs inside the nucleus.

Supplementary information

Supplementary information accompanies this paper at <https://doi.org/10.1186/s12859-020-03841-7>.

Additional file 1. This file contains Tables S1–S2 and Figures S1–S11.

Abbreviations

CT: Chromosome territory; PCA: Principal component analysis; FISH: Fluorescence in situ hybridization.

Acknowledgements

We thank Jacob Sanders for pre-publication access to his BJ1-hTERT and GM12878 Hi-C datasets and Dr. Rebeca San Martin and Rosela Gollosi for their valuable suggestions and feedback in the development of the project.

Authors' contributions

P.D. and R.P.M. conceived and designed the project with conceptual input from T.S. P.D. performed all computational analysis, data processing, and prepared all figures. P.D. and R.P.M. wrote the manuscript with input from T.S. All authors have read and approved the manuscript.

Funding

This research was supported in part by NIH NIGMS grant R35GM133557 to R.P.M. Salaries, publication costs, and computing resources were funded by this NIH grant. The funding body did not play any role in the design of the study, the collection, analysis, and interpretation of data, or in writing the manuscript.

Availability of data and materials

The source code of the analysis technique is available at <https://github.com/rpmccordlab/Radial-CT-Analysis-HiC>. Source Hi-C data are publicly available from GEO and EGA.

Ethics approval and consent to participate

The human neutrophil Hi-C data utilized here was obtained by the PCHI-C Consortium and analyzed with permission of the European Genome Archive.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Author details

¹ UT-ORNL Graduate School of Genome Science and Technology, University of Tennessee, Knoxville, TN 37996, USA.

² Department of Biochemistry and Cellular and Molecular Biology, University of Tennessee, Knoxville, TN 37996, USA.

Received: 9 January 2020 Accepted: 27 October 2020

Published online: 10 November 2020

References

- McCord RP, Kaplan N, Giorgetti L. Chromosome conformation capture and beyond: toward an integrative view of chromosome structure and function. *Mol Cell*. 2020;77(4):688–708.
- Cremer T, Cremer M. Chromosome territories. *Cold Spring Harbor Perspect Biol*. 2010;2(3):003889.
- Cremer T, Kurz A, Zirbel R, Dietzel S, Rinke B, Schröck E, Speicher MR, Mathieu U, Jauch A, Emmerich P, et al. Role of chromosome territories in the functional compartmentalization of the cell nucleus. In: *Cold spring harbor symposia on quantitative biology*, vol. 58, 1993; pp. 777–792. Cold Spring Harbor Laboratory Press.
- Rosin LF, Crocker O, Isenhardt RL, Nguyen SC, Xu Z, Joyce EF. Chromosome territory formation attenuates the translocation potential of cells. *eLife*. 2019;8:49553.
- Kuroda M, Tanabe H, Yoshida K, Oikawa K, Saito A, Kiyuna T, Mizusawa H, Mukai K. Alteration of chromosome positioning during adipocyte differentiation. *J Cell Sci*. 2004;117(24):5897–903.
- Mehta IS, Eskiw CH, Arican HD, Kill IR, Bridger JM. Farnesyltransferase inhibitor treatment restores chromosome territory positions and active chromosome dynamics in Hutchinson–Gilford progeria syndrome cells. *Genome Biol*. 2011;12(8):74.
- Bercht Pflieghaar K, Taimen P, Butin-Israeli V, Shimi T, Langer-Freitag S, Markaki Y, Goldman AE, Wehnert M, Goldman RD. Gene-rich chromosomal regions are preferentially localized in the lamin b deficient nuclear blebs of atypical progeria cells. *Nucleus*. 2015;6(1):66–76.
- Fritz AJ, Stojkovic B, Ding H, Xu J, Bhattacharya S, Gaile D, Berezney R. Wide-scale alterations in interchromosomal organization in breast cancer cells: defining a network of interacting chromosomes. *Human Mol Genet*. 2014;23(19):5133–46.
- Boyle S, Gilchrist S, Bridger JM, Mahy NL, Ellis JA, Bickmore WA. The spatial organization of human chromosomes within the nuclei of normal and emerin-mutant cells. *Human Mol Genet*. 2001;10(3):211–20.
- Sun HB, Shen J, Yokota H. Size-dependent positioning of human chromosomes in interphase nuclei. *Biophys J*. 2000;79(1):184–90.
- Van Steensel B, Belmont AS. Lamina-associated domains: links with chromosome architecture, heterochromatin, and gene repression. *Cell*. 2017;169(5):780–91.
- Amendola M, van Steensel B. Mechanisms and dynamics of nuclear lamina–genome interactions. *Curr Opin Cell Biol*. 2014;28:61–8.
- Guelen L, Pagie L, Brasset E, Meuleman W, Faza MB, Talhout W, Eussen BH, de Klein A, Wessels L, de Laat W, et al. Domain organization of human chromosomes revealed by mapping of nuclear lamina interactions. *Nature*. 2008;453(7197):948.
- Bridger J, Boyle S, Kill I, Bickmore W. Re-modelling of nuclear architecture in quiescent and senescent human fibroblasts. *Curr Biol*. 2000;10(3):149–52.

15. Tanabe H, Habermann FA, Solovei I, Cremer M, Cremer T. Non-random radial arrangements of interphase chromosome territories: evolutionary considerations and functional implications. *Mutat Res Fundam Mol Mech Mutagenes*. 2002;504(1–2):37–45.
16. Tanabe H, Küpper K, Ishida T, Neusser M, Mizusawa H. Inter- and intra-specific gene-density-correlated radial chromosome territory arrangements are conserved in old world monkeys. *Cytogenet Genome Res*. 2005;108(1–3):255–61.
17. Bolzer A, Kreth G, Solovei I, Koehler D, Saracoglu K, Fauth C, Müller S, Eils R, Cremer C, Speicher MR, et al. Three-dimensional maps of all chromosomes in human male fibroblast nuclei and prometaphase rosettes. *PLoS Biol*. 2005;3(5):157.
18. Kind J, Pagie L, de Vries SS, Nahidiar L, Dey SS, Bienko M, Zhan Y, Lajoie B, de Graaf CA, Amendola M, et al. Genome-wide maps of nuclear lamina interactions in single human cells. *Cell*. 2015;163(1):134–47.
19. Lucas J, Cervantes E. Significant large-scale chromosome territory movement occurs as a result of mitosis, but not during interphase. *Int J Radiat Biol*. 2002;78(6):449–55.
20. Walter J, Schermelleh L, Cremer M, Tashiro S, Cremer T. Chromosome order in hela cells changes during mitosis and early g1, but is stably maintained during subsequent interphase stages. *J Cell Biol*. 2003;160(5):685–97.
21. Kind J, Pagie L, Ortobozkoyun H, Boyle S, de Vries SS, Janssen H, Amendola M, Nolen LD, Bickmore WA, van Steensel B. Single-cell dynamics of genome–nuclear lamina interactions. *Cell*. 2013;153(1):178–92.
22. Strickfaden H, Zunhammer A, van Koningsbruggen S, Köhler D, Cremer T. 4d chromatin dynamics in cycling cells: Theodor Boveri's hypotheses revisited. *Nucleus*. 2010;1(3):284–97.
23. Solovei I, Cremer M. 3D-FISH on cultured cells combined with immunostaining. In: Bridger J, Volpi E, editors. *Fluorescence in situ hybridization (FISH)*. Berlin: Springer; 2010. p. 117–26.
24. Shachar S, Voss TC, Pegoraro G, Sciascia N, Misteli T. A high-throughput imaging-based mapping platform for the systematic identification of gene positioning factors. *Cell*. 2015;162(4):911.
25. Iannuccelli E, Mompert F, Gellin J, Lahbib-Mansais Y, Yerle M, Boudier T. NEMO: a tool for analyzing gene and chromosome territory distributions from 3D-fish experiments. *Bioinformatics*. 2010;26(5):696–7.
26. Ollion J, Cochenne J, Loll F, Escudé C, Boudier T. TANGO: a generic tool for high-throughput 3d image analysis for studying nuclear organization. *Bioinformatics*. 2013;29(14):1840–1.
27. Gué M, Messaoudi C, Sun JS, Boudier T. Smart 3D-FISH: automation of distance analysis in nuclei of interphase cells by image processing. *Cytom Part A J Int Soc Anal Cytol*. 2005;67(1):18–26.
28. Fritz AJ, Barutcu AR, Martin-Buley L, van Wijnen AJ, Zaidi SK, Imbalzano AN, Lian JB, Stein JL, Stein GS. Chromosomes at work: organization of chromosome territories in the interphase nucleus. *J Cell Biochem*. 2016;117(1):9–19.
29. Lieberman-Aiden E, Van Berkum NL, Williams L, Imakaev M, Ragozcy T, Telling A, Amit I, Lajoie BR, Sabo PJ, Dorschner MO, et al. Comprehensive mapping of long-range interactions reveals folding principles of the human genome. *Science*. 2009;326(5950):289–93.
30. Gollosi R, Sanders JT, McCord RP. Iteratively improving Hi-C experiments one step at a time. *Methods*. 2018;142:47–58.
31. Dekker J, Belmont AS, Guttman M, Leshyk VO, Lis JT, Lomvardas S, Mirny LA, O'shea CC, Park PJ, Ren B, et al. The 4D nucleome project. *Nature*. 2017;549(7671):219–26.
32. Schoenfelder S, Javierre B-M, Furlan-Magaril M, Wingett SW, Fraser P. Promoter capture HI-C: high-resolution, genome-wide profiling of promoter interactions. *J Vis Exp*. 2018;136:e57320.
33. Maass PG, Barutcu AR, Rinn JL. Interchromosomal interactions: a genomic love story of kissing chromosomes. *J Cell Biol*. 2019;218(1):27–38.
34. Zhang X, Zhang Y, Zhu X, Purmann C, Haney MS, Ward T, Khechaduri A, Yao J, Weissman SM, Urban AE. Local and global chromatin interactions are altered by large genomic deletions associated with human brain development. *Nat Commun*. 2018;9(1):1–15.
35. Pal K, Forcato M, Jost D, Sexton T, Vaillant C, Salviato E, Mazza EMC, Lugli E, Cavalli G, Ferrari F. Global chromatin conformation differences in the drosophila dosage compensated chromosome x. *Nat Commun*. 2019;10(1):1–16.
36. Xiong K, Ma J. Revealing Hi-C subcompartments by imputing inter-chromosomal chromatin interactions. *Nat Commun*. 2019;10:1–12.
37. Horta A, Monahan K, Bashkirova E, Lomvardas S. Cell type-specific interchromosomal interactions as a mechanism for transcriptional diversity. 2018; bioRxiv 287532
38. Steininger A, Ebert G, Becker BV, Assaf C, Möbs M, Schmidt CA, Grabarczyk P, Jensen LR, Przybylski GK, Port M, et al. Genome-wide analysis of interchromosomal interaction probabilities reveals chained translocations and overrepresentation of translocation breakpoints in genes in a cutaneous t-cell lymphoma cell line. *Front Oncol*. 2018;8:183.
39. Gollosi R, San Martin R, Das P, Raines TI, Thurston DM, Freeman TF, McCord RP. Constricted migration contributes to persistent 3d genome structure changes associated with an invasive phenotype in melanoma cells. 2019; bioRxiv 856583.
40. Oluwadare O, Highsmith M, Cheng J. An overview of methods for reconstructing 3-D chromosome and genome structures from Hi-C data. *Biol Proced Online*. 2019;21(1):7.
41. Meluzzi D, Arya G. Computational approaches for inferring 3D conformations of chromatin from chromosome conformation capture data. *Methods*. 2019.
42. Kalhor R, Tjong H, Jayathilaka N, Alber F, Chen L. Genome architectures revealed by tethered chromosome conformation capture and population-based modeling. *Nat Biotechnol*. 2012;30(1):90–8.
43. Chiariello AM, Annunziatella C, Bianco S, Esposito A, Nicodemi M. Polymer physics of chromosome large-scale 3d organisation. *Sci Rep*. 2016;6:29775.
44. Di Pierro M, Zhang B, Aiden EL, Wolynes PG, Onuchic JN. Transferable model for chromosome architecture. *Proc Natl Acad Sci*. 2016;113(43):12168–73.
45. Serra F, Baù D, Goodstadt M, Castillo D, Filion GJ, Marti-Renom MA. Automatic analysis and 3D-modelling of Hi-C data using tadbit reveals structural features of the fly chromatin colors. *PLoS Comput Biol*. 2017;13(7):1005665.
46. Zhang B, Wolynes PG. Genomic energy landscapes. *Biophys J*. 2017;112(3):427–33.

47. Wettermann S, Brems M, Siebert J, Vu G, Stevens T, Virnau P. A minimal gö-model for rebuilding whole genome structures from haploid single-cell Hi-C data. *Comput Mater Sci.* 2020;173:109178.
48. Gibcus JH, Samejima K, Goloborodko A, Samejima I, Naumova N, Nuebler J, Kanemaki MT, Xie L, Paulson JR, Earnshaw WC, et al. A pathway for mitotic chromosome formation. *Science.* 2018;359(6376):1–12.
49. Tjong H, Li W, Kalhor R, Dai C, Hao S, Gong K, Zhou Y, Li H, Zhou XJ, Le Gros MA, et al. Population-based 3D genome structure analysis reveals driving forces in spatial genome organization. *Proc Natl Acad Sci.* 2016;113(12):1663–72.
50. Ganai N, Sengupta S, Menon GI. Chromosome positioning from activity-based segregation. *Nucl Acids Res.* 2014;42(7):4145–59.
51. Agrawal A, Ganai N, Sengupta S, Menon GI. Chromatin as active matter. *J Stat Mech Theory Exp.* 2017;2017(1):014001.
52. Agrawal A, Ganai N, Sengupta S, Menon GI. Nonequilibrium biophysical processes influence the large-scale architecture of the cell nucleus. *Biophys J.* 2020;118(9):2229–44.
53. Qi Y, Reyes A, Johnstone SE, Aryee MJ, Bernstein BE, Zhang B. Data-driven polymer model for mechanistic exploration of diploid genome organization. *bioRxiv.* 2020; <https://doi.org/10.1101/2020.02.27.968735>.
54. Stevens TJ, Lando D, Basu S, Atkinson LP, Cao Y, Lee SF, Leeb M, Wohlfahrt KJ, Boucher W, O'Shaughnessy-Kirwan A, et al. 3D structures of individual mammalian genomes studied by single-cell Hi-C. *Nature.* 2017;544(7648):59–64.
55. Paulsen J, Sekelja M, Oldenburg AR, Barateau A, Briand N, Delbarre E, Shah A, Sørensen AL, Vigouroux C, Buendia B, et al. Chrom3D: three-dimensional genome modeling from Hi-C and nuclear lamin-genome contacts. *Genome Biol.* 2017;18(1):21.
56. Di Stefano M, Paulsen J, Lien TG, Hovig E, Micheletti C. Hi-C-constrained physical models of human chromosomes recover functionally-related properties of genome organization. *Sci Rep.* 2016;6:35985.
57. Lajoie BR, Dekker J, Kaplan N. The Hitchhiker's guide to Hi-C analysis: practical guidelines. *Methods.* 2015;72:65–75.
58. Sanders JT, Freeman TF, Xu Y, Golloshi R, Stallard MA, Martin RS, Balajee AS, McCord RP. Radiation-induced DNA damage and repair effects on 3D genome organization. *bioRxiv.* 2019; <https://doi.org/10.1101/740704>.
59. Ay F, Bailey TL, Noble WS. Statistical confidence estimation for Hi-C data reveals regulatory chromatin contacts. *Genome Res.* 2014;24(6):999–1011.
60. Lindsay RJ, Pham B, Shen T, McCord RP. Characterizing the 3D structure and dynamics of chromosomes and proteins in a common contact matrix framework. *Nucl Acids Res.* 2018;46(16):8143–52.
61. Das P, Golloshi R, McCord RP, Shen T. Using contact statistics to characterize structure transformation of biopolymer ensembles. *Phys Rev E.* 2020;101(1):012419.
62. Imakaev M, Fudenberg G, McCord RP, Naumova N, Goloborodko A, Lajoie BR, Dekker J, Mirny LA. Iterative correction of Hi-C data reveals hallmarks of chromosome organization. *Nat Methods.* 2012;9(10):999.
63. Luperchio T, Sauria M, Hoskins V, Wong X, DeBoy E, Gaillard M-C, Tsang P, Pekrun K, Ach R, Yamada N, Taylor J, Reddy K. The repressive genome compartment is established early in the cell cycle before forming the lamina associated domains. *bioRxiv.* 2018; <https://doi.org/10.1101/481598>.
64. Meaburn KJ. Spatial genome organization and its emerging role as a potential diagnosis tool. *Front Genet.* 2016;7:134.
65. Chandra T, Ewels PA, Schoenfelder S, Furlan-Magaril M, Wingett SW, Kirschner K, Thuret J-Y, Andrews S, Fraser P, Reik W. Global reorganization of the nuclear landscape in senescent cells. *Cell Rep.* 2015;10(4):471–83.
66. McCord RP, Nazario-Toole A, Zhang H, Chines PS, Zhan Y, Erdos MR, Collins FS, Dekker J, Cao K. Correlated alterations in genome organization, histone methylation, and DNA–lamin A/C interactions in Hutchinson–Gilford progeria syndrome. *Genome Res.* 2013;23(2):260–9.
67. Fruchterman TM, Reingold EM. Graph drawing by force-directed placement. *Softw Pract Exp.* 1991;21(11):1129–64.
68. Mehta IS, Amira M, Harvey AJ, Bridger JM. Rapid chromosome territory relocation by nuclear motor activity in response to serum removal in primary human fibroblasts. *Genome Biol.* 2010;11(1):1–17.
69. Cremer M, Von Hase J, Volm T, Brero A, Kreth G, Walter J, Fischer C, Solovei I, Cremer C, Cremer T. Non-random radial higher-order chromatin arrangements in nuclei of diploid human cells. *Chromosome Res.* 2001;9(7):541–67.
70. Mora L, Sánchez I, García M, Ponsà M. Chromosome territory positioning of conserved homologous chromosomes in different primate species. *Chromosoma.* 2006;115(5):367–75.
71. Cleveland WS. Robust locally weighted regression and smoothing scatterplots. *J Am Stat Assoc.* 1979;74(368):829–36.
72. Rao SS, Huntley MH, Durand NC, Stamenova EK, Bochkov ID, Robinson JT, Sanborn AL, Machol I, Omer AD, Lander ES, et al. A 3D map of the human genome at kilobase resolution reveals principles of chromatin looping. *Cell.* 2014;159(7):1665–80.
73. Barutcu AR, Lajoie BR, McCord RP, Tye CE, Hong D, Messier TL, Browne G, van Wijnen AJ, Lian JB, Stein JL, et al. Chromatin interaction analysis reveals changes in small chromosome and telomere clustering between epithelial and breast cancer cells. *Genome Biol.* 2015;16(1):214.
74. Nandakumar V, Kelbauskas L, Hernandez KF, Lintecum KM, Senechal P, Bussey KJ, Davies PC, Johnson RH, Meldrum DR. Isotropic 3d nuclear morphometry of normal, fibrocystic and malignant breast epithelial cells reveals new structural alterations. *PLoS One.* 2012;7(1):29230.
75. Javierre BM, Burren OS, Wilder SP, Kreuzhuber R, Hill SM, Sewitz S, Cairns J, Wingett SW, Várnai C, Thiecke MJ, et al. Lineage-specific genome architecture links enhancers and non-coding disease variants to target gene promoters. *Cell.* 2016;167(5):1369–84.
76. Zhu Y, Gong K, Denholtz M, Chandra V, Kamps MP, Alber F, Murre C. Comprehensive characterization of neutrophil genome topology. *Genes Dev.* 2017;31(2):141–53.
77. Sanchez JA, Karni RJ, Wangh LJ. Fluorescent in situ hybridization (fish) analysis of the relationship between chromosome location and nuclear morphology in human neutrophils. *Chromosoma.* 1997;106(3):168–77.
78. Fatakia SN, Kulashreshtha M, Mehta IS, Rao BJ. Chromosome territory relocation paradigm during dna damage response: some insights from molecular biology to physics. *Nucleus.* 2017;8(5):449–60.
79. Ranade D, Koul S, Thompson J, Prasad KB, Sengupta K. Chromosomal aneuploidies induced upon lamin b2 depletion are mislocalized in the interphase nucleus. *Chromosoma.* 2017;126(2):223–44.

80. Todd MJ, Yildirim EA. On Khachiyan's algorithm for the computation of minimum-volume enclosing ellipsoids. *Discrete Appl Math.* 2007;155(13):1731–44.
81. Lloyd S. Least squares quantization in PCM. *IEEE Trans Inform Theory.* 1982;28(2):129–37.
82. Nazeer KA, Sebastian M. Improving the accuracy and efficiency of the k-means clustering algorithm. In: *Proceedings of the World Congress on Engineering*, vol. 1, 2009; pp. 1–3. Association of Engineers London.
83. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R, Dubourg V, Vanderplas J, Passos A, Cournapeau D, Brucher M, Perrot M, Duchesnay E. *Scikit-learn: machine learning in Python.* *J Mach Learn Res.* 2011;12:2825–30.
84. Satopaa V, Albrecht J, Irwin D, Raghavan B. Finding a “kneedle” in a haystack: Detecting knee points in system behavior. In: *2011 31st international conference on distributed computing systems workshops*, 2011; pp. 166–171. IEEE
85. Scherer S. *Short guide to the human genome.* Cold Spring Harbor: Cold Spring Harbor Laboratory Press; 2008.
86. Haeussler M, Zweig AS, Tyner C, Speir ML, Rosenbloom KR, Raney BJ, Lee CM, Lee BT, Hinrichs AS, Gonzalez JN, et al. The ucsc genome browser database: 2019 update. *Nucl Acids Res.* 2018;47(D1):853–8.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more biomedcentral.com/submissions

