

# 2D and 3D Genome Structure Modeling

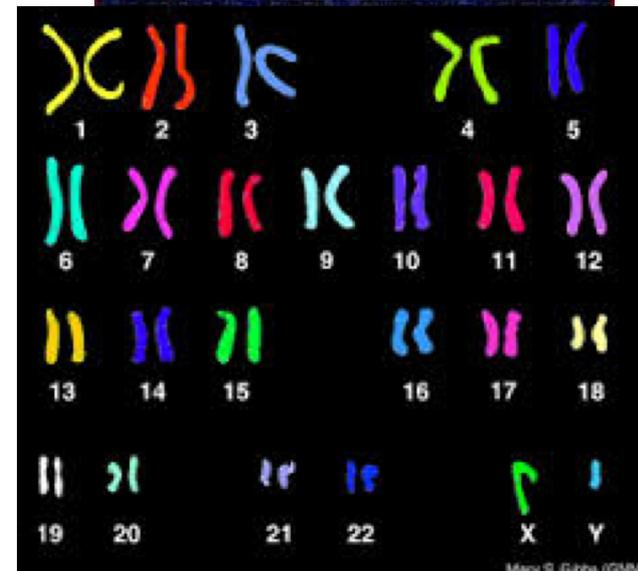
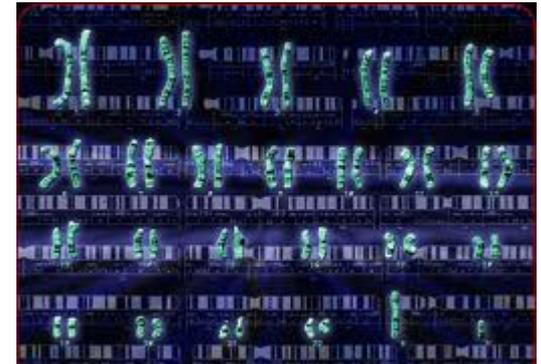
**Jianlin Cheng, PhD**

Department of EECS  
Informatics Institute  
University of Missouri, Columbia  
2018

# Genome – Code of Life



```
CTCCGGACAAATCGATGGCTTTTCG
AATCCGTCGTTTTTTATGCTTTTA
AGGAGAAAGCTATGCGTGTATGCCGT
ATTCCGCGTCCCATGATCCGGCGCCG
ATCGCTCCACATCAACAAAACATAT
GGTAAGAATGAGGGCAATTGTATAG
GCAGATGTGCCCGAGATCGGGTGG
CTATATAACGATCTAACCGTTGGG
GAACAGATGTAGCTCACCATACCT
CGTGGAGAGGTAGCAAACCTGCGC
AAAGAGTGTAAAGCATCGCCCGAT
GCGTAGTTAACTCAAAGGGGAGG
GTGTGAGTAAAGCTTAAAGACGCA
TAAGCTTGAGATCAATAGTTAATT
CGCTGACGCCCATAAAGGCGAGGGG
ACGCCCTGAGCGAATCTAATGGATG
ATTGCTAGGGCTGGGATTCACTTC
ATCGTTTTATCCACACCCAAAGCGAA
```

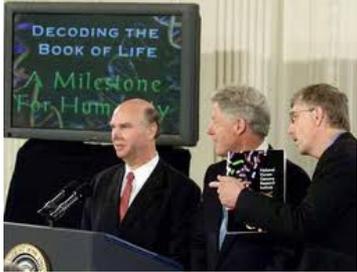


# Genome Sequencing (1D)



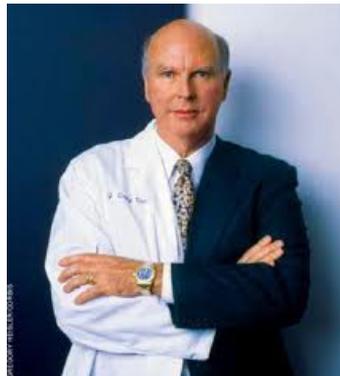
# The Genomic Era

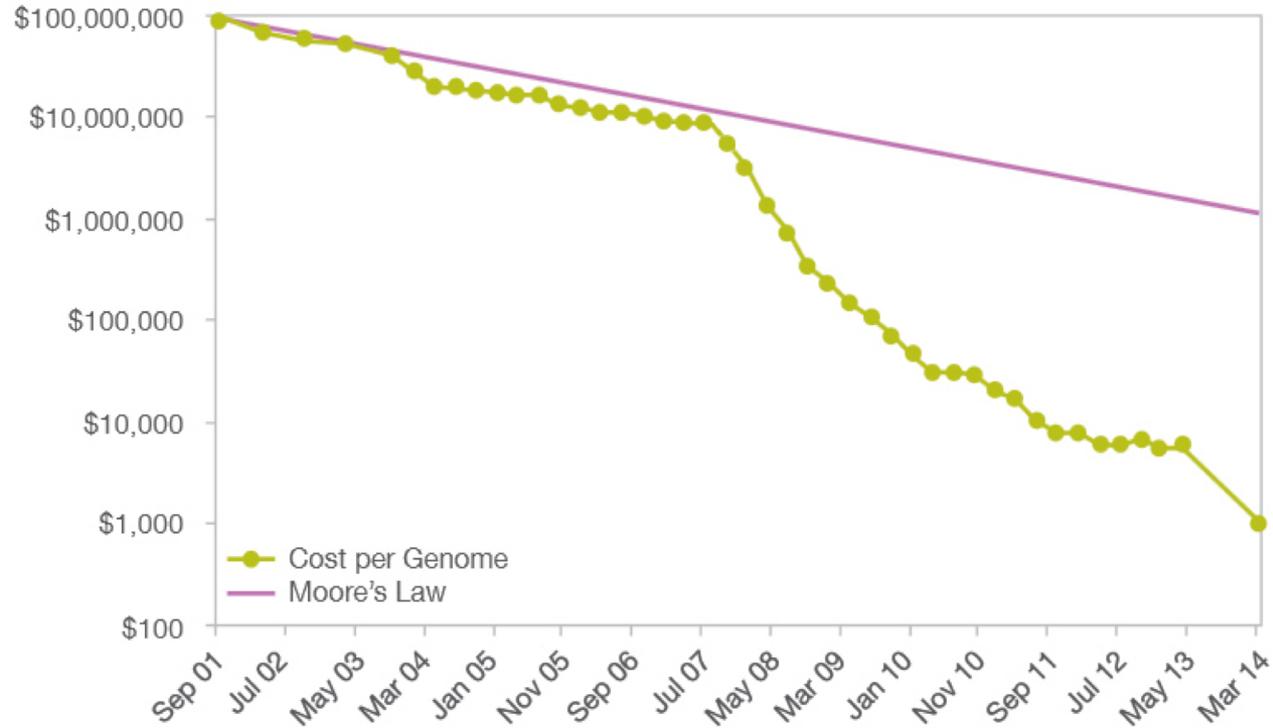
Collins, Venter, Human Genome, 2000



Personal Genome

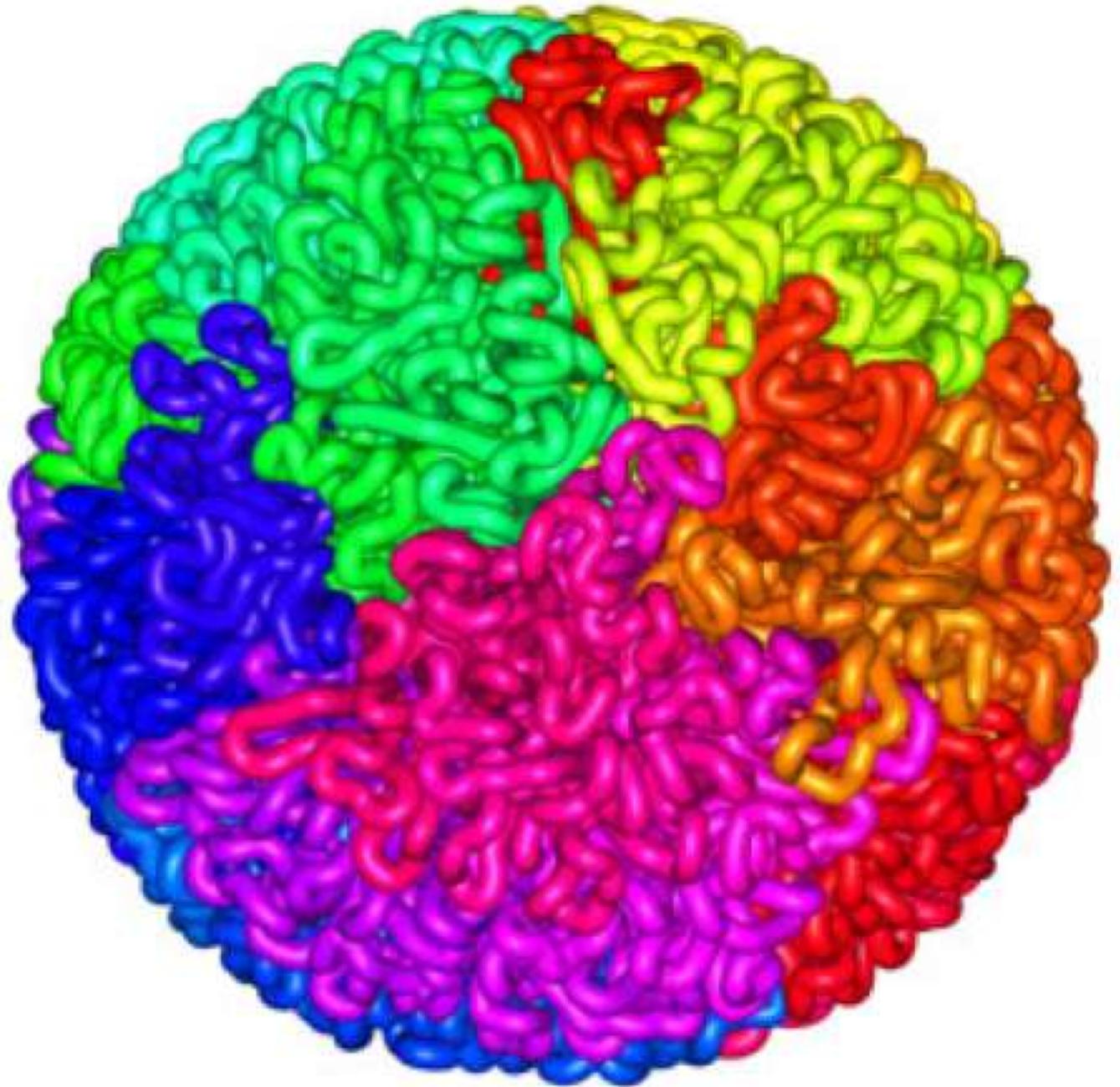
# Sequencing Revolution





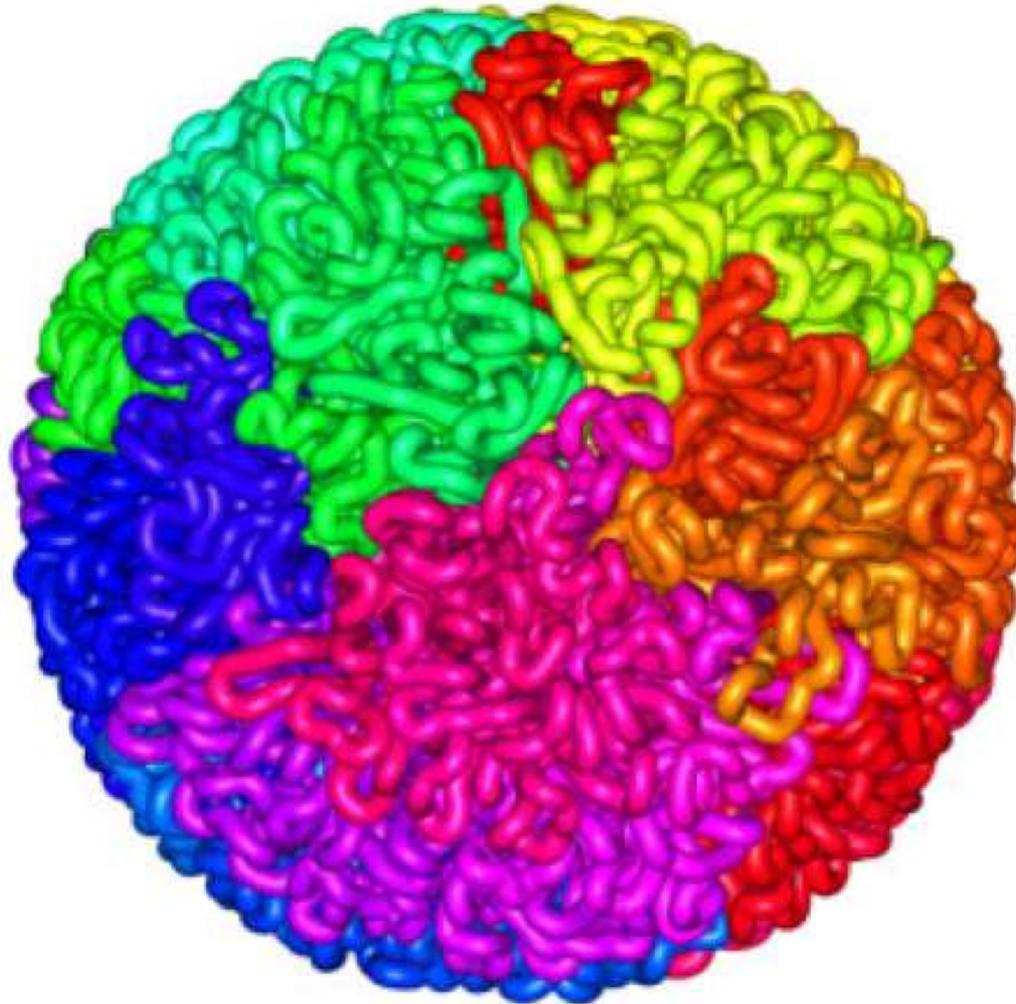
# \$1000 Personal Genome!

# 3D Shape





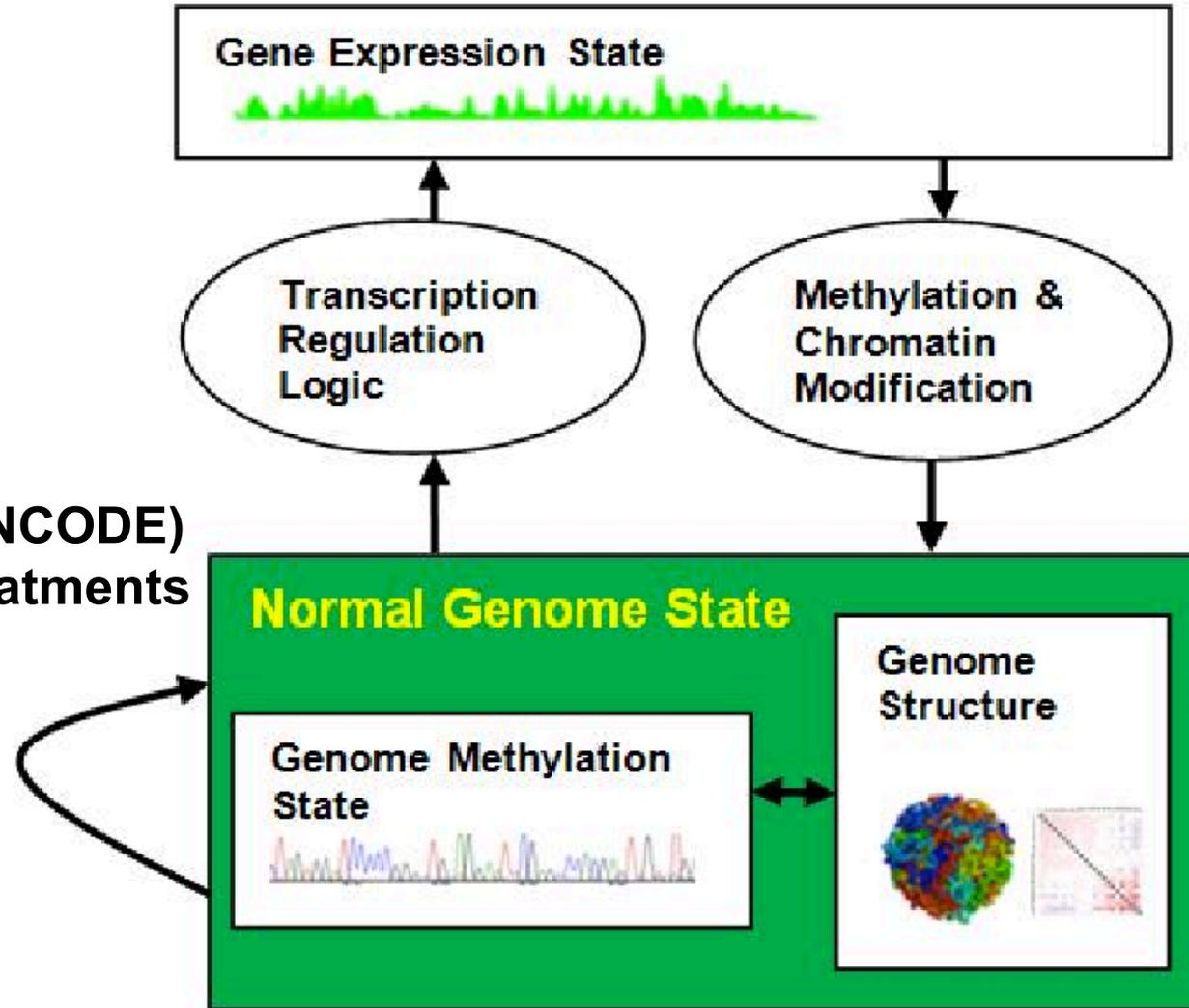
# Genome Conformation



```
ACAACAACCGCTGGTGGTAATAATGCTACCAA  
AGTATGGACAGCAATATGGTCAGCAATATGAA  
AAATGATCAGCAATTCAGTCAGCAATATGCT  
ATAACAGGGCTGTGTATCCCCCCCCCTCAATTC  
SCAACGGGTACAACAATCCCTAATGTAAACGCA  
TGTCAATCCTCCACCCTCAAACACAACATATT  
CTCAATGTACTGGGGCTAGAAAGGCTTTGATT  
ATCAACTGGGTGGTTGATCAATGATGCTCAT  
ACGGTACAGTTCAGATGACATTTGTCATATTA  
TTCOCCTAGGGCTAATATGATTAGGGCCATG  
ATGATTCCTTGGTCCCTCATEATCTGAGCAT  
ACGAAGAAGATGGGATGGATGATTTATATAT  
TTATCGACGATGAAATGCACGATATAATGGTG  
CAGCATTGTTGACTCTTGTGATTCGGGTACA  
DTAAGGGTATTATTAAAGGAGCCCAATATTGG  
CAGCTATTTTCATATGCCACAGGAAACAGGGCT  
TCAAGACCCTTAAAGGAGGTATGGGCAATAAT  
AATTCACAGCAGCATGTTGTTATGTTATCA  
ATGCTGTGAGATGGGCAAAATACAGGTTCA  
CTTTACAACCACAGCAATCATATTATCTCTT  
ATAAGTATTCTCAAAAACCAAAATATCATG  
TTATTATGTAG
```

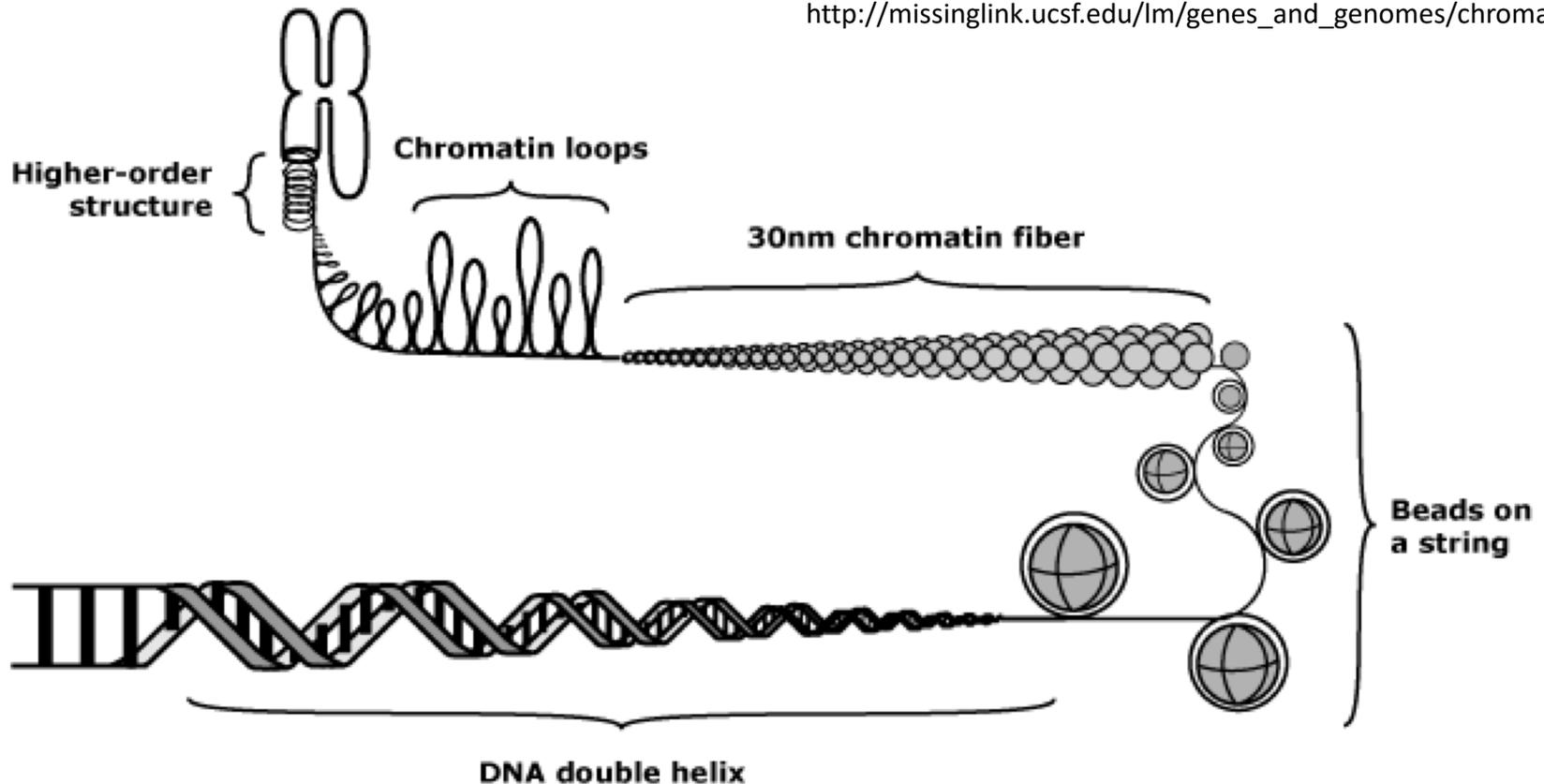
# 3D Genome Structure is Important

- Spatial gene regulation
- Transcription efficiency
- Genome interpretation
- Function implication (ENCODE)
- Disease diagnosis & treatments
- Drug design



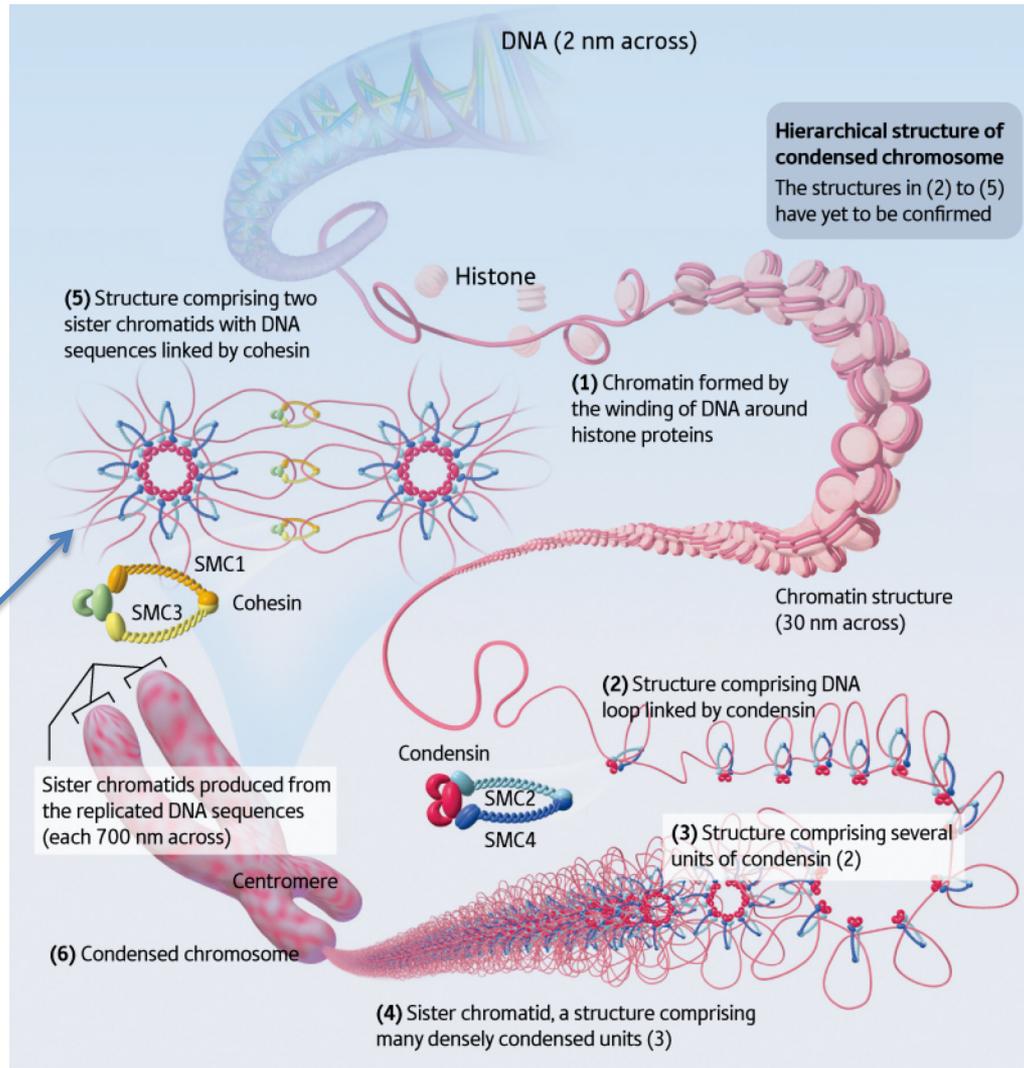
# Multi-Level Chromosome Structure

[http://missinglink.ucsf.edu/lm/genes\\_and\\_genomes/chromatin.html](http://missinglink.ucsf.edu/lm/genes_and_genomes/chromatin.html)



DNA is further packaged. **Nucleosomes** are arranged together into a **fiber** approximately 30 nanometers in diameter. The *precise structure of the chromatin fiber is not known*. Chromatin fiber is further organized into **chromatin loops**, and chromatin loops are further organized into higher-order structures. It has been suggested that **packaging plays a role in gene expression** (gene expression may require associated DNA to open up and acquire an unpackaged conformation). The fully condensed chromosome structure is only seen during mitosis.

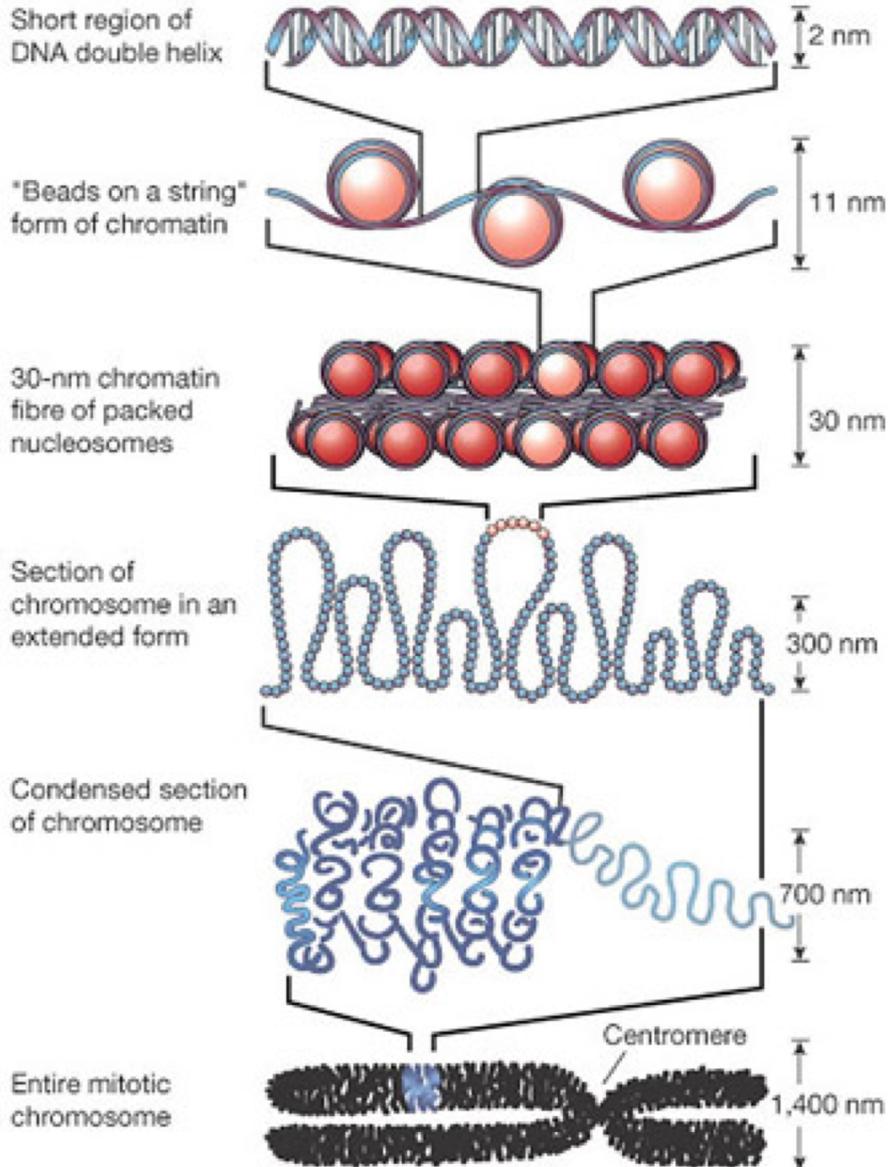
# Model 1 (Riken)



**Globular Structure**

The DNA is a remarkable molecule in many ways. It encodes our entire genome, and if stretched out in a thin thread would measure **1.8 m** in length.

# Size of Elements

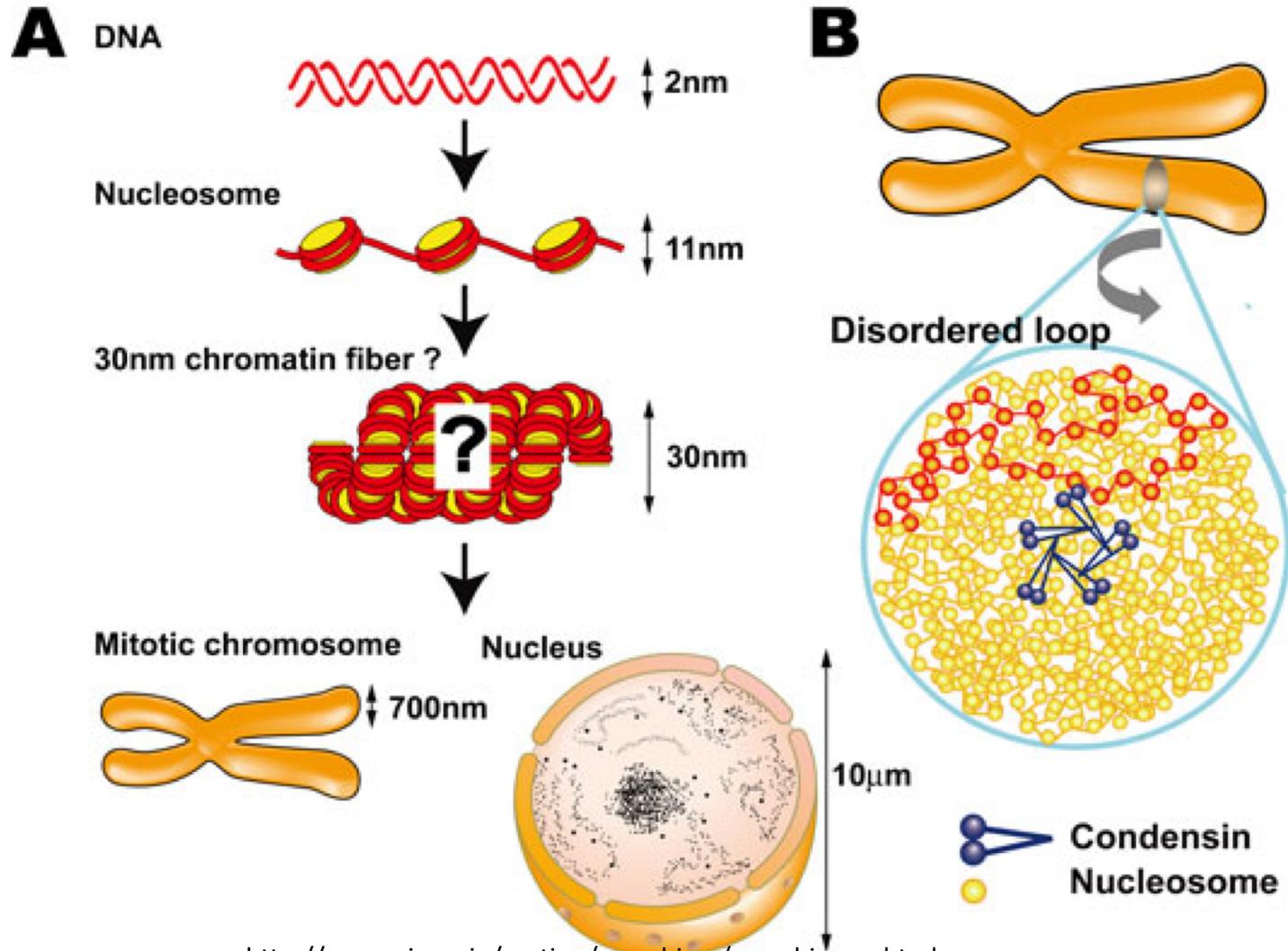


**NM:  $10^{-9}$  meter**

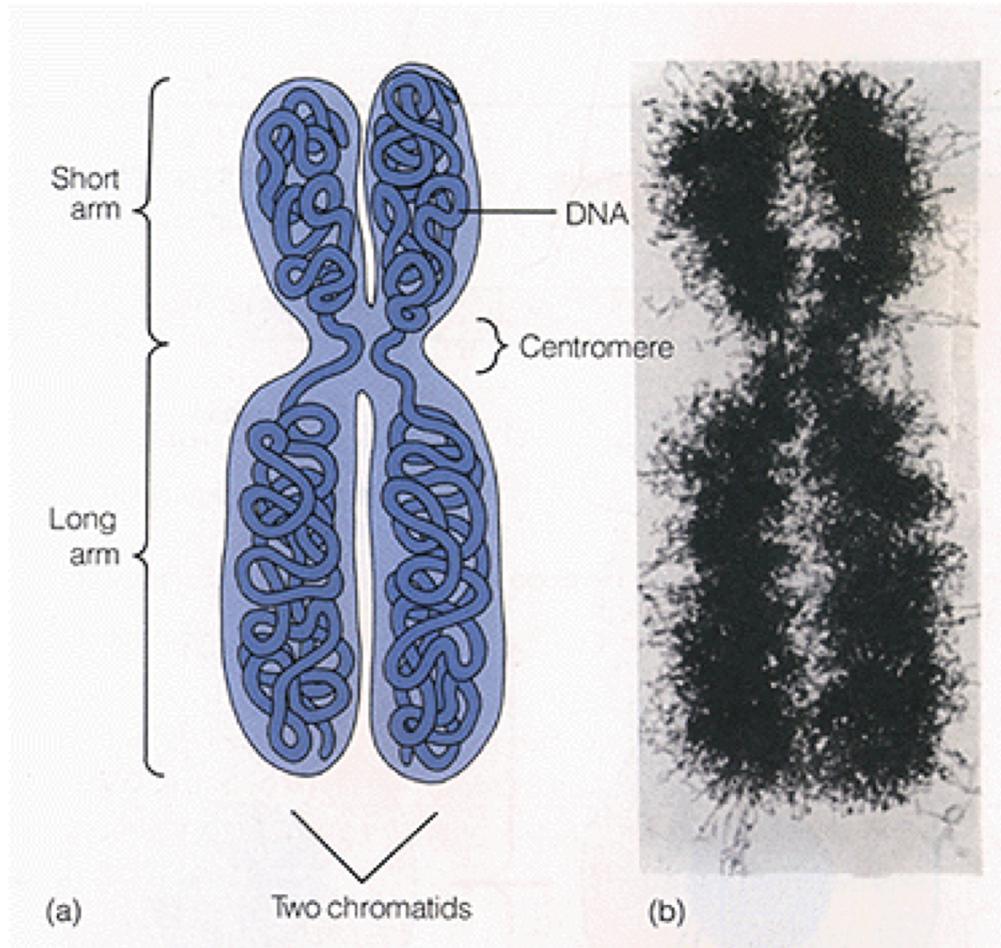
A complex of DNA and basic proteins (such as **histone**) in eukaryotic cells that is condensed into chromosomes in mitosis and meiosis.

There are two types. **Heterochromatic** is densely-coiled chromatin that appears as nodules in or along chromosomes and contains relatively few genes. **Euchromatic** is the less-coiled and genetically active portion of chromatin that is largely composed of genes.

# Size of Elements



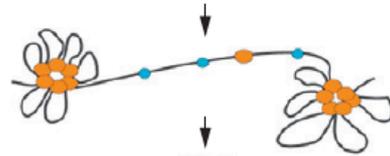
# Chromosomes



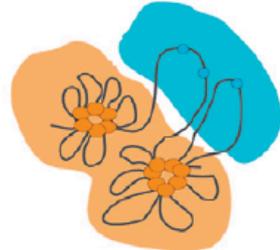
<http://www.copernicusproject.ucr.edu/ssi/HSBiologyResources.htm>

# Two Compartments

Linear chromosome:  
Active genes: orange  
Inactive genes: blue



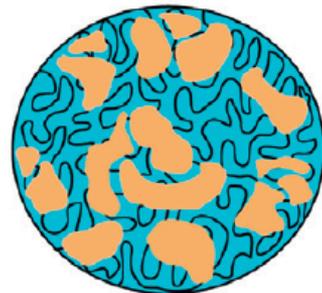
Chromatin looping between active genes and regulatory elements and clustering of genes at the site of active transcription facilitates formation of chromatin globules.



Active gene cluster associate with other expressed genes in active neighborhoods while inactive genes cluster in silent neighborhoods.



Active and silent neighborhoods associate in cis and in trans to form larger active and inactive compartments.

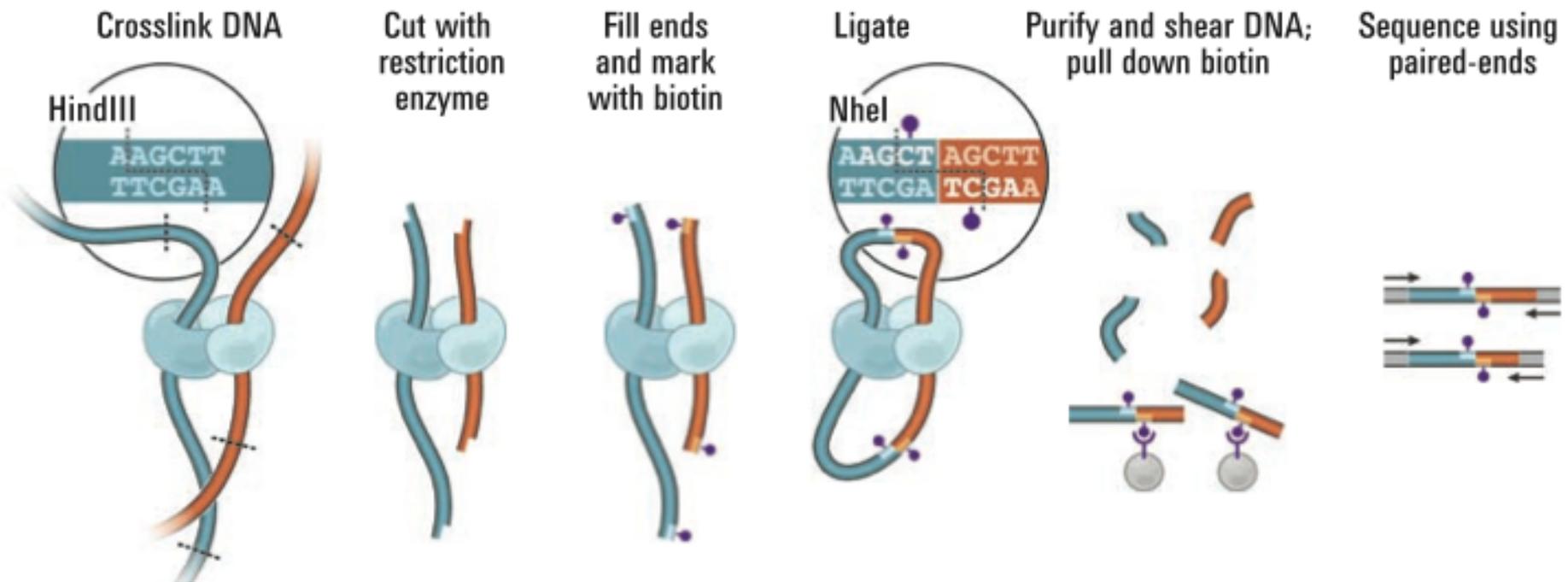


Nuclear organization reflects clustering of active and inactive loci in distinct compartments forming a fractal globule.

# Chromosome Conformation Capturing Techniques

# Chromosome Conformation Capturing (Hi-C)

**A**



# Hi-C Protocol

- Cells are cross-linked with formaldehyde
- DNA is digested with a restriction enzyme and ligated, resulting enriched with cross-linked elements with a biotin marked at the junction
- Shearing DNA and selecting biotin-containing fragments to create a Hi-C library
- Sequence the library to create a catalog of interacting fragments

# Genomic Spatial Interaction (Contact) Data



```

ATGTATCCAGGTAGTGGAGTTACACCTACAACAACCTGGTGGTAATAATGGCTACCAA
CGGCCCATGGCTCCTCCCTTAACCAGCAGTATGGACAGCAATATGGTCAGCAATATGAA
CASCAGTATGGACAGTAAATATGGGCAACAATAATGATCAACAATTCAATCAGTCASCAATATGCT
CCACCACCAGSTCCCTCCCTATGGCTTATAACAGGCTTGTGATCCCTCCCTCAATTC
CAGCAGGAACAGGTAAGGGCACAAATTAAGCAACGGCTACACAATCCCTAATGTAAACGCA
TCCAATATGTAAGTCCACCCAGAAATATGTCATTACCTCCCTCCAAACACAAACTATF
CAAGGTACAGACCACCTTATCAGTATTCTCAATGTACTGGGCTAGAAAAGGCTTTGATT
ATCGGTATAAACTACATAGSTTCAAAAAATCAACTGCGTGGTGTATCAATGATGCTCAT
AACATCTTCAACTTTTTTACTAATGGSTACGGTTACAGTTCAGATACATTGTTCATATTA
ACTGATGATCAGAACGATTTGGTCAGGGTTCCCACTAGGGCTAATGATGATTAGGGCCATG
CAATGGTTGGTCAAGSATGGCCAAACCCAATGATTCCTTGGTCCCTCAATATTCTGGACAT
GGTGGCCAAACTGAAGATTGGATGSGGACGAAGAAGATGGGATGGATGATGTTATATAT
CCGSTCGATTTCGAAACTCAAGGGCCCAATTATCGACGATGAAATGCCACTATAATGGTG
AAGCCCTTACAACAAGGTGTTAGACTAACAGCATTGTTTACTCTTGTCAATTCGGGTACA
GTGTTGGATCTTCCATATACCTATTCTACTAAAGGGTATTATTAAGGAGCCCATATTGG
AAGSATGTTGGCCAAAGATGGCTGCAAGCAGCTATTTTCATATGCCACAGGAACAGGGCT
GCTTTSATTTGGTCTTTAGSTTCTATATTCAAGACCSTTAAGGSAAGGTATGGGCAATAAT
GTGGATAGAGAACCSTGAGACAGATCAAATTCACSCASCAGATGTTGTTATGATATCA
GGTTCGAAGGATAATCAAACCTTCTGCAAGATGCTGTGCAAGATGGGCAAAATACACTTCA
ATGTCOCACGGCTTCATCAAGGTTATGACTTACAACCACAGCAATCATATTTATCTTT
TTACAGAACATGAGGAAAGAAATGGCTGGTAAGTATCTCAAAAAACCAATTATCATCG
TCACACCCTATTGACGTAATCTGCAATTTATTTTATAG
    
```

Map reads to the human genome sequence

1 2 ..... n

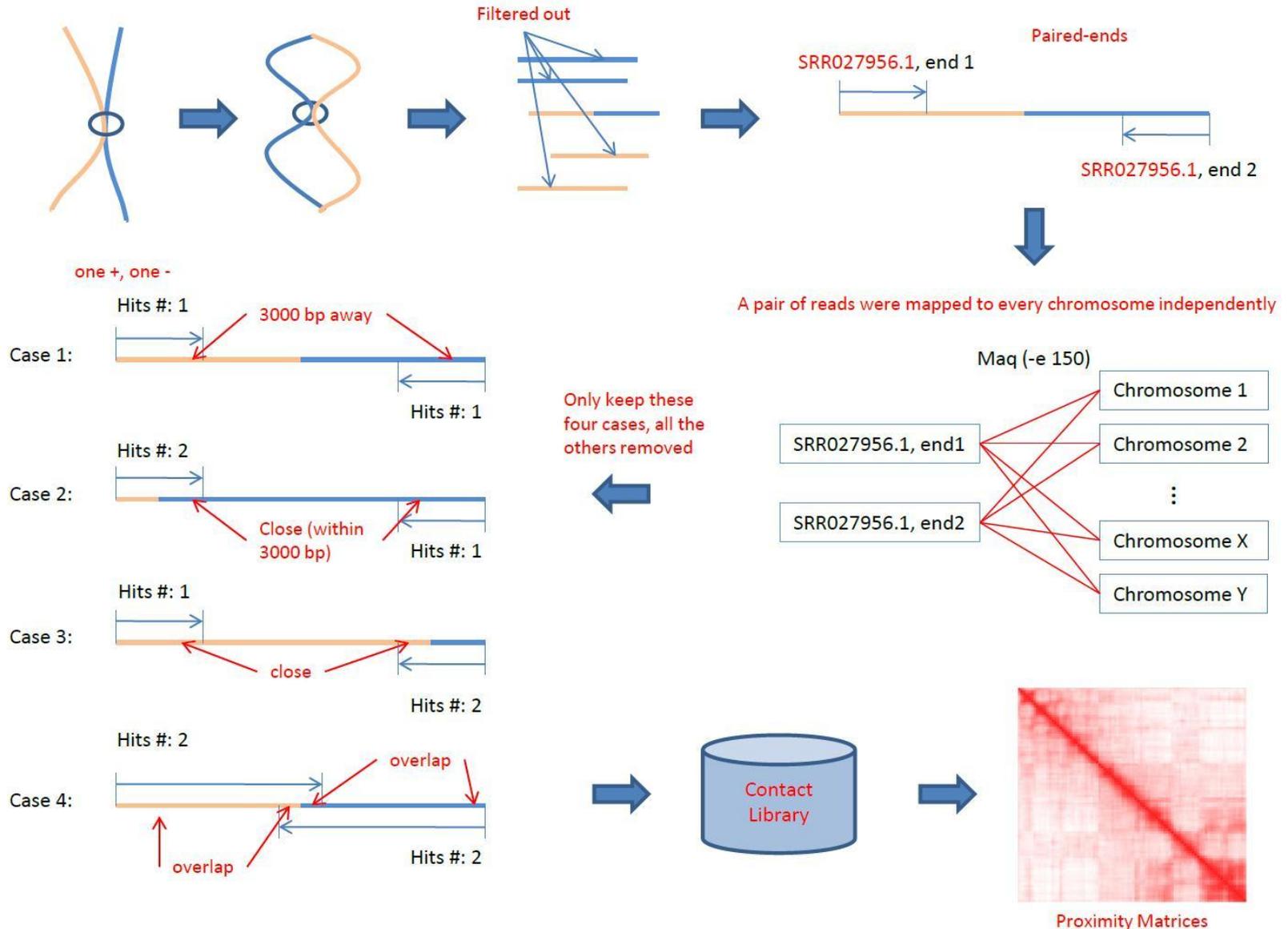
1	40				1
2		3		8	
.			10		
.		8		15	
n	1				

$C[i,j] = \# \text{ contacts}$

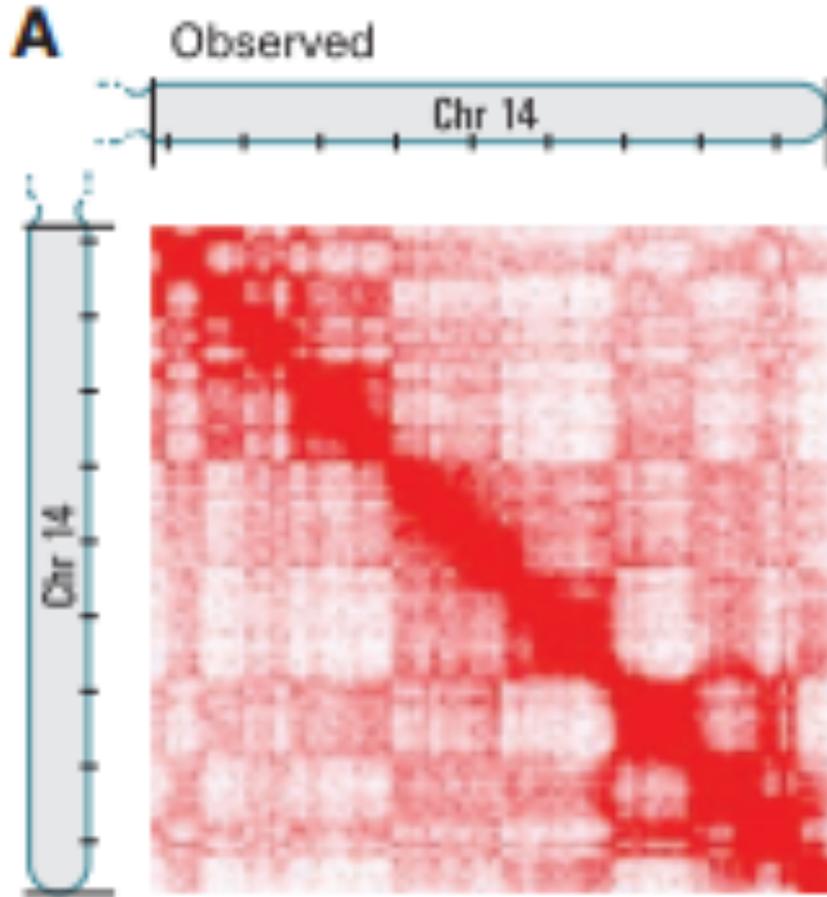
Intra- / inter-chromosome contact map

Wang et al., 2013

# Hi-C Data Analysis Pipeline

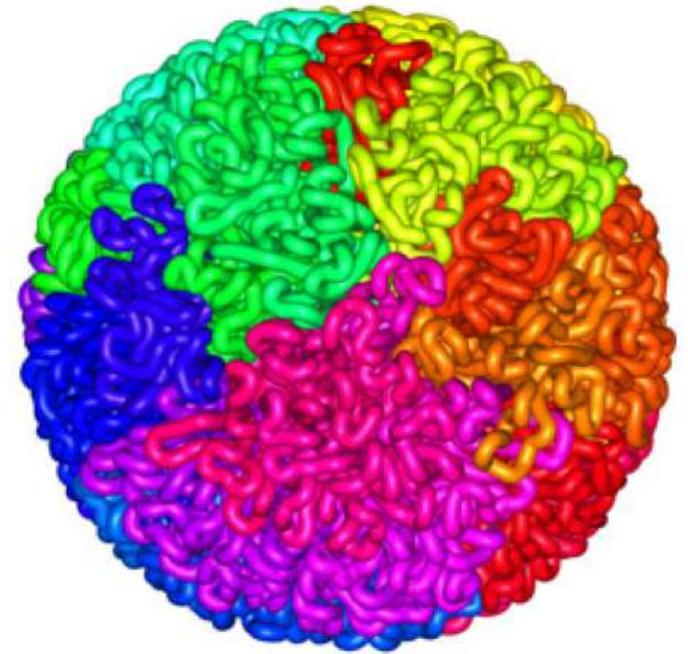


# 2D Chromosome Contact Map



**Chromosome  
Conformation  
Capturing**

# Construct 3D Shape of Genome



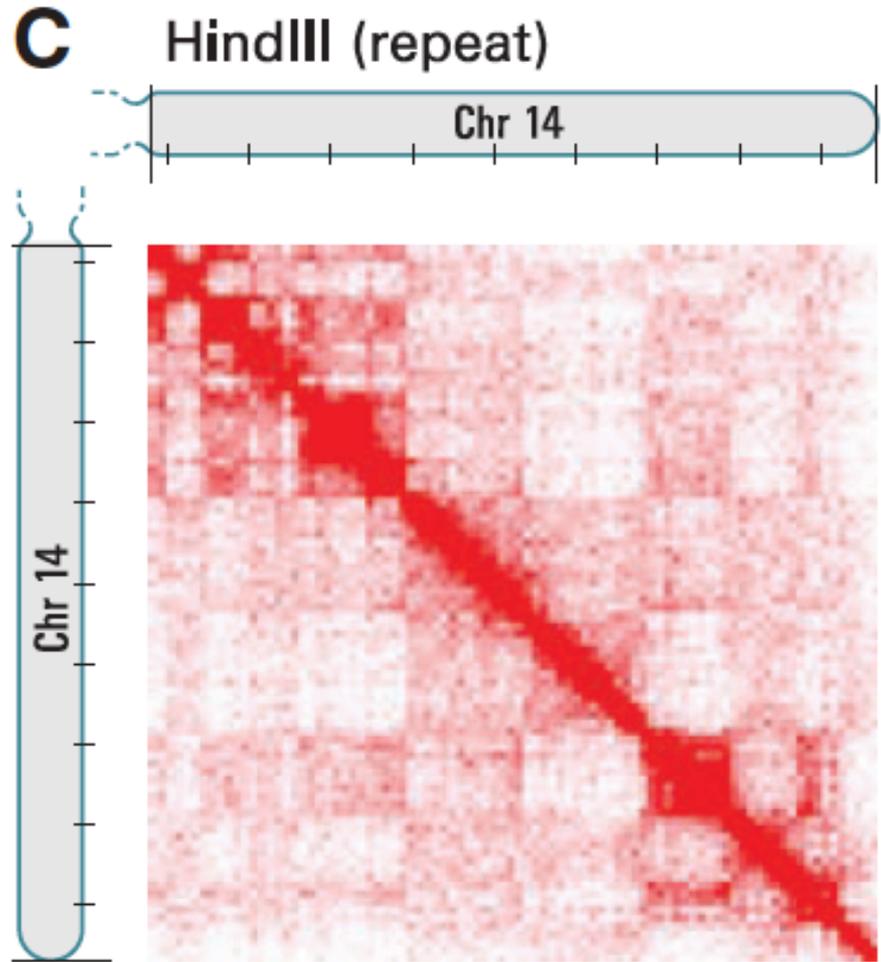
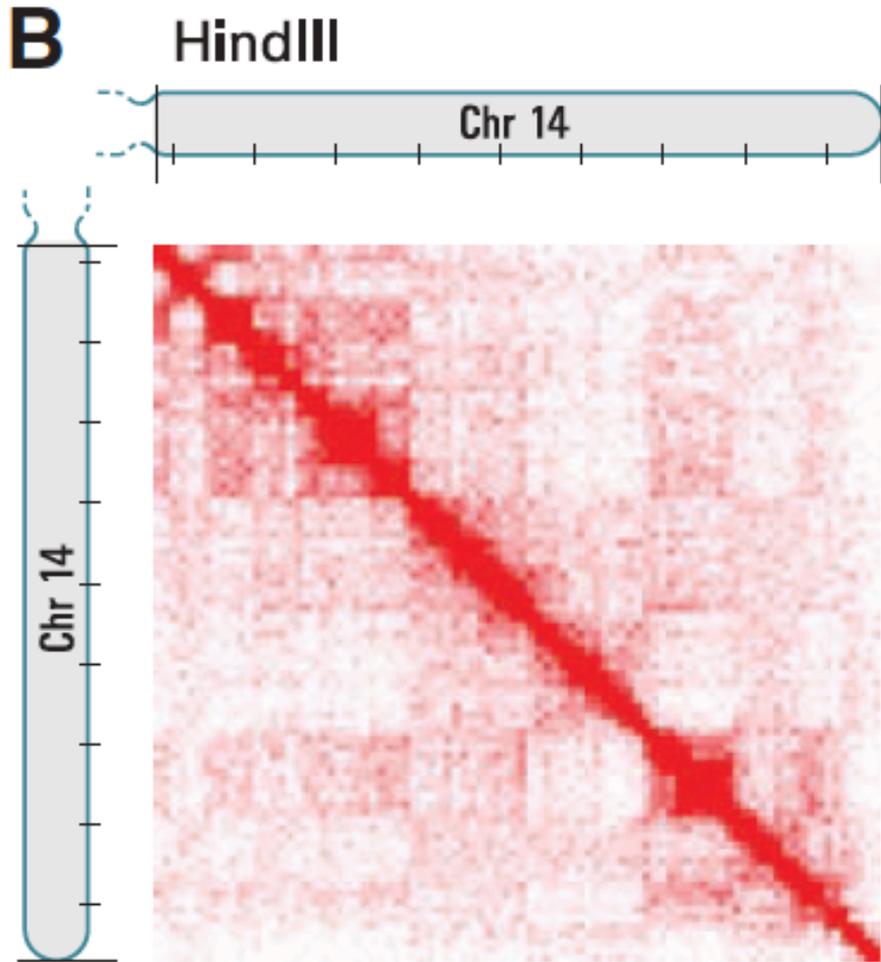
Images.google.com

# Data Set I

- A normal human lymphoblastoid cell line (B-cell)
- 8.4 million read pairs uniquely mapped to the human genome reference sequence
- 6.7 million corresponded to long-range contacts between segments  $> 20$  kb apart

# Genome Wide Contact Map

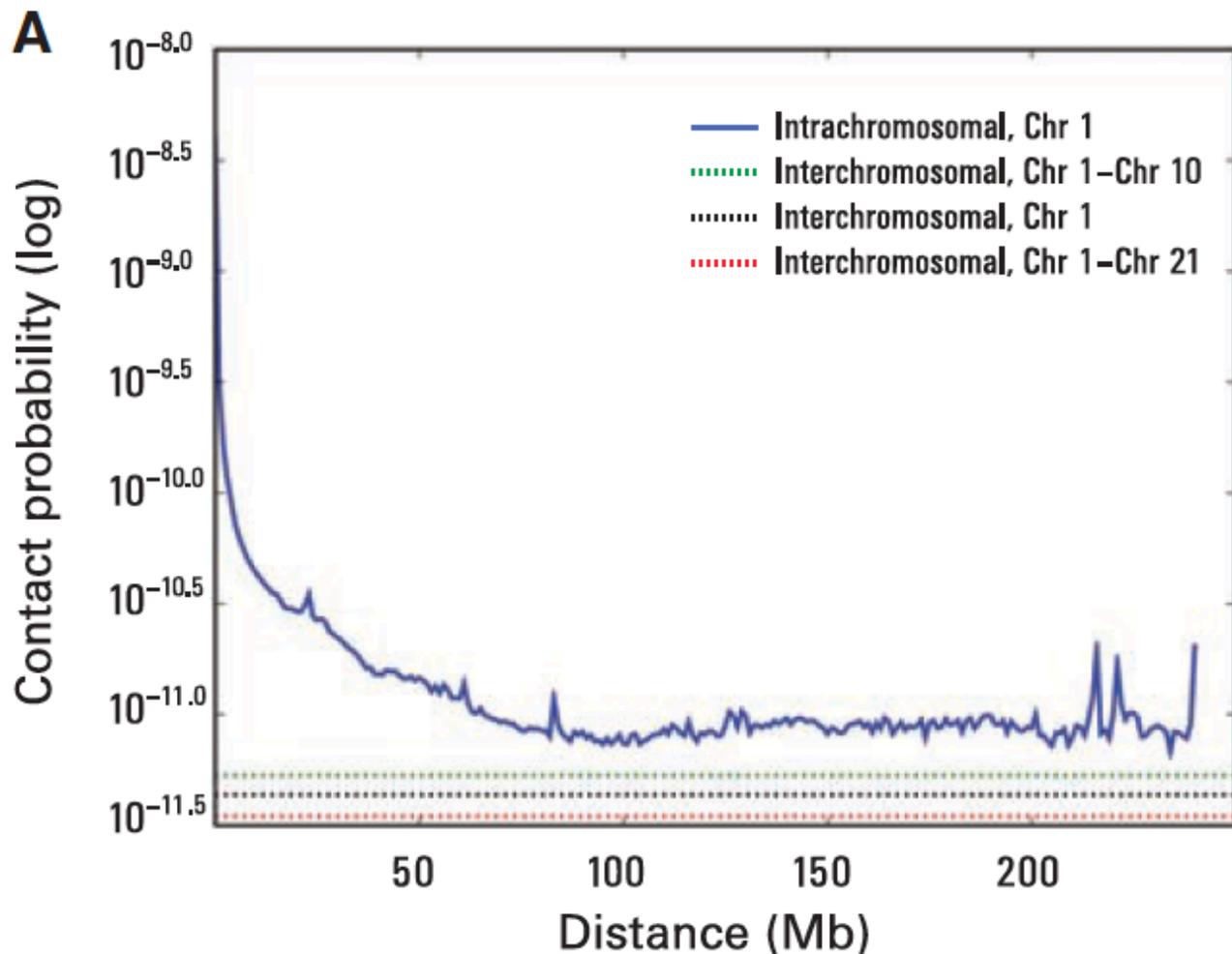
- Divide genome into 1-Mb regions (loci)
- $M_{ij}$ : number of contacts between loci  $i$  and  $j$
- The matrix reflects an ensemble average of the interactions in the original sample of cells
- Represented as a heat map



(B) Hi-C produces a genome-wide contact matrix. The submatrix shown here corresponds to **intrachromosomal interactions on chromosome 14**. (Chromosome 14 is acrocentric; the short arm is not shown.) Each pixel represents all interactions between a 1-Mb locus and another 1-Mb locus; intensity corresponds to the total number of reads (0 to 50). Tick marks appear every 10 Mb. (C) We compared the original experiment with results from a biological repeat using the same restriction enzyme [(C), range from 0 to 50 reads]. Correlation is 0.99.

# Relation between Euclidean Distance and Genomic Distance

- Average intrachromosomal contact probability  $I_n(s)$  for pairs of loci separated by a genomic distance  $s$  on chromosome  $n$ .
- $I_n(s)$  decreases monotonically on every chromosome
- Even at distances  $> 200$  Mb,  $I_n(s)$  is always much greater than the average contact probability between different chromosomes

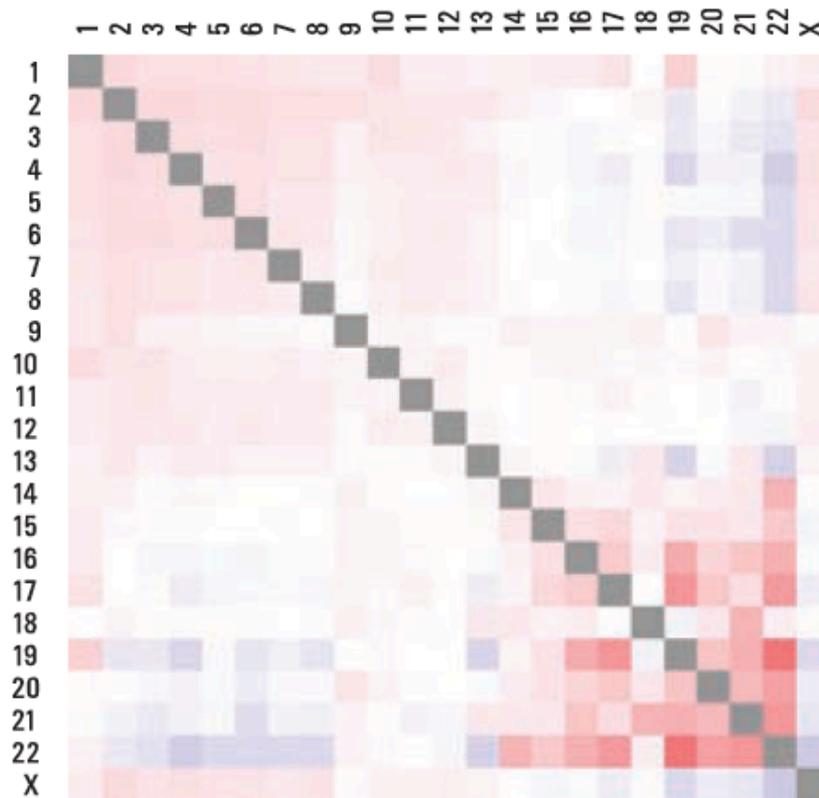


**Implication:  
Chromosome  
territory**

Probability of contact decreases as a function of genomic distance on chromosome 1, eventually reaching a plateau at ~90 Mb (blue). The level of interchromosomal contact (black dashes) differs for different pairs of chromosomes; loci on chromosome 1 are most likely to interact with loci on chromosome 10 (green dashes) and least likely to interact with loci on chromosome 21 (red dashes). Interchromosomal interactions are depleted relative to intrachromosomal interactions.

**B**

Human chromosomes



### Observed/expected number of interchromosomal contacts

between all pairs of chromosomes. Red indicates enrichment, and blue indicates depletion (range from 0.5 to 2). Small, gene-rich chromosomes tend to interact more with one another, suggesting that they cluster together in the nucleus.

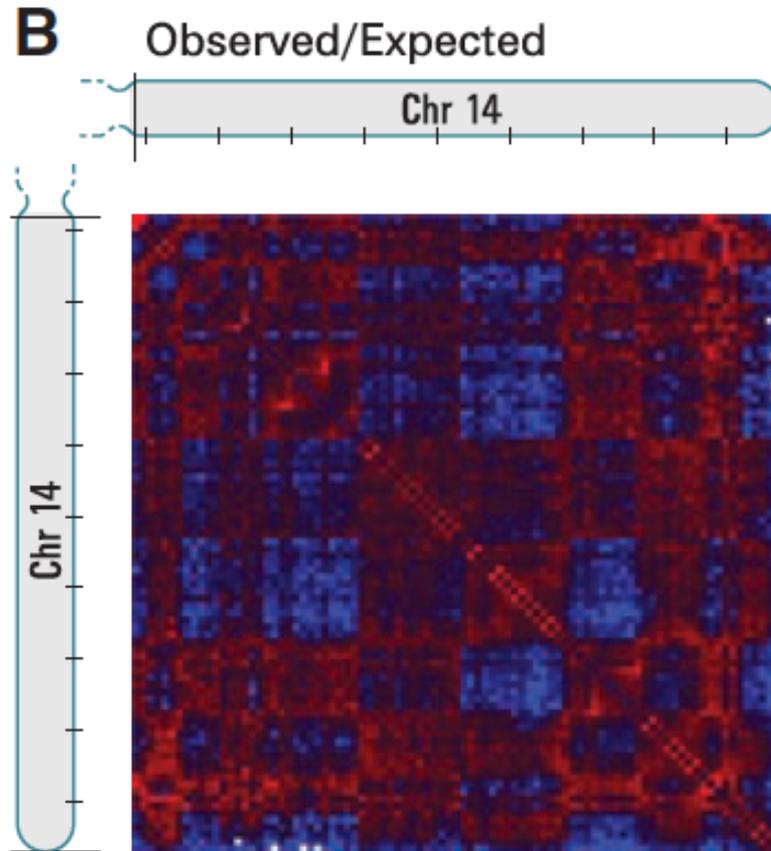
### Human chromosomes

expected number of contacts between chromosome  $i$  and  $j$  was calculated by:

$$E_{i,j} = R_i \times R_j \times N_{INTER},$$

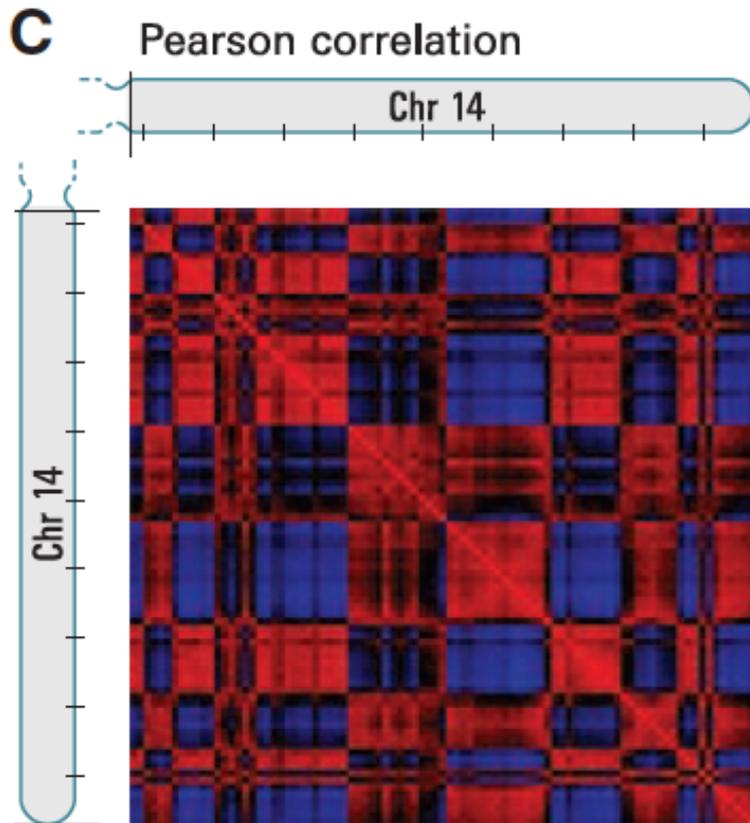
where  $R_i$  and  $R_j$  are the fractions of inter-chromosomal reads associated with  $i$  and  $j$ , respectively, and  $N_{INTER}$  is the total number of inter-chromosomal reads for a cell sample. The actual observed number of inter-chromosomal contacts between chromosomes  $i$  and  $j$  divided by the expected number  $E_{i,j}$  indicates the enrichment or depletion of inter-chromosomal contacts between them.

# Normalizing Contact Map by Expected Number of Contacts at Genomic Distance



Dividing each entry in the contact matrix ( $M$ ) by the genome-wide average contact probability for loci at that genomic distance. The normalized matrix shows many large blocks of enriched and depleted interactions, generating a plaid pattern

# Pearson Correlation Map

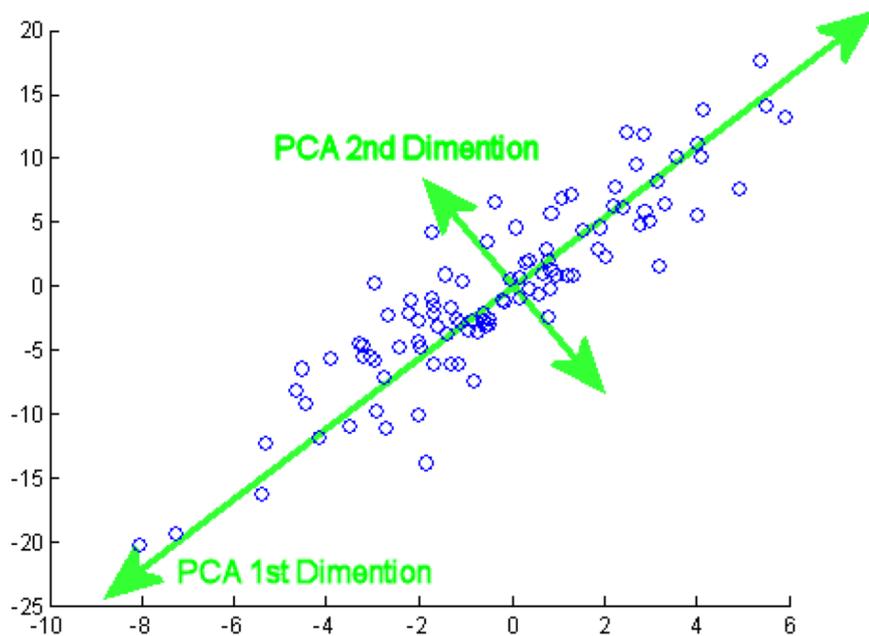


If two loci (here 1-Mb regions) are nearby in space, we reasoned that they will share neighbors and have correlated interaction profiles.

This process dramatically **sharpened the plaid pattern**; 71% of the resulting matrix entries represent statistically significant correlations ( $P \leq 0.05$ ).

The **plaid pattern** suggests that each chromosome can be decomposed into two sets of loci (arbitrarily labeled A and B) such that contacts within each set are enriched and contacts between sets are depleted.

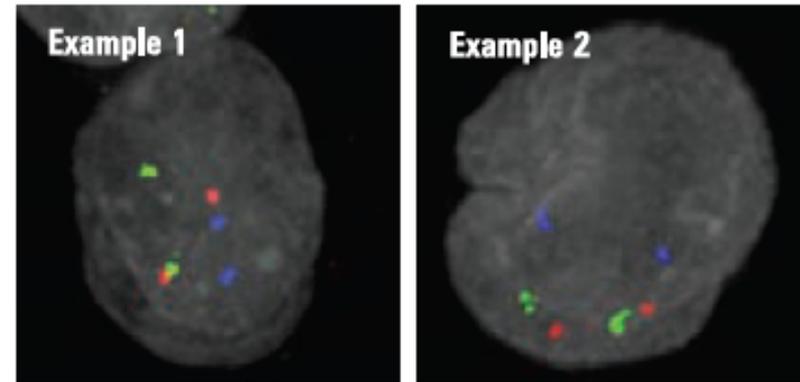
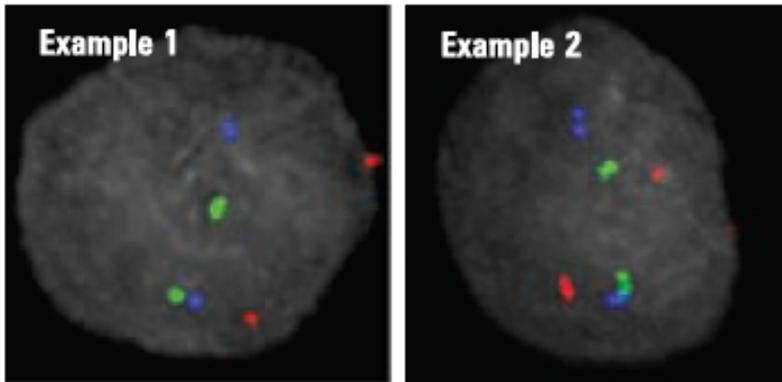
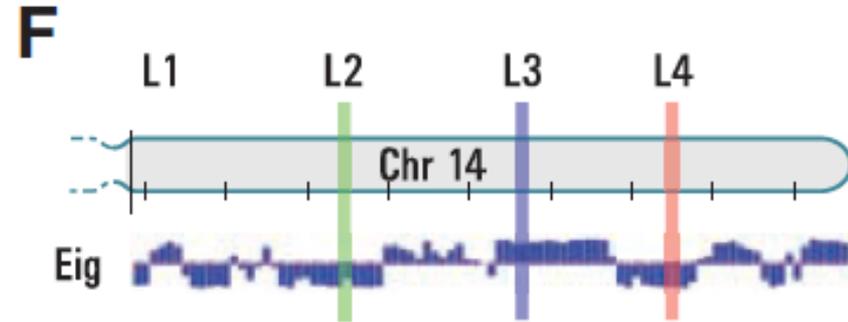
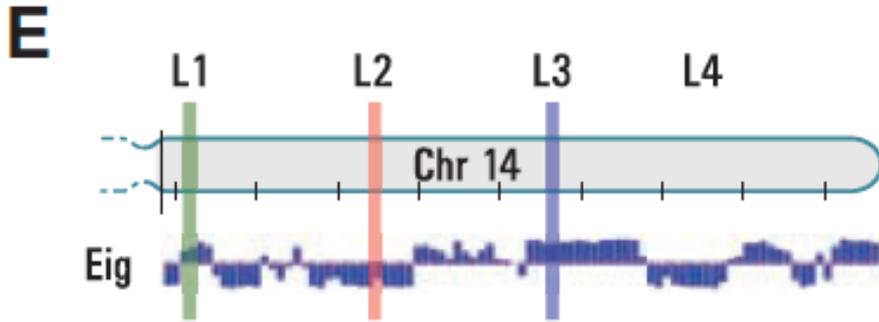
# Principle Component Analysis



The first two principal components (PC) clearly corresponded to the plaid pattern (positive values defining one set, negative values the other).

The entire genome can be partitioned into two spatial compartments such that greater interaction occurs within each compartment rather than across compartments.

# FISH Validation of Two Compartments



L1 – L3: Compartment A  
L2 – L 4: Compartment B

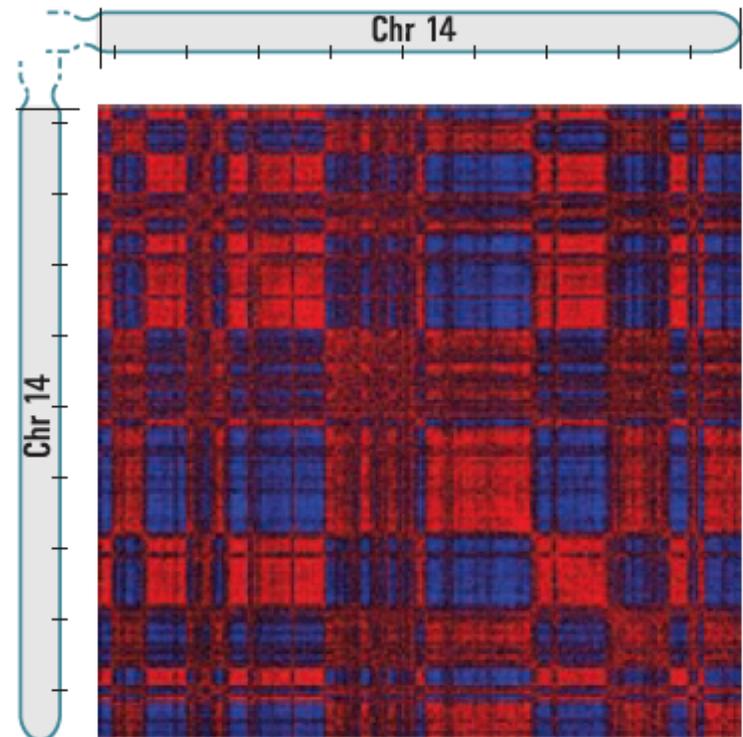
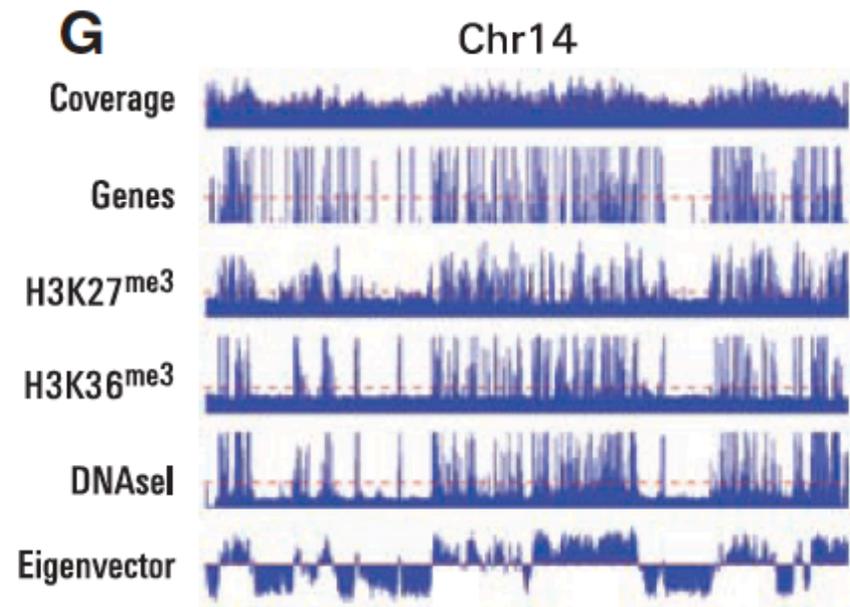
# # Contacts and Physical Distance Measured by FISH

- A strong correlation was observed between the number of Hi-C reads  $m_{ij}$  and the 3D distance between locus  $i$  and locus  $j$  as measured by FISH [Spearman's  $r = -0.916$ ,  $P = 0.00003$ ], suggesting that Hi-C read count may serve as a proxy for distance.
- Pairs of loci in compartment B showed a consistently higher interaction frequency at a given genomic distance than pairs of loci in compartment A. This suggests that compartment B is more densely packed.

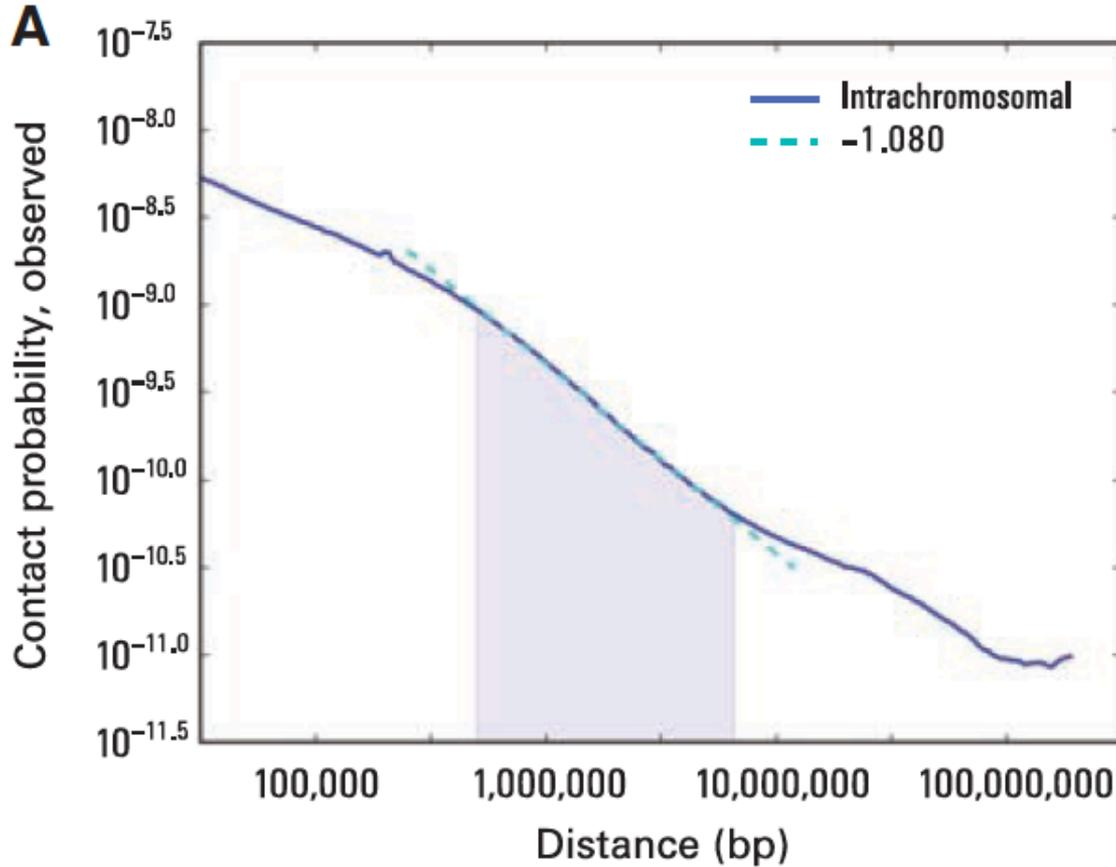
# Open More Accessible, Gene Rich Compartment VS less Accessible

H3K27/H3K36: activating/repressive chromatin Marks

DNAseI: deoxyribonuclease I, sensitivity, measure chromatin accessibility



# Contact Probability VS Genomic Distance (Power Law Distribution)



Contact probability as a function of genomic distance averaged across the genome (blue) shows a power law scaling between 500 kb and 7 Mb (shaded region) with a slope of  $-1.08$  (fit shown in cyan).

When plotted on log-log axes,  $I(s)$  exhibits a prominent power law scaling between  $\sim 500$  kb and  $\sim 7$  Mb, where contact probability scales as  $s^{-1}$ . This range corresponds to the known size of open and closed chromatin domains.

# Genome 3D Model

- **Power-law** dependencies can arise from polymer like behavior.
- **Equilibrium globule:** a compact, densely knotted configuration originally used to describe a polymer in a poor solvent at equilibrium
- **Fractal Model:** This highly compact state is formed by an unentangled polymer when it crumples into a series of small globules in a “beads-on-a-string” configuration. These beads serve as monomers in subsequent rounds of spontaneous crumpling until only a single globule of globules-of-globules remains.

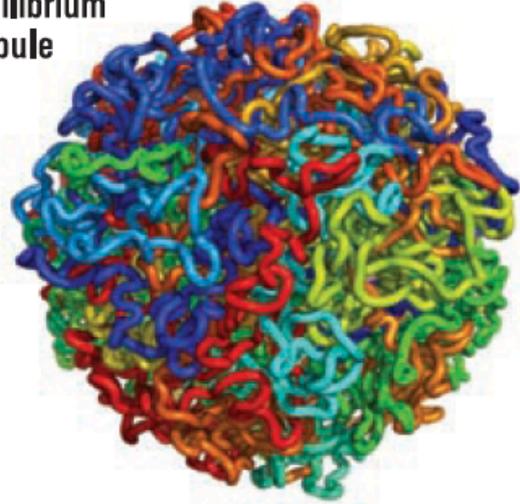
# Fractal Model

- Lack knots and would facilitate unfolding and refolding, for example, during gene activation, gene repression, or the cell cycle.
- In a fractal globule, contiguous regions of the genome tend to form spatial sectors whose size corresponds to the length of the original region.
- In contrast, an equilibrium globule is highly knotted and lacks such sectors; instead, linear and spatial positions are largely decorrelated

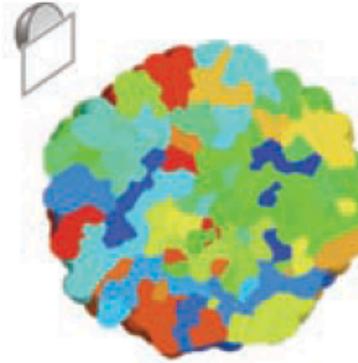
# Comparison of Two Models

**FOLDED POLYMER**

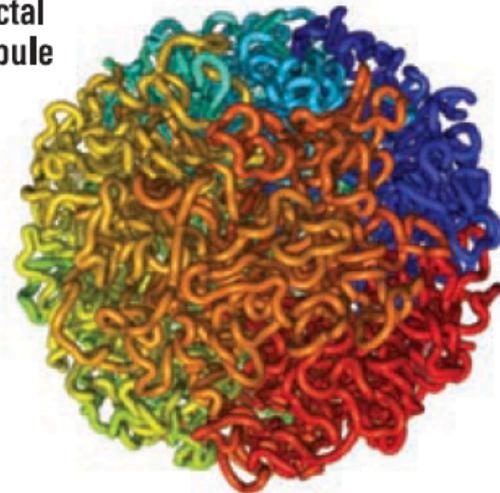
**Equilibrium  
globule**



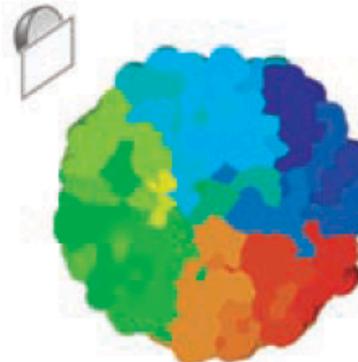
**Cross-section view**



**Fractal  
globule**



**Cross-section view**

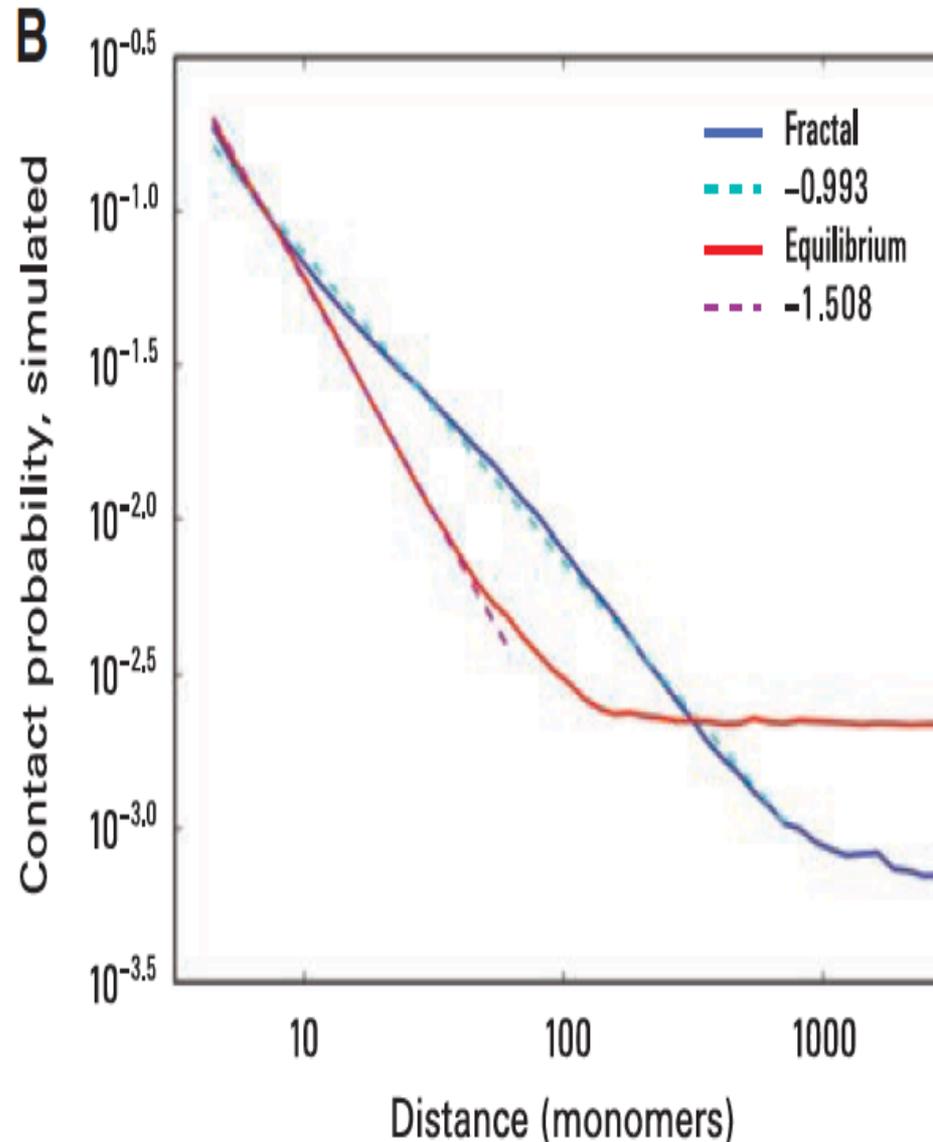


# Consistency Checking

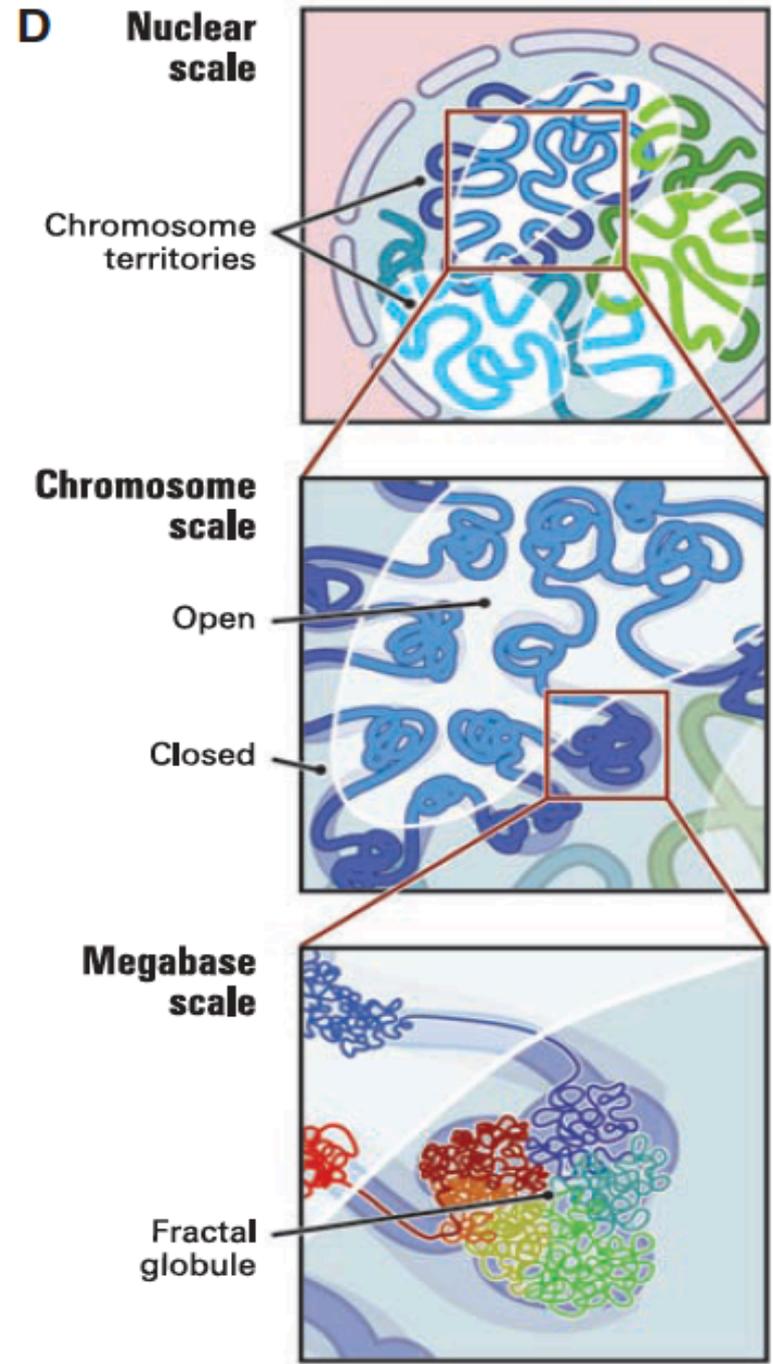
- The **equilibrium globule model predicts that contact probability will scale as  $s^{-3/2}$** , which we do not observe in our data.
- **We analytically derived the contact probability for a fractal globule and found that it decays as  $s^{-1}$** ; this corresponds closely with the prominent scaling we observed ( $s^{-1.08}$ ).
- **3D distance between pairs of loci:  $s^{1/2}$  for an equilibrium globule,  $s^{1/3}$  for a fractal globule.** Although 3D distance is not directly measured by Hi-C, we note that a recent paper using 3D-FISH reported an  $s^{1/3}$  scaling for genomic distances between 500 kb and 2 Mb

# MCMC Simulation Validation

- Monte Carlo simulations to construct ensembles of fractal globules and equilibrium globules (500 each).
- Contact probability (for fractal globules,  $s^{-1}$ , and for equilibrium globules,  $s^{-3/2}$ )
- 3D distance (for fractal globules  $s^{1/3}$ , for equilibrium globules  $s^{1/2}$ )
- Lack of entanglements and the formation of spatial sectors within a fractal globule.



# Multi-Scale Fractal Model



# Hi-C Data Analysis II

Z. Wang, R. Cao, K. Taylor, A. Briley, C. Caldwell, J. Cheng. **The Properties of Genome Conformation and Spatial Gene Interaction and Regulation Networks of Normal and Malignant Human Cell Types.** PLoS ONE. 2013

# Data Sets

- Primary human acute lymphoblastic leukemia (B-ALL) B-cell
- The MHH-CALL-4 B-ALL cell line (CALL4)
- The follicular lymphoma cell-line (RL)
- Sequenced by Illumina HiSeq 2000

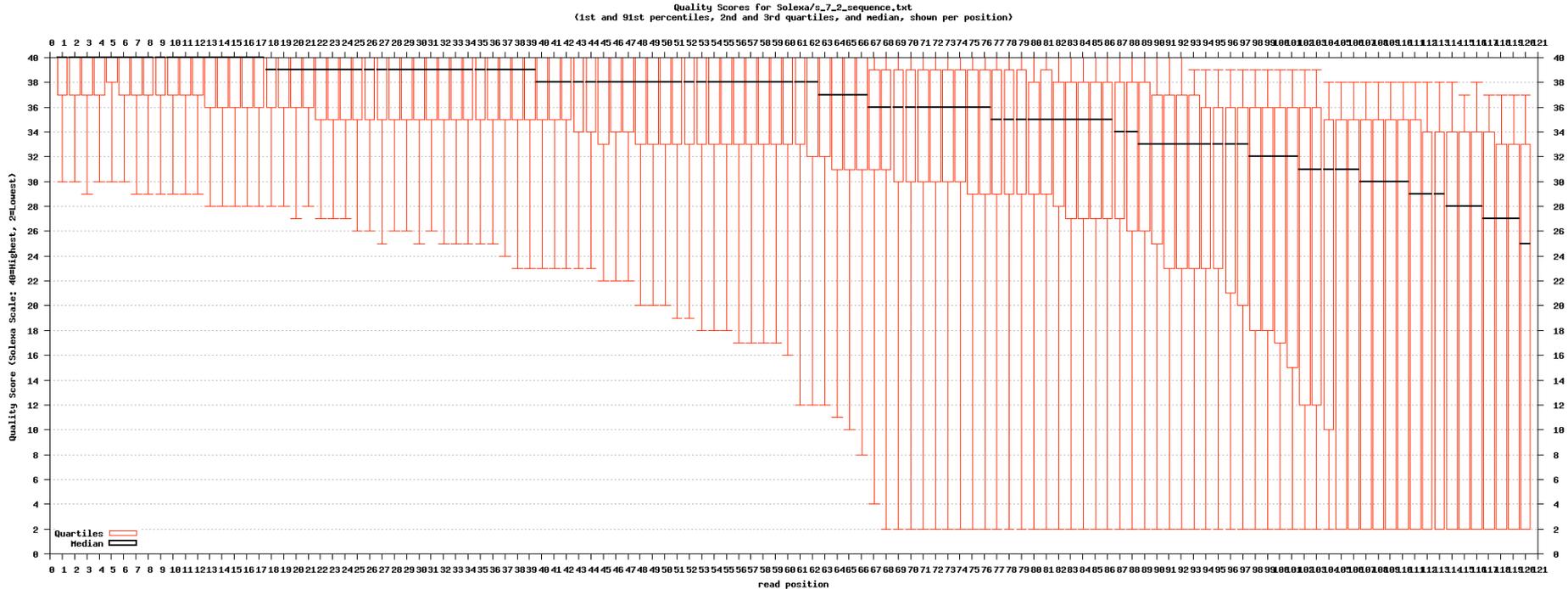
**Number of Reads of the samples**

<b>Samples</b>	<b>Total number of reads</b>	<b>Utilized for analysis</b>
Normal B cell	12,887,282	YES
RL1	60,272,006	NO
RL2	61,043,078	NO
RL3	65,579,872	NO
RL4	125,256,746	YES
Call4_1	62,741,712	NO
Call4_2	62,607,906	NO
Call4_3	133,542,778	YES
ALL B-Cell	77,888,742	YES

**Read coverage**

	Read coverage of gene region	Read coverage of non-gene region	Read length
Call4 cell line	2.81129121903121	2.35780895	100
RL cell line	1.47413416423764	1.14648699	100
Normal B-cell	0.186290512630489	0.1725446	76
ALL B-cell	1.788587532589	1.49037874	120

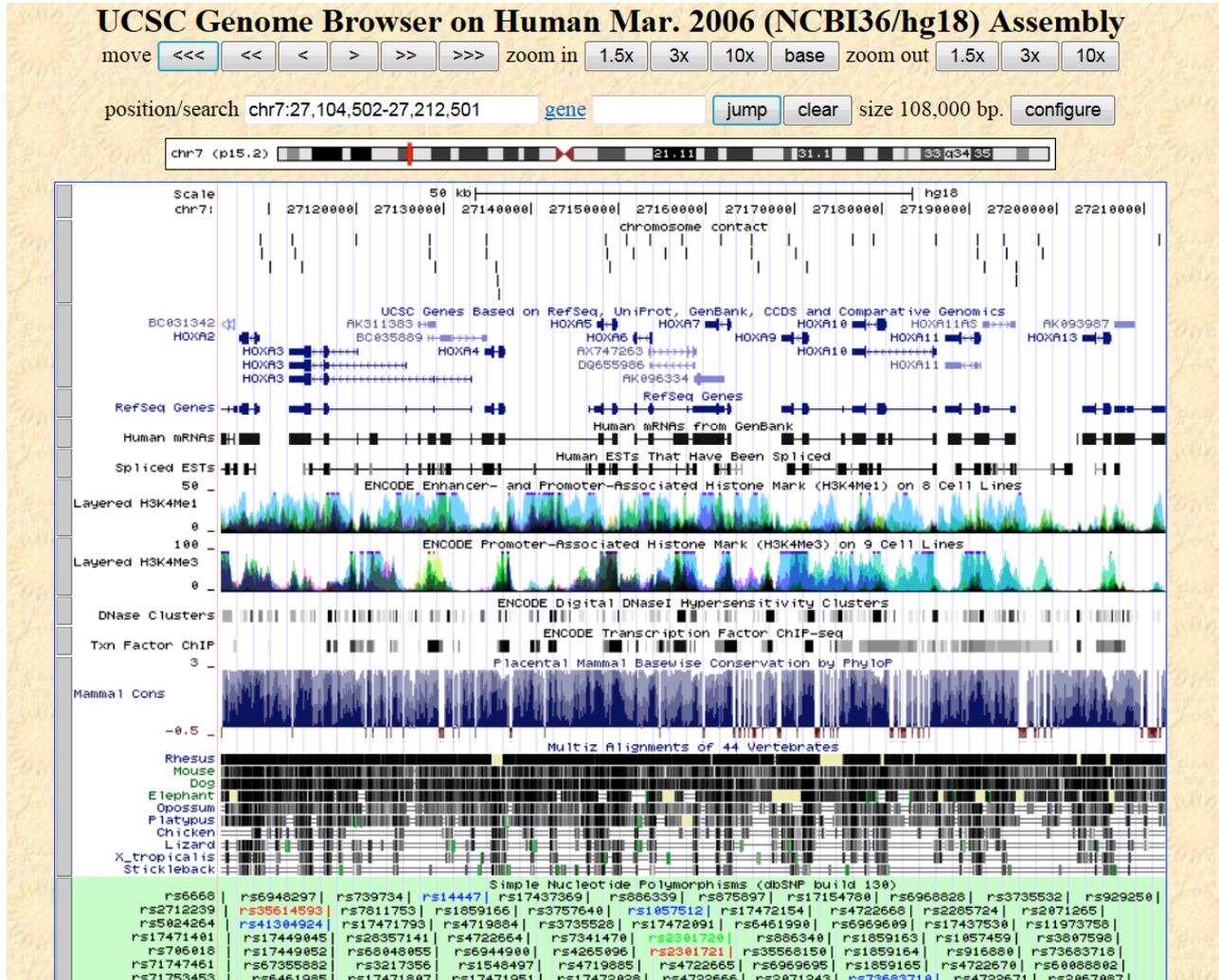
# Read Quality



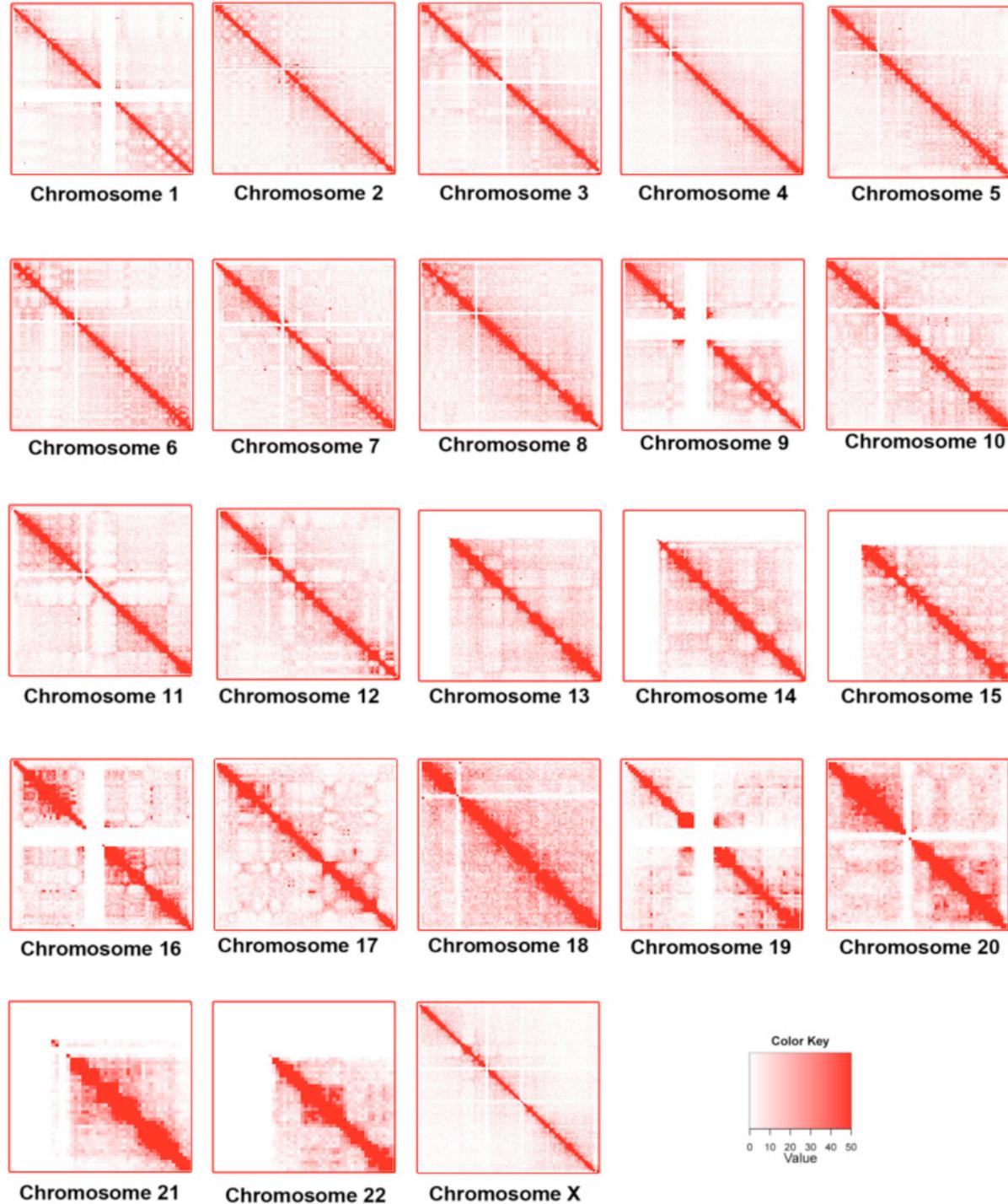
The sequencing quality score at a position is calculated as  $Q_{solexa} = -10 \log_{10} \frac{p}{1-p}$ , where  $p$  is

the probability of a sequencing error at the position. A score 30 means the probability of a sequencing error at the position is  $\sim 0.001$ . A score 20 or above may be considered acceptable.

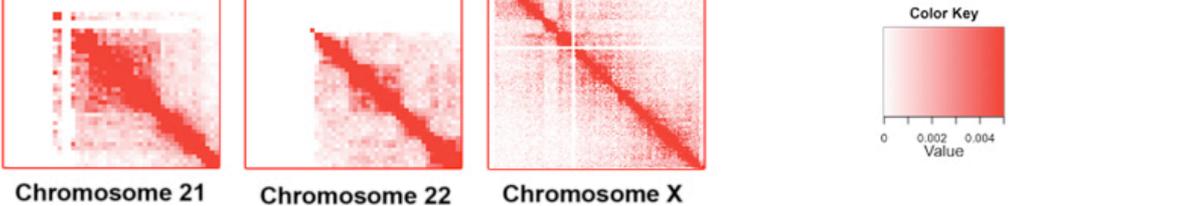
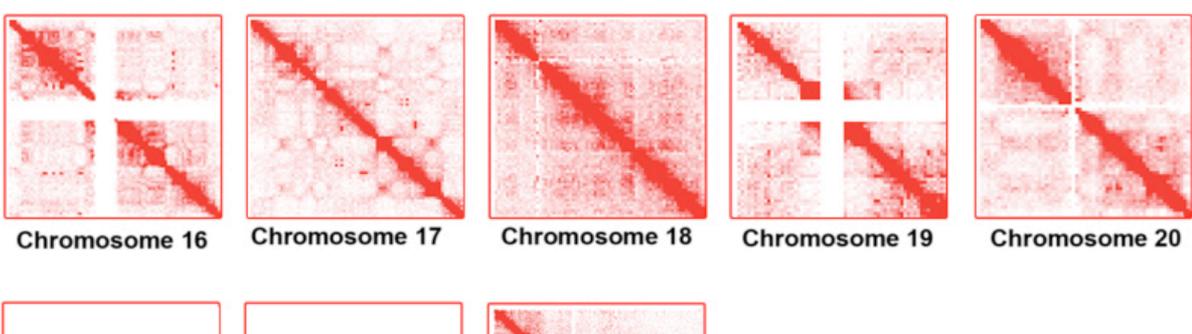
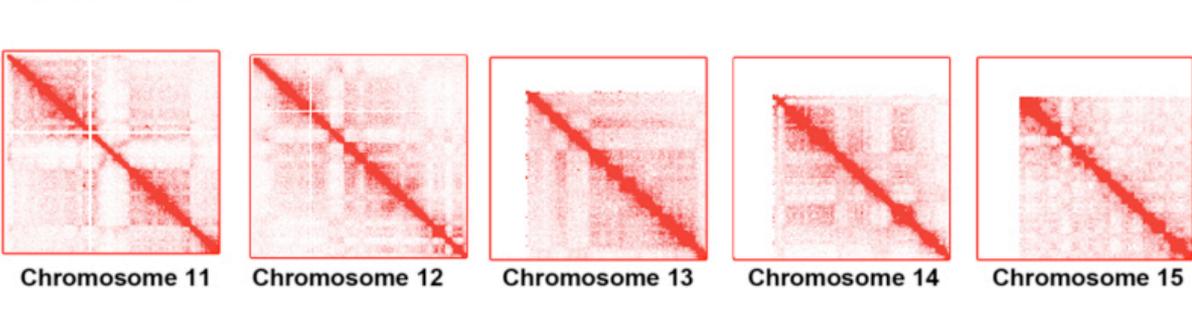
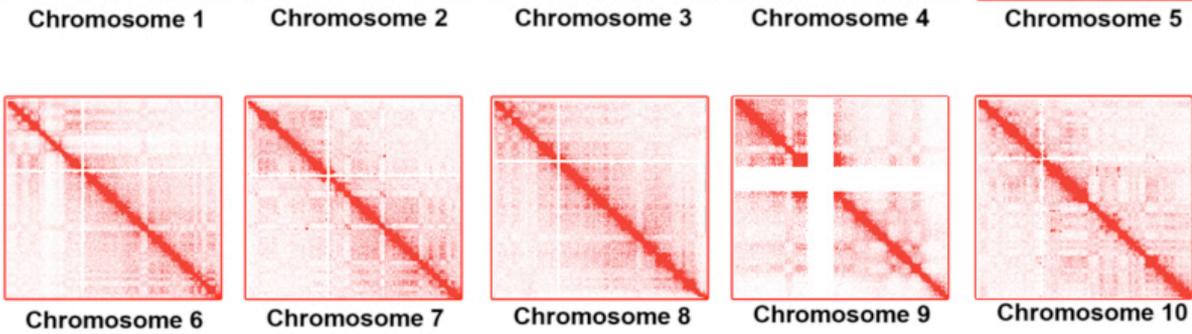
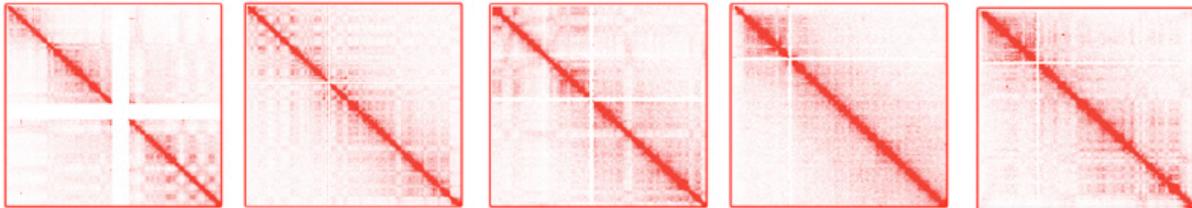
# Mapping Reads to Human Genome



# Un-normalized intra- chromosomal heat maps for primary ALL B- Cell



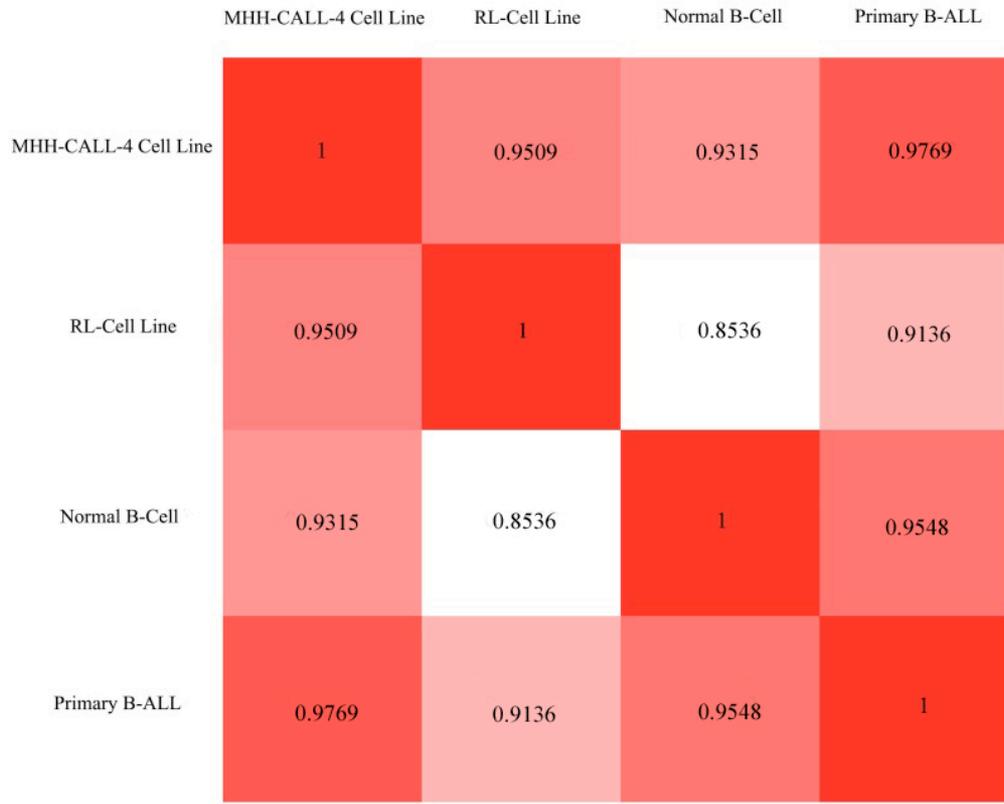
# chromosomal heat maps for primary ALL B-Cell



Sequential  
Component  
Normalization:  
 $M[i,j] / |M[i]|$

$M[i,j] / |M[j]|$   
Repeat until  
symmetric

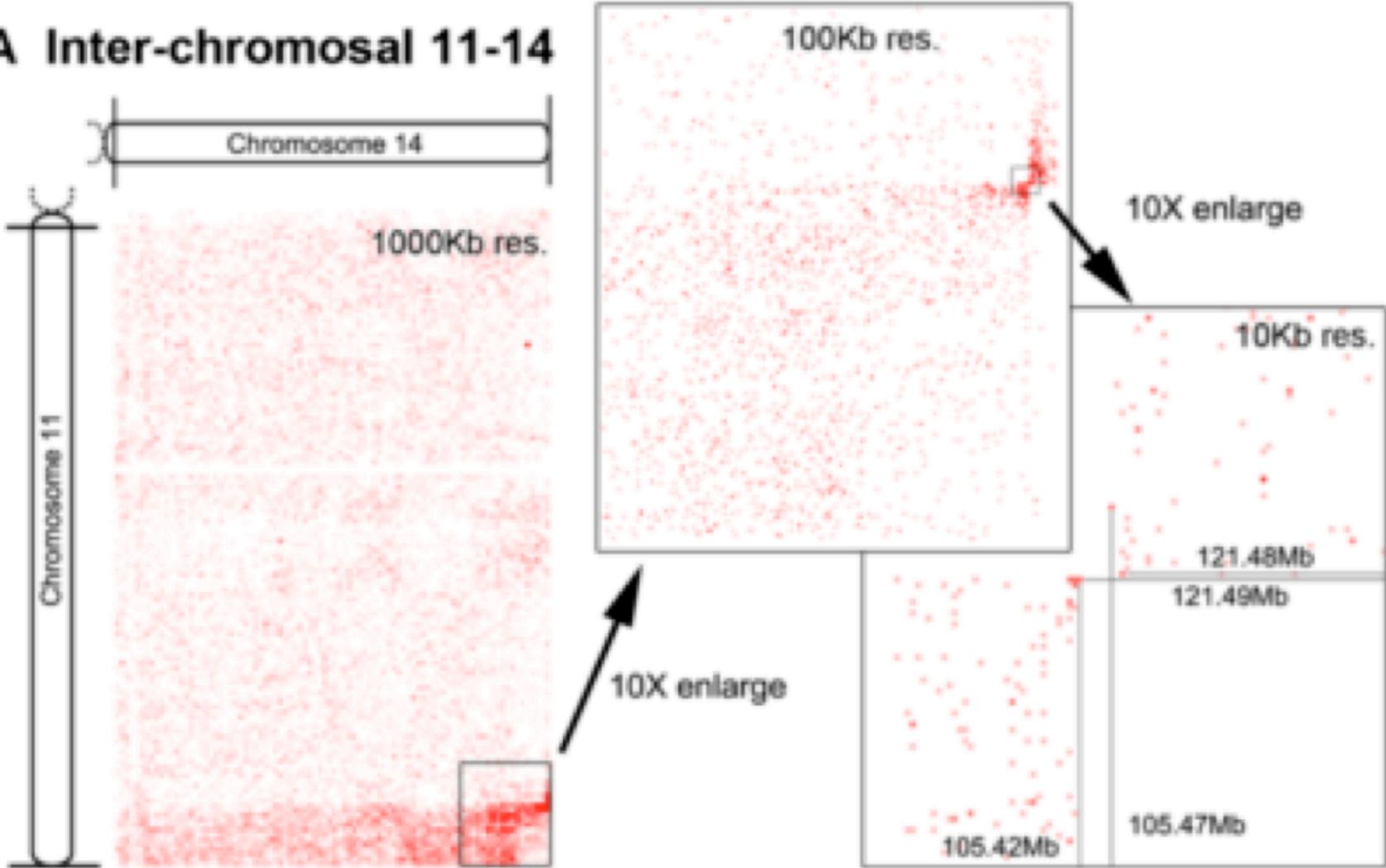
# Contact Correlation Between Cell Types



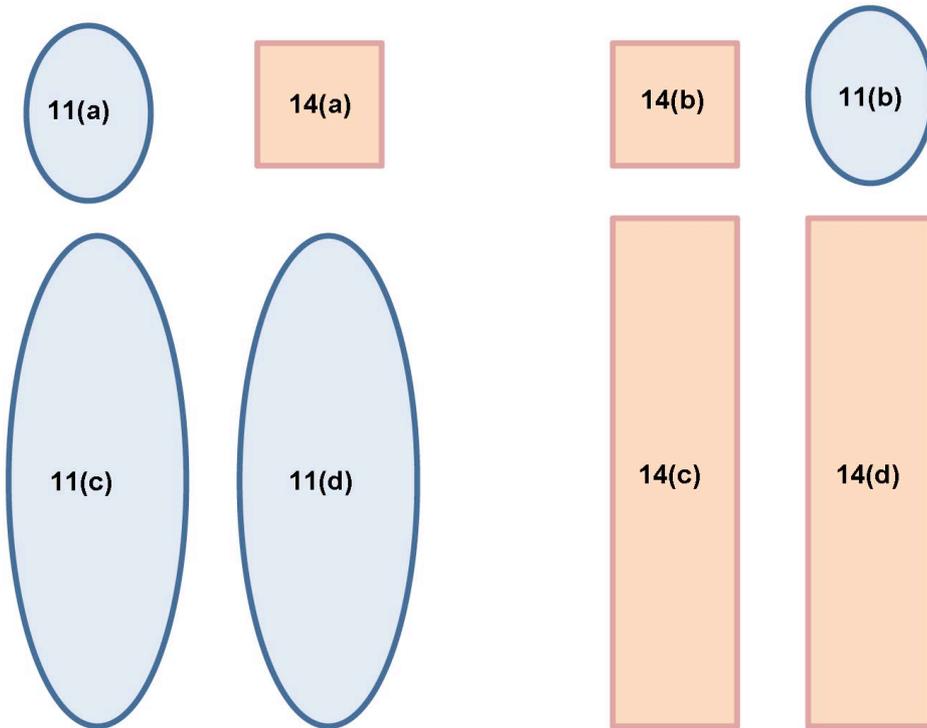
**Correlate intra-contact numbers of 23 pairs of chromosomes**

# Inter-Chromosome Contacts and Chromosome Translocation

**A Inter-chromosomal 11-14**



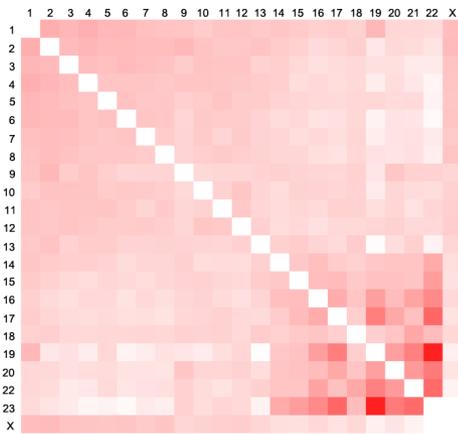
# Cancer Causing Chromosome Translocation



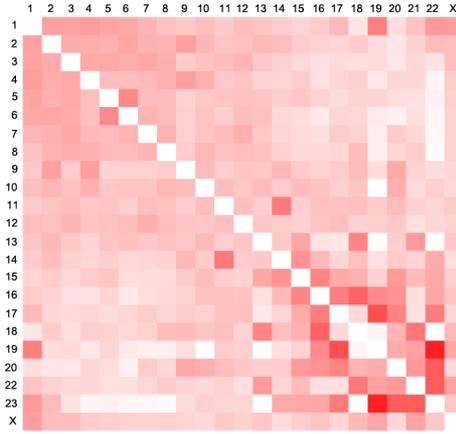
- Reconstruction of translocated chromosomes
- Two cancer related genes

# Comparison of Inter-Chromosome Contact Profiles of Different Cells

**A Normal B-Cell Line**



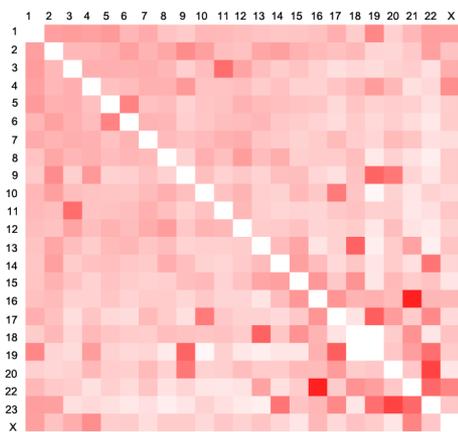
**B Primary B-ALL**



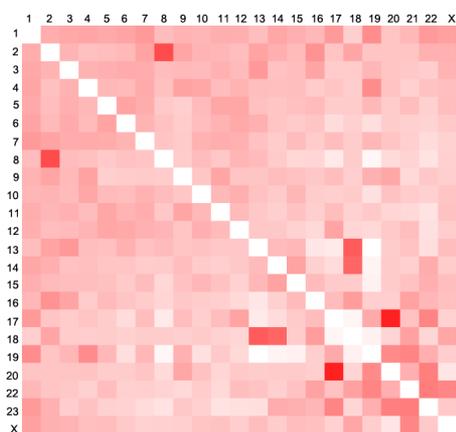
the expected number of contacts between chromosome  $i$  and  $j$  was calculated by

$$E_{i,j} = R_i \times R_j \times N_{INTER}$$

**C MHH-CALL-4 Cell Line**



**D RL-Cell Line**

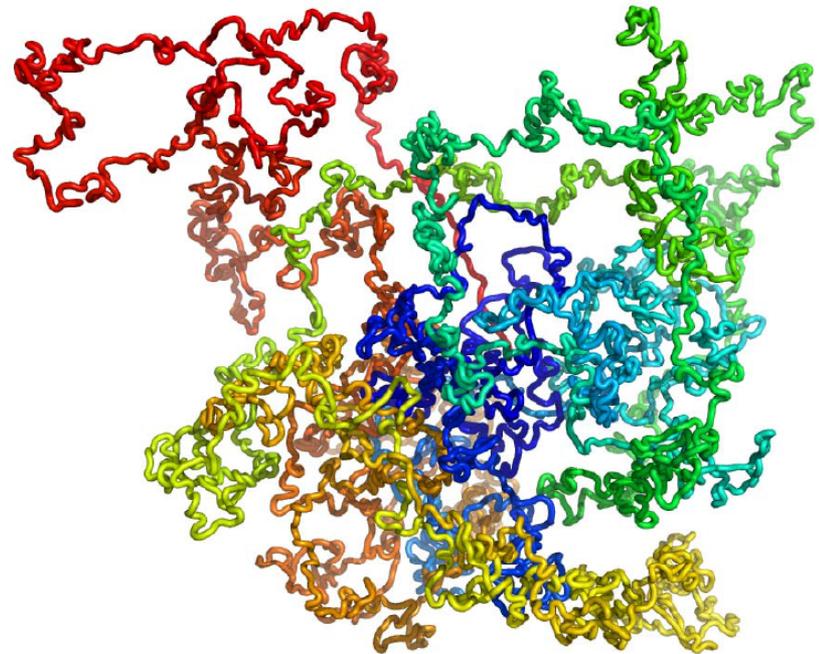
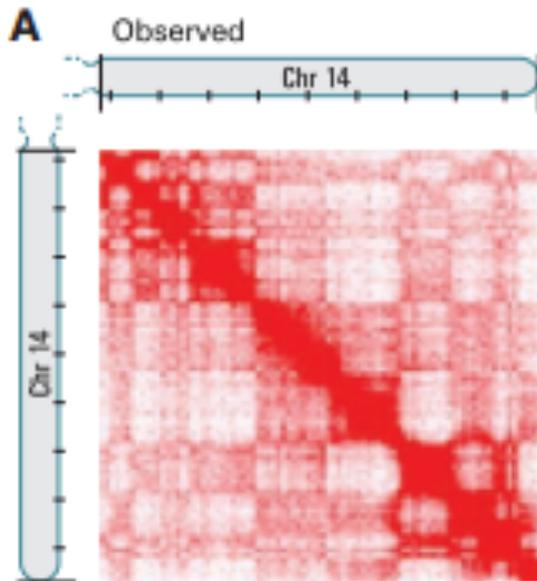




# Hi-C Databases

- 4DN data portal: <https://data.4dnucleome.org>
- ENCODE database:  
<https://www.encodeproject.org>

# 3D Genome Structure Modeling



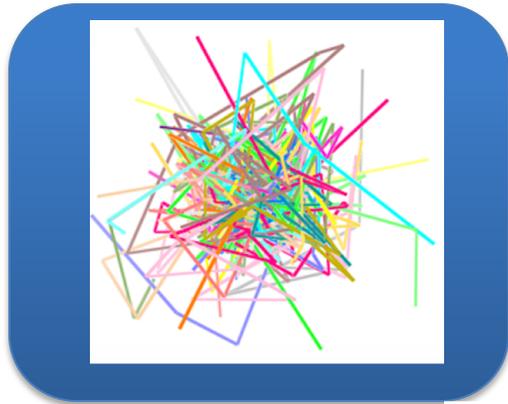
# Spatial Representation of Chromosome or Genome



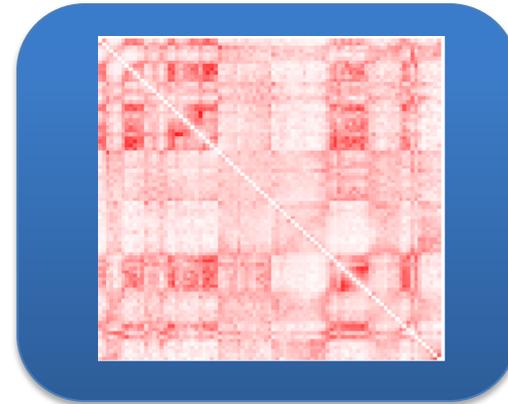
- Divided into  $N$  equal-size (e.g. 1MB) consecutive units sequentially
- The center of a unit is denoted by a point and its coordinate  $(x, y, z)$

# Contact Driven Structure Modeling

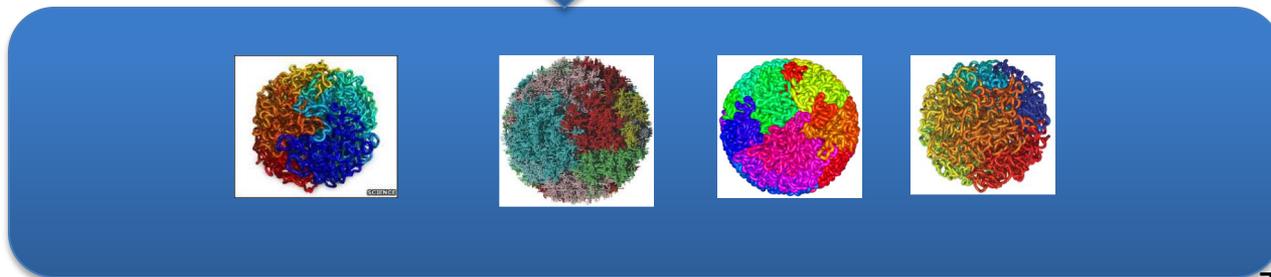
Random Structure



Normalized Contact Map



Permutation & Optimization



# Opportunities

- Chromosome contact data can be generated easily and cheaply
- Chromosome contact data is rather reliable
- A 3D model of a genome is very valuable in studying spatial regulation of gene expression and methylation

# Challenges

- Genome and chromosome is very large (3 billion nucleotide of human genome)
- Genome structure is very dynamic
- No known experimental genome structure other than some point distance data generated by FISH
- Relationship between contact and distance is not deterministic
- Hi-C data is noisy

# Maximum Likelihood or Distance-Based Approaches

- **Convert interaction frequency (contact number) into distance**

$$d \propto 1 / IF^a$$

**or**

$$IF \propto 1 / d^a$$

- **Translate distance into (x, y, z) coordinates**

# A MCMC Approach

M. Rousseau, J. Fraser, M.A. Ferraiuolo, J. Dostie, M. Blanchette. ***Three-dimensional modeling of chromatin structure from interaction frequency data using Markov chain Monte Carlo sampling***. BMC Bioinformatics, 2011.

# MCMC5C

- Formulate a probabilistic model linking 5C/Hi-C data to physical distances
- **Markov chain Monte Carlo (MCMC)** approach called MCMC5C to generate a representative sample from the posterior distribution over structures from IF data.

# MCMC5C

- Structural properties (base looping, condensation, and local density) were defined in the models
- Applied these methods to a biological model of human myelomonocyte cellular differentiation and identified distinct chromatin conformation signatures corresponding to each of the cellular states.
- run on Hi-C data and produce a model of human chromosome 14 at 1Mb resolution that is consistent with previously observed structural properties as measured by 3D-FISH.

# Other Existing Methods

- **5C3D (Fraser et al.):** translates IF values into physical distance estimates and then uses a gradient descent approach to find the 3D conformations.

# Other Existing Methods

- **Bau et al.** Interactions are modeled with springs whose equilibrium length depends on the observed IF values, subject to certain constraints based on the structure of the 30-nm fiber, optimized by Integrative Modeling Platform.

# Other Existing Methods

- **Duan et al.**, convert interaction frequencies to Euclidean distances and then seek conformation minimizing the misfit, with addition of a set of clash avoidance constraints and a few prior known knowledge about the yeast genome organization. The constrained optimization problem is solved to find the best structure.

# Possible Drawbacks of Existing Methods

- Objective function (sum of square difference between predicted and derived distance) is debatable.
- Assume each IF is equally reliable.
- The absence of an underlying probabilistic model, preventing the calculation of confidence intervals on specific structural properties (e.g. distance between two genomic sites)

# Probabilistic Model of Chromatin Conformation

- A chromosome is modeled as a continuous piece-wise linear curve in 3D.
- Theoretical interaction frequency between fragment  $i$  and  $j$ , denoted  $IF(i, j)$ , is inversely correlated with the distance between two fragments in 3D conformation:  **$IF(i, j) = f(D_s(i, j))$** , where  $D_s(i, j)$  is the Euclidean distance between sites  $i$  and  $j$  and  $f$  is an appropriately chosen function.

$$f(D_s(i, j)) \propto 1/D_s(i, j)^\alpha$$

# Probabilistic Model of Observed IF and Theoretical IF

$$\hat{IF}(i, j) | IF(i, j), \sigma(i, j) ] = N(\hat{IF}(i, j); IF(i, j), \sigma(i, j)^2)$$

re  $N(x; \mu, \sigma^2)$  is the normal density function

# Observed IF and Theoretical IF (Hi-C)

$$\Pr[\hat{r}(i, j) | r(i, j)] = N(\hat{r}(i, j); r(i, j), r(i, j) + \kappa). \quad (2)$$

The role of  $\kappa$ , which we set to 10, is to avoid having small read counts being assigned too low a variance.

# Posterior Probability (Structure | IF data)

The observed data  $\hat{IF}$  defines a posterior distribution over the set of possible conformations of the chromatin:  $\Pr[\mathbf{S}|\hat{IF}] = \Pr[\hat{IF}|\mathbf{S}] \cdot \Pr[\mathbf{S}] / \Pr[\hat{IF}]$ . Since there are no constraints imposed on the structure space and the probability of the observed data ( $\hat{IF}$ ) is constant with respect to  $\mathbf{S}$ , we get  $\Pr[\mathbf{S}|\hat{IF}] = \zeta \cdot \Pr[\hat{IF}|\mathbf{S}]$ , for some constant  $\zeta$ , and thus

$$\Pr[\mathbf{S}|\hat{IF}] = \zeta \cdot \prod_{i,j} \Pr[\hat{IF}(i,j) | IF(i,j) = f(D_{\mathbf{S}}(i,j), \sigma(i,j))].$$

# Sampling Conformations from Posterior Distribution to maximize probability

- A random structure  $R_0$  is initially chosen to seed the process ( $t=0$ ), where each point is placed randomly in a cube of side length  $10 * \text{avg}(f(\text{IF}))$ .
- Repeat: The current structure  $R_t$  is randomly perturbed to obtain a new structure  $R_{t'}$ . If  $\text{Pr}[R_{t'} | \text{IF}] > \text{Pr}[R_t | \text{IF}]$ , the perturbation is obtained and we set  $R_{t+1} = R_{t'}$ . Otherwise, we set  $R_{t+1} = R_t$ .
- For values of  $t$  sufficiently large, find a structure with high probability:  $\text{Pr}[S | \text{IF}]$ , let  $R_t = S$ .

# Random Structure Perturbation

- Randomly choose one point  $S(i)$  along the structure and moving it by a vector  $v$  randomly choosing within a sphere of radius  $r$  (e.g.  $r = 0.25$  nm)
- The likelihood of the resulting structure is then quickly obtained from that of the old by updating the terms corresponding to the pairs of points involving  $i$ .

# Assessing Mixing

- $R_1, \dots, R_k$  of early iterations are highly dependent on  $R_0$ .
- Determine at what point  $m$ , the Markov process has mixed, i.e.,  $R_m$  is independent of  $R_0$
- After mixing, i.e. for  $k \geq m$ , any sample  $R_k$  is representative of the target distribution. For  $d$  sufficiently large,  $R_k$  and  $R_{k+d}$  are independent.

# Convergence Determination

- Run two independent chains  $R$  and  $R'$  in parallel, from independently chosen initial conformations  $R_0$  and  $R_0'$ .
- Mixing is achieved if the samples  $\{R_{k/2}, \dots, R_k\}$  and  $\{R'_{k/2}, \dots, R'_k\}$  cannot be distinguished from each other. Specifically, the average pairwise structural distances among  $R_k$  is compared to that between  $R_k$  and  $R_k'$ .
- After mixing is achieved, collect samples every  $d=k/20$  iterations.

# Clustering of Structure Ensembles

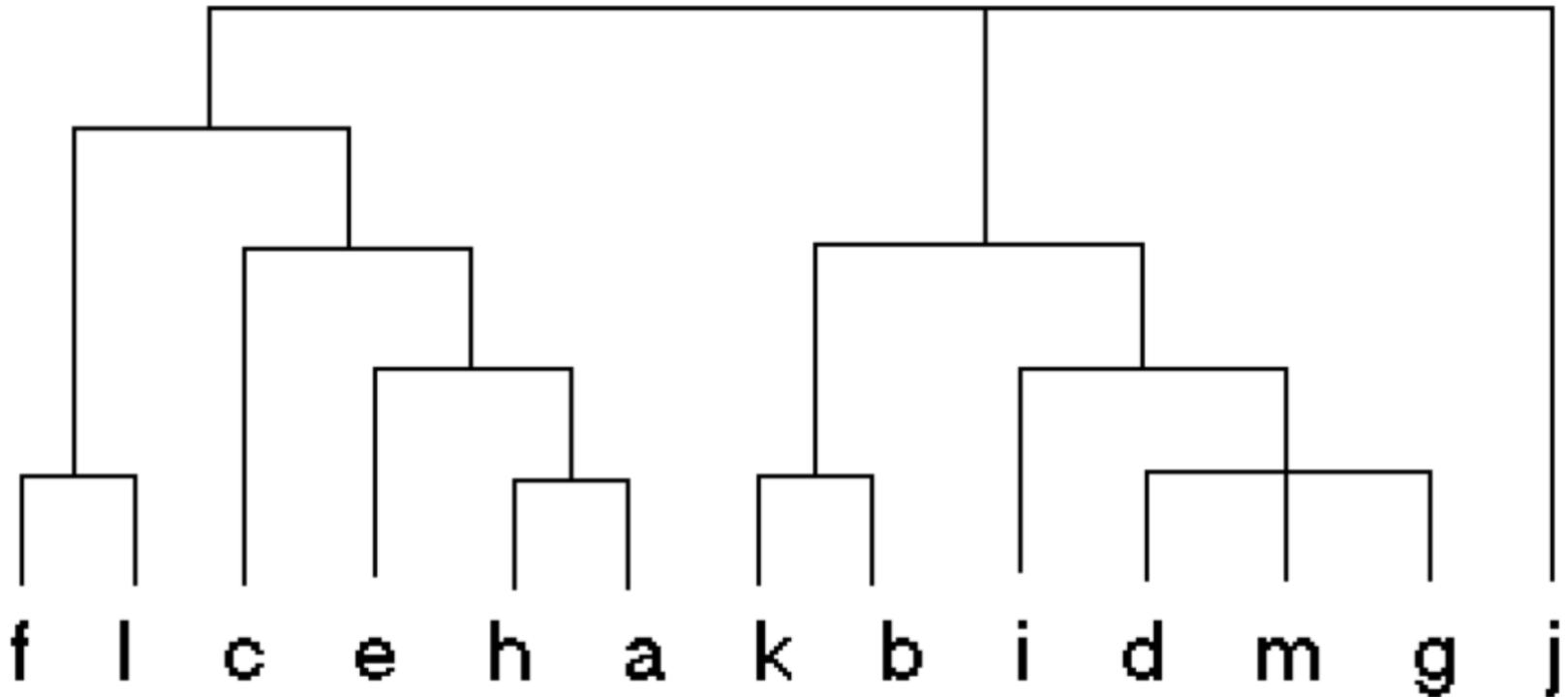
- Distance metric:  $N \times N$  intra-structure distance matrix  $D_S$ .
- The distance  $(S, T)$  between two structures  $S$  and  $T$  is:

$$\text{dist}(S, T) = \sqrt{\sum_{i,j} (D_S(i, j) - D_T(i, j))^2}$$

# Structure Clustering

- Hierarchical clustering
- Visualization: tree dendrogram
- Visual inspection is performed to determine the tree height cutoff and number of subfamilies
- Choose maximum likelihood structure from each cluster as representative and assigning it a weight proportional to the number of structures in its cluster.

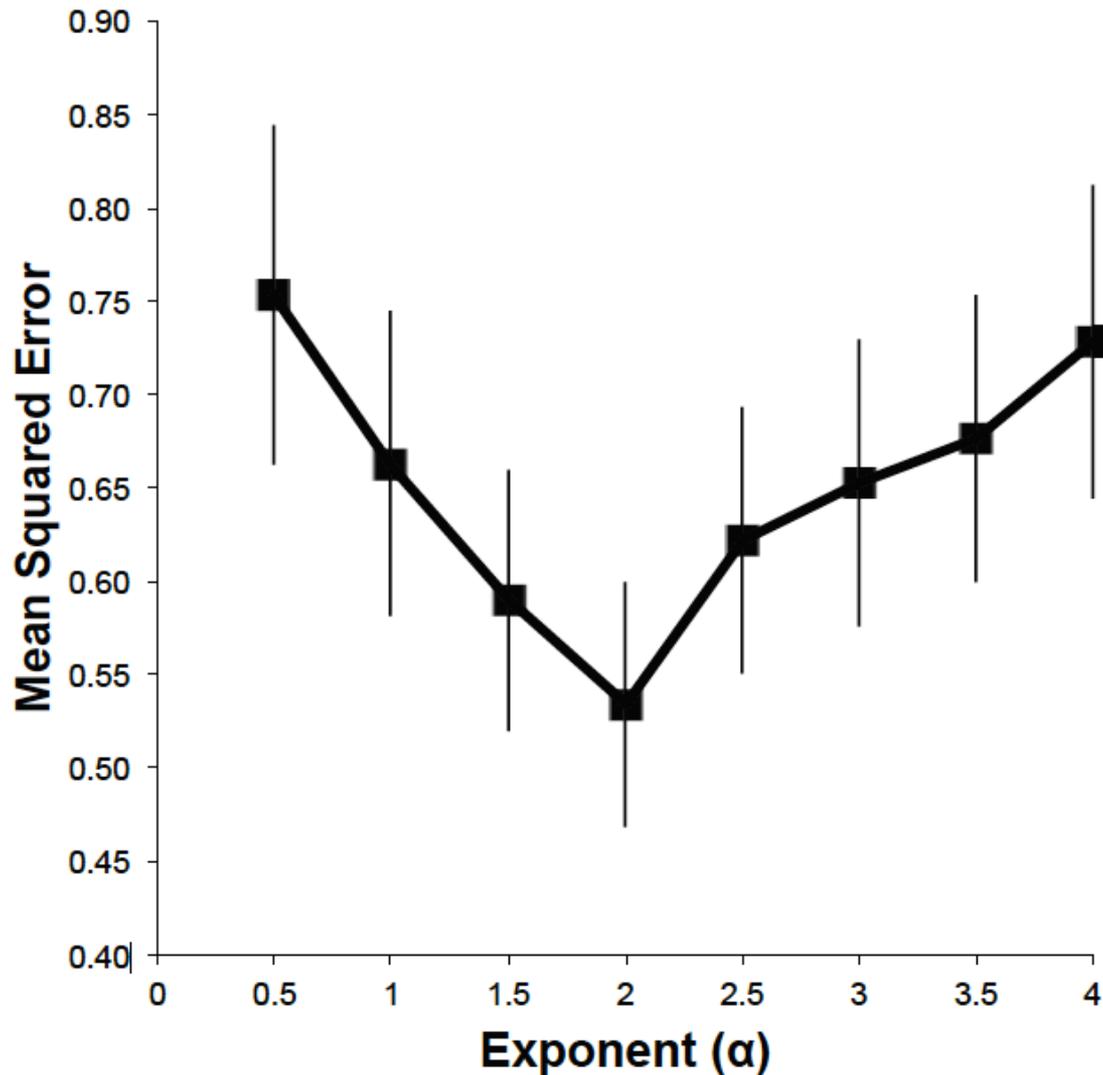
# Hierarchical Clustering



# Convert IF to distance

- The most accurate model is the one that is best able to predict unseen pairwise interaction frequencies. For each of a set of possible a leave-one-out cross-validation was performed. Find  $\alpha$  to minimize

$$MSE(\alpha) = \frac{1}{n} \sum_{(i,j)} (D_{\mathcal{S}_{(i,j)}^*; \alpha} (i, j)^{-\alpha} - \widehat{IF}(i, j))^2.$$

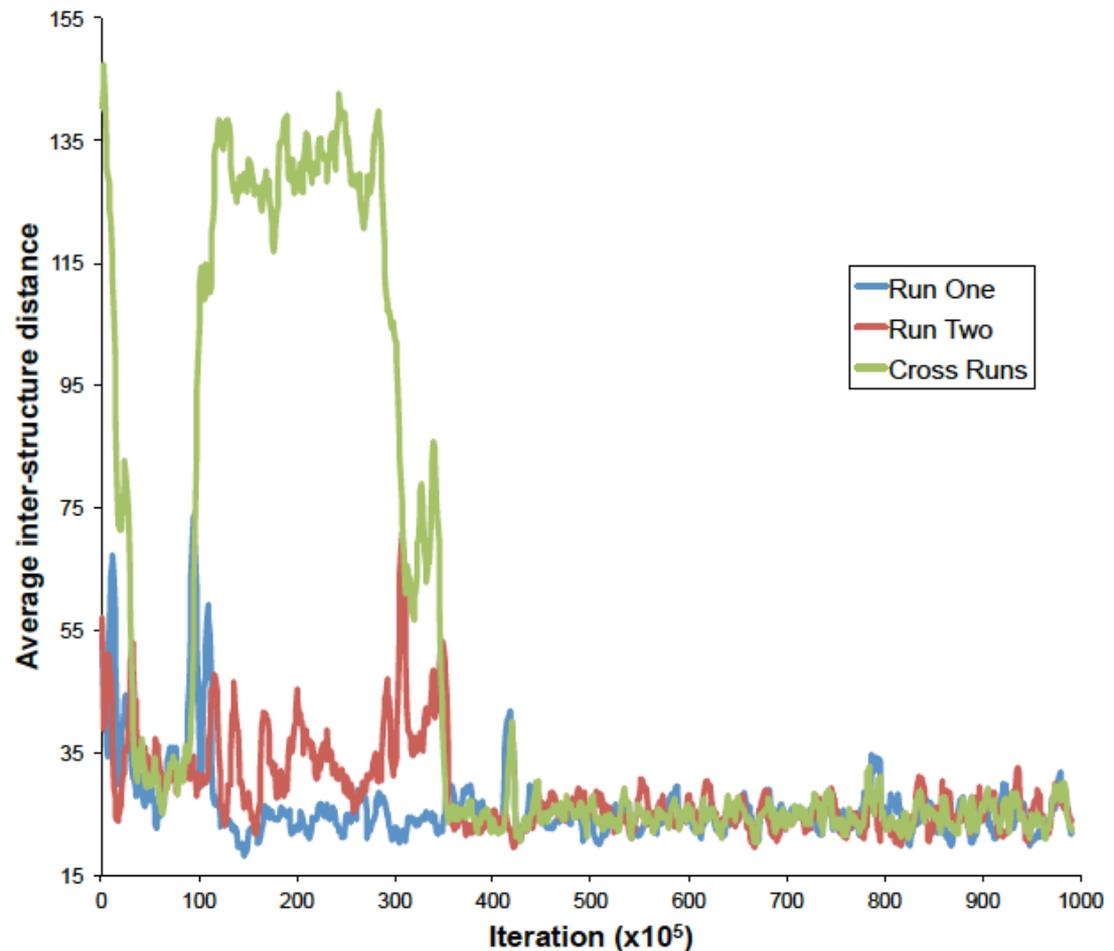


**Figure 3 Leave-one-out cross-validation.** Value of the mean-squared-errors as a function of  $\alpha$ , obtained for a leave-one-out cross-validation on the HB-1119 dataset. The minimum error is found for an exponent of 2.0, although values of  $\alpha$  between 1 and 3 do not produce significantly worse errors.

## C (Distance Scaling Factor) Calibration

Without physical measurement of the distance between pairs of points along the sequence, it is difficult to accurately estimate the value of  $C$ . However, based on the average IF value of pairs of fragments located less than 5kb apart along the sequence and following Bystricky et al . [51] that packed chromatin has a physical length of 1 nm for every 110-150bp,  $C$  was estimated as approximately 50 nm.

## Assessment of Mixing



**Figure 4** Mixing of parallel *MCMC5C* runs (HB-1119 dataset).

Distance between consecutive structures (sampled every  $10^6$  iterations) from within one of two parallel *MCMC5C* runs (blue and red curves) or across the two runs (green curve), on the HB-1119 5C dataset. The runs converge to the same distribution very rapidly (in less than 250 seconds) and the cross-run distance (green) drops to within the same range as the within-run distances (blue and red curves) after  $350 \times 10^5$  iterations.

# Experiment

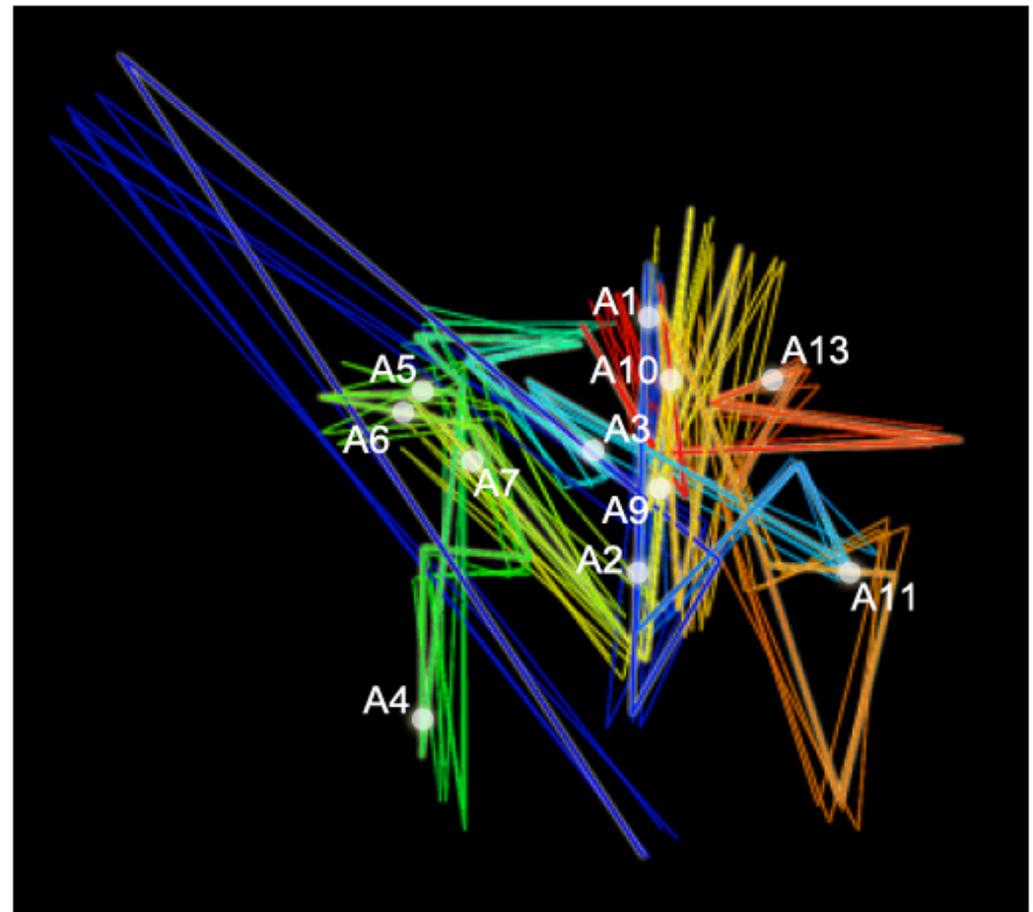
- Figure 4 shows that mixing is achieved after approximately  $350 * 10^5$  iterations, which requires less than 250 seconds of running time. Passed this point, structures sampled every  $10^6$  steps from the two parallel runs are undistinguishable from each other and sample structures from the same distribution.
- 250 structures were sampled after burn-in from each of the two runs. The two ensembles of structures were then combined and the 500 structures were clustered based on their structural similarity
- Analysis of the two THP-1 5C datasets produced similar results, and runs of a larger number of parallel MCMC chains confirm that they all sample similar structures.

# Simulation Verification

**Gold structure:** a computationally constructed 3D structure used to generate IF data.

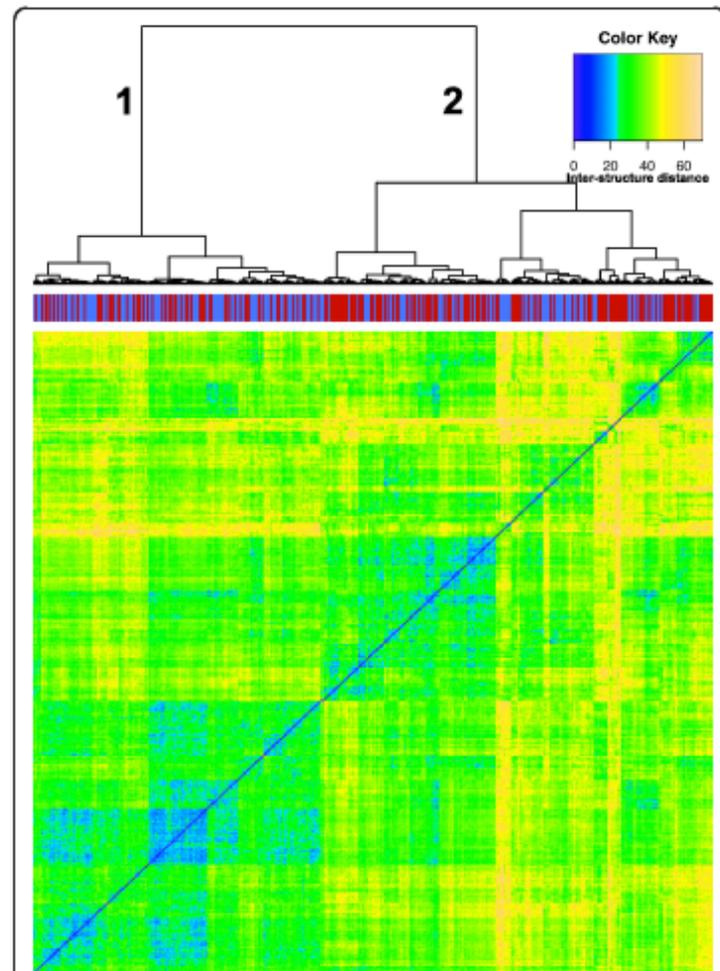
Simulated structure: models constructed from the IF data of the gold structure.

# Verification by Sampling from Simulated True Structures



**Figure 6** HB-1119 Structures from simulated data aligned to gold standard structure. The “gold standard” structure is used as a reference structure to which structures from four different parallel *MCMC5C* runs on simulated data generated from the gold standard structure are aligned. The gold standard structure is shown highlighted with a white glow and the transcription start sites for the HoxA genes are annotated. The structures found from the simulated data are shown in superimposition to the gold standard structure and show a high degree of alignment.

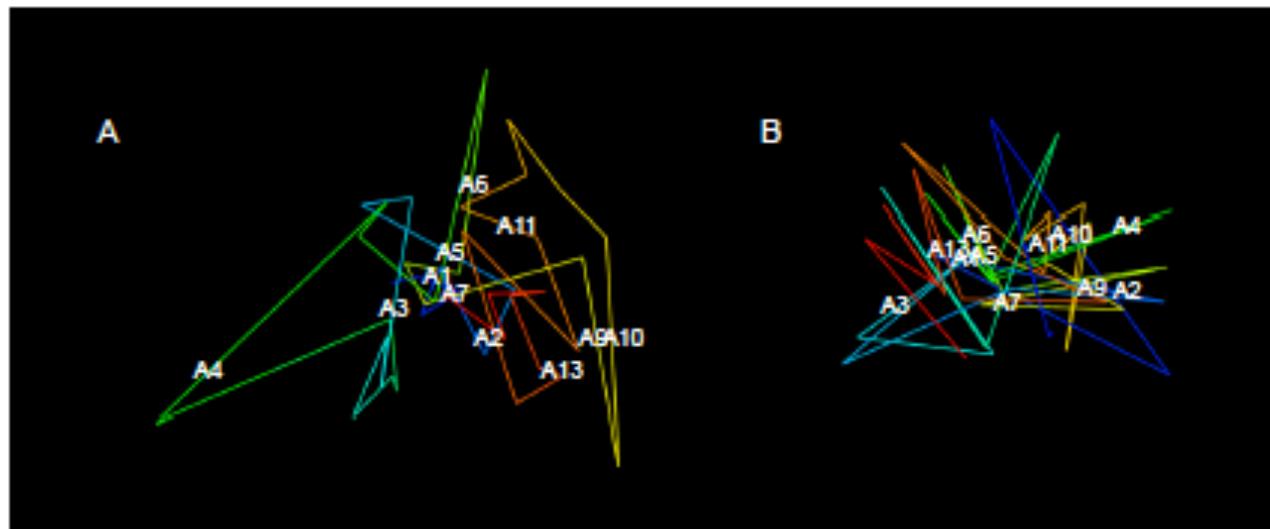
Structure  
Clustering and  
Sub-Structure  
Families.  
Sub-structure  
families may  
correspond to  
chromatin  
structures of cells  
in different  
stages



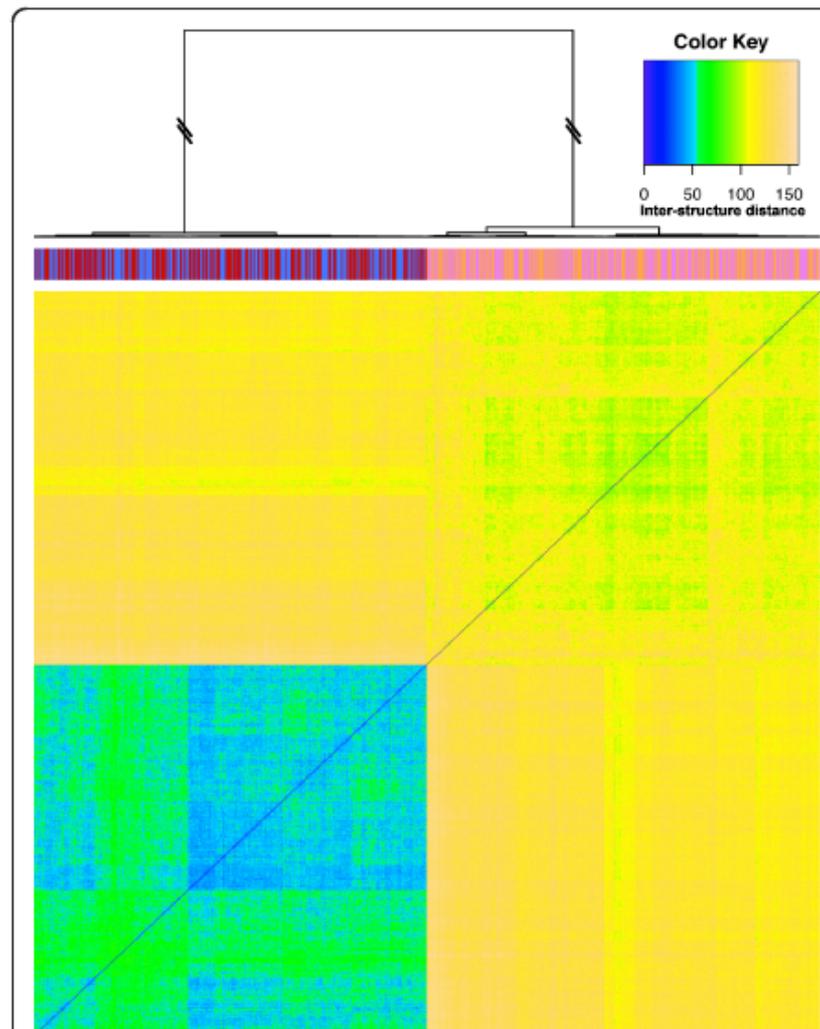
**Figure 5** Mixing and subclustering of HB-1119 structures.

Mixing and hierarchical clustering (Ward's method) of structure similarity. The five-hundred structures come from two parallel MCMCSC runs on the HB-1119 dataset (pools of 250 structures from each run were used). The colors along the top indicate which run each structure originated from (run one = blue, run two = red) and demonstrates that the sampling process has successfully mixed. The blocks in the heatmap and the dendrogram indicate the presence of sub-clusters of structures (numbered in the dendrogram). The two clusters (numbered 1 and 2) both contain structures from the two parallel runs (blue and red vertical bars), indicating that the structures are conserved across runs and are not an artifact of the burn-in process.

# Conformations of HoxA in Undifferentiated and Differentiated Conditions



**Figure 7 Models of HoxA cluster before and after differentiation.** Maximum likelihood structures found by *MCMC5C* from the undifferentiated and differentiated THP-1 datasets (A and B, respectively). The HoxA gene transcription start sites are annotated on each of the structures.



**Figure 8 THP-1 clustering of undifferentiated and differentiated structures.** Hierarchical clustering (Ward's method) of one-thousand structures from four parallel *MCMCSC* runs, two on the undifferentiated THP-1 dataset and two on the differentiated THP-1 dataset (250 structures each). The colors along the top indicate which state each structure originated from (undifferentiated run one = blue, run two = red; differentiated run one = pink, run two = orange) and demonstrate a clear distinction between the two states, indicating that the undifferentiated and differentiated cell states specify different structure signatures.

# Structural Variation and Conservation

The subset of fragments that are the most conserved across the ensemble of structures are found to lie within the central core region of the structures. These fragments are spatially close to each other and may be involved in looping contacts that are important for the maintenance of the chromatin structure and are therefore highly conserved.

## On Hi-C Data (Data Set II)

Model the long arm of human chromosome 14 (88.4 Mb region) from Hi-C data published by Lieberman-Aiden et al . [18] at a 1Mb resolution (89 fragments in total). Lieberman- Aiden et al . [18] proposed the existence of two physically disjoint compartments, whereby compartment A was found to correlate with open and actively transcribed chromatin, while compartment B was found to be more densely packed and repressed.

# MCMC5C Availability

- MCMC5C modeling movie:

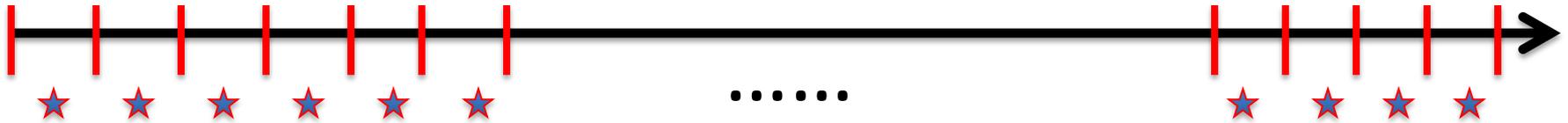
<https://www.youtube.com/watch?v=2aIQvU>

DIGNc

# Distance-based approach: convert IF to distance

- Duan et al., the resulting conversion approximately follows  $d \propto 1/IF$  .
- Mateos-Langerak et al . [50] also suggest a relationship of the form  $d \propto 1 / IF^a$  .
- Bau et al . [28] convert their IF via a linear transformation of the IF' s z-score.

# Spatial Representation of A Genome Region at a Scale



- A genome region (e.g. a chromosome) is divided into  $N$  equal / variable size units sequentially
- The spatial position of the center of a unit is denoted by a point and its coordinate  $(x, y, z)$  and constraints on the size of unit (radius of sphere)
- The consecutive points are joined into fragments forming the folding trace of the region.

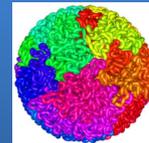
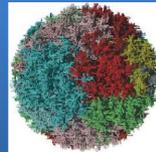
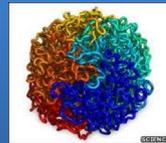
# Structure Construction at One Scale

Initial Structure Representation of Units

Contact Map Between Units

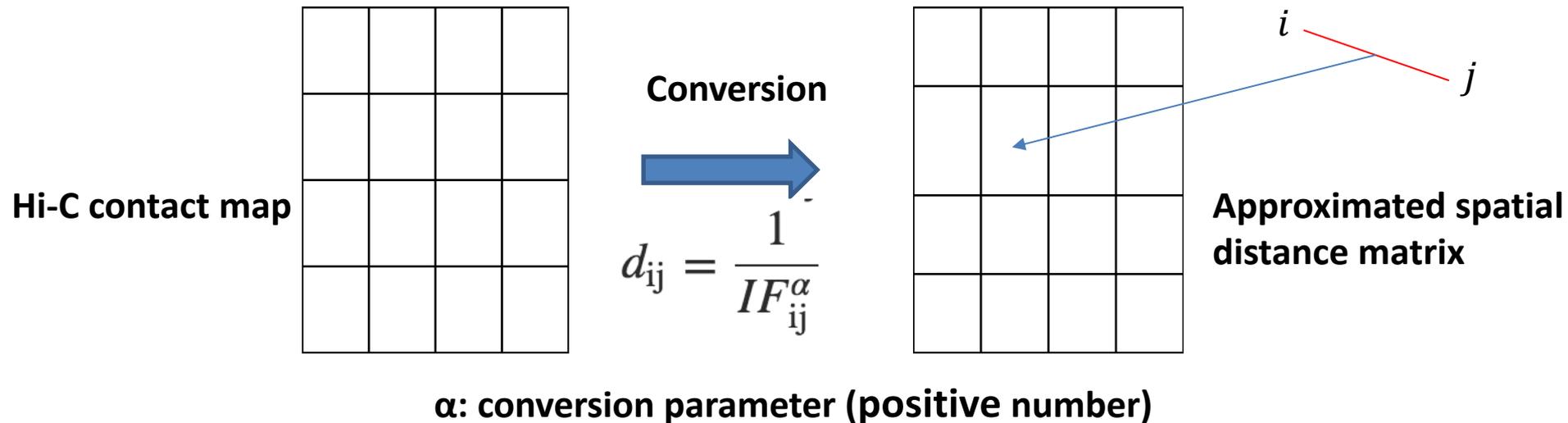


**Sampling & Optimization**



# A Data Driven Optimization Approach

- Convert observed interaction frequencies into distance map



- Objective:** Adjust positions ( $x, y, z$  coordinates) of beads in 3D space to satisfy converted expected distances as well as possible
- Challenges:** approximated distances, noise, conflicts

# Problem of Traditional Objective Functions

Square Error

$$\sum_{i < n; j < n} (x_{ij} - d_{ij})^2$$

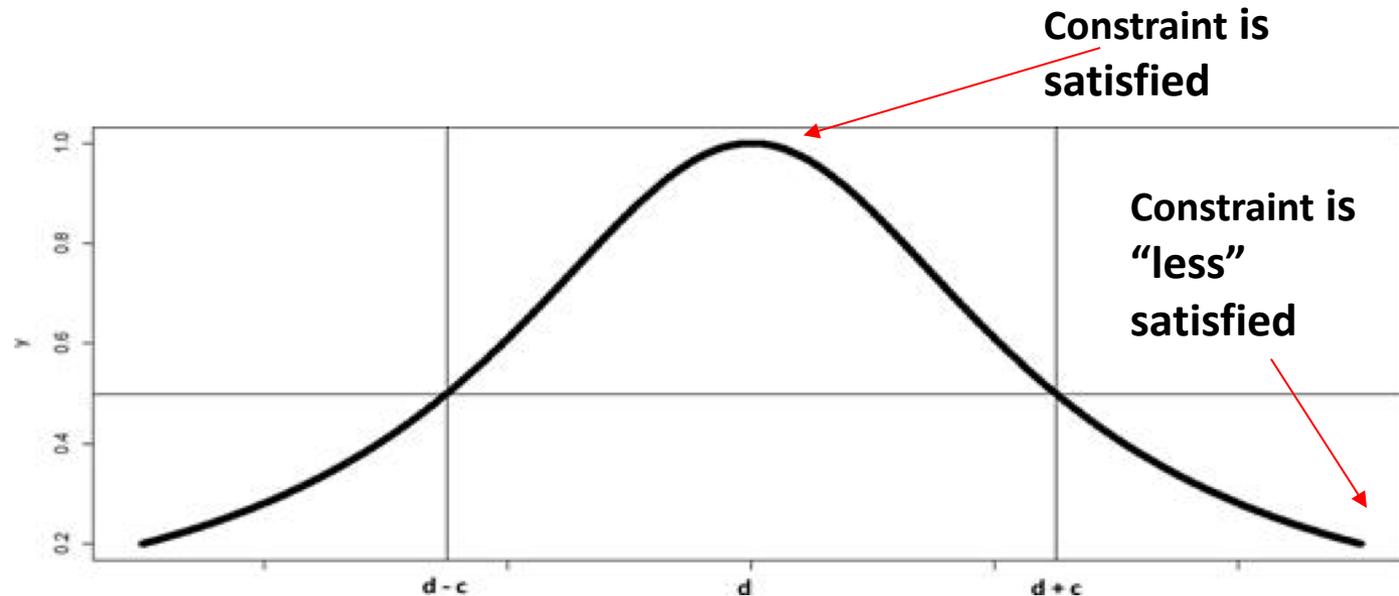
- Very sensitive to noisy / conflicting distances
- Cannot build 3D models of large genomes involving noisy inter-chromosomal contacts

# Deriving Soft Constraints by Lorentzian Function

- Use Lorentzian function to measure the satisfaction of a distance restraint

$$\frac{c * c}{c * c + (x - d)^2}$$

**[0, 1]**



- Tolerate noisy restraints and maximize the satisfaction of feasible restraints

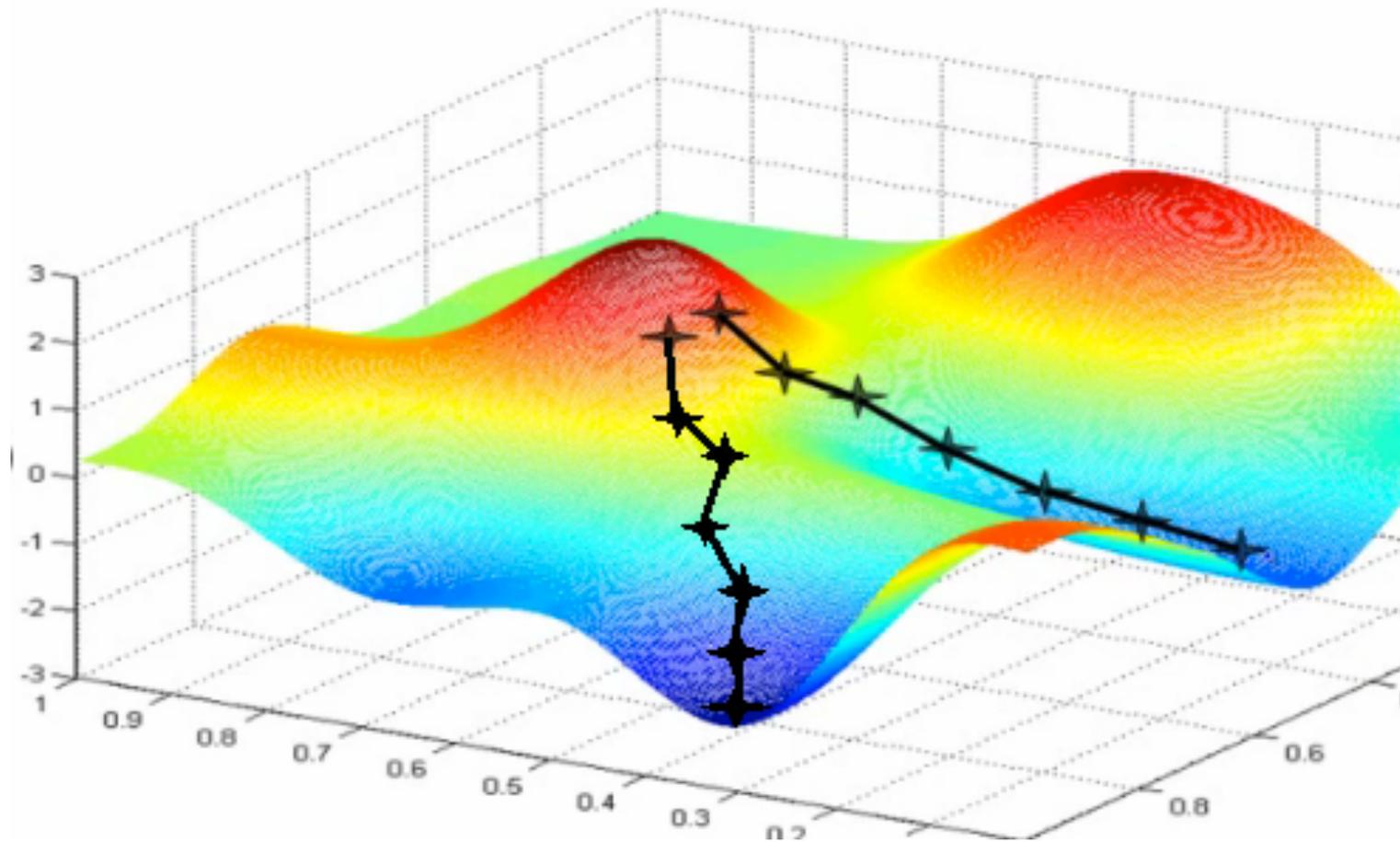
# Objective Function

- Objective function
  - Sum of all distance constraints
  - Adjacent beads have maximum weight to enforce their proximity

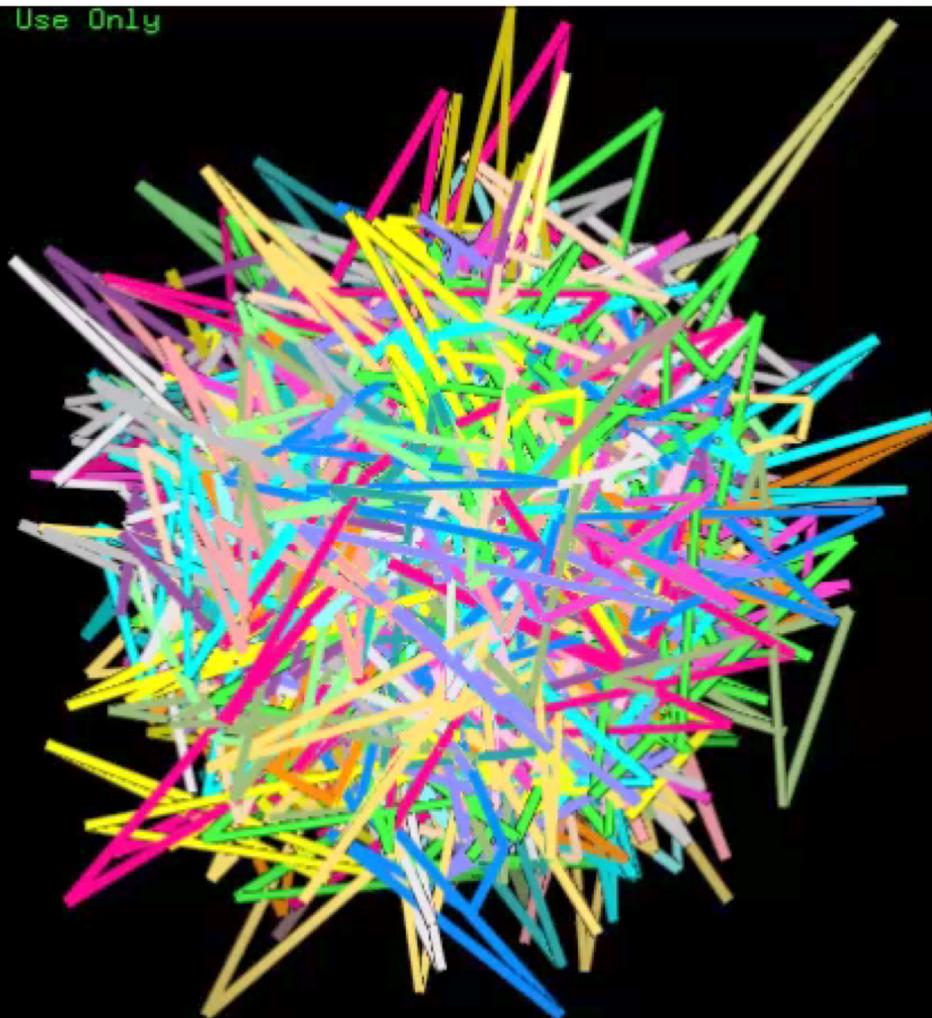
$$fn = \sum_{|i-j|=1} \frac{c * c * IF_{\max}}{c * c + (x_{ij} - d_{ij})^2} + \sum_{|i-j| \neq 1} \frac{c * c * IF_{ij}}{c * c + (x_{ij} - d_{ij})^2}$$

- $c$  is set to average of  $d_{ij}$  to control gradient vanishing during optimization

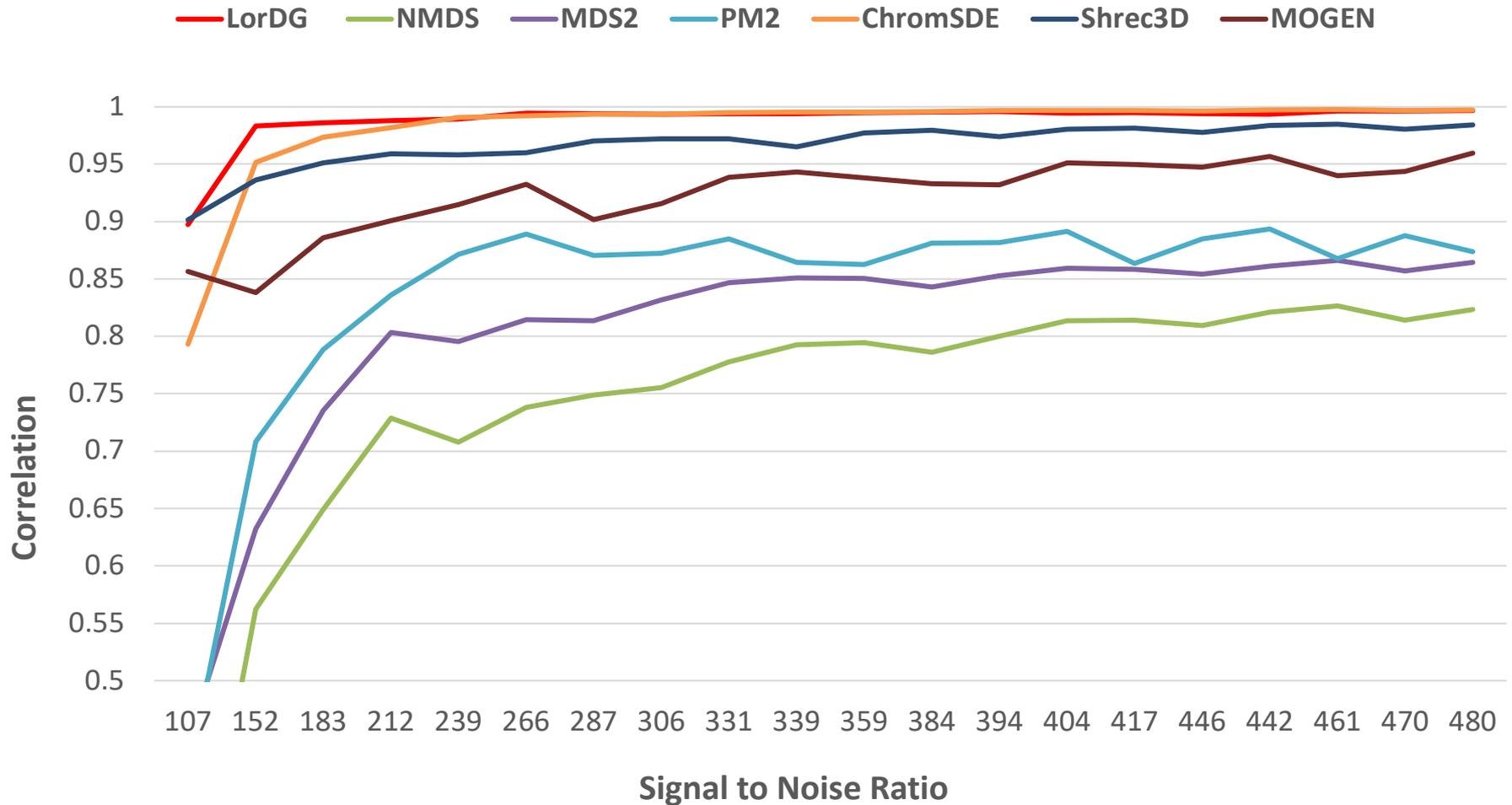
# Large-Scale Adaptive Step-Size Gradient Ascent



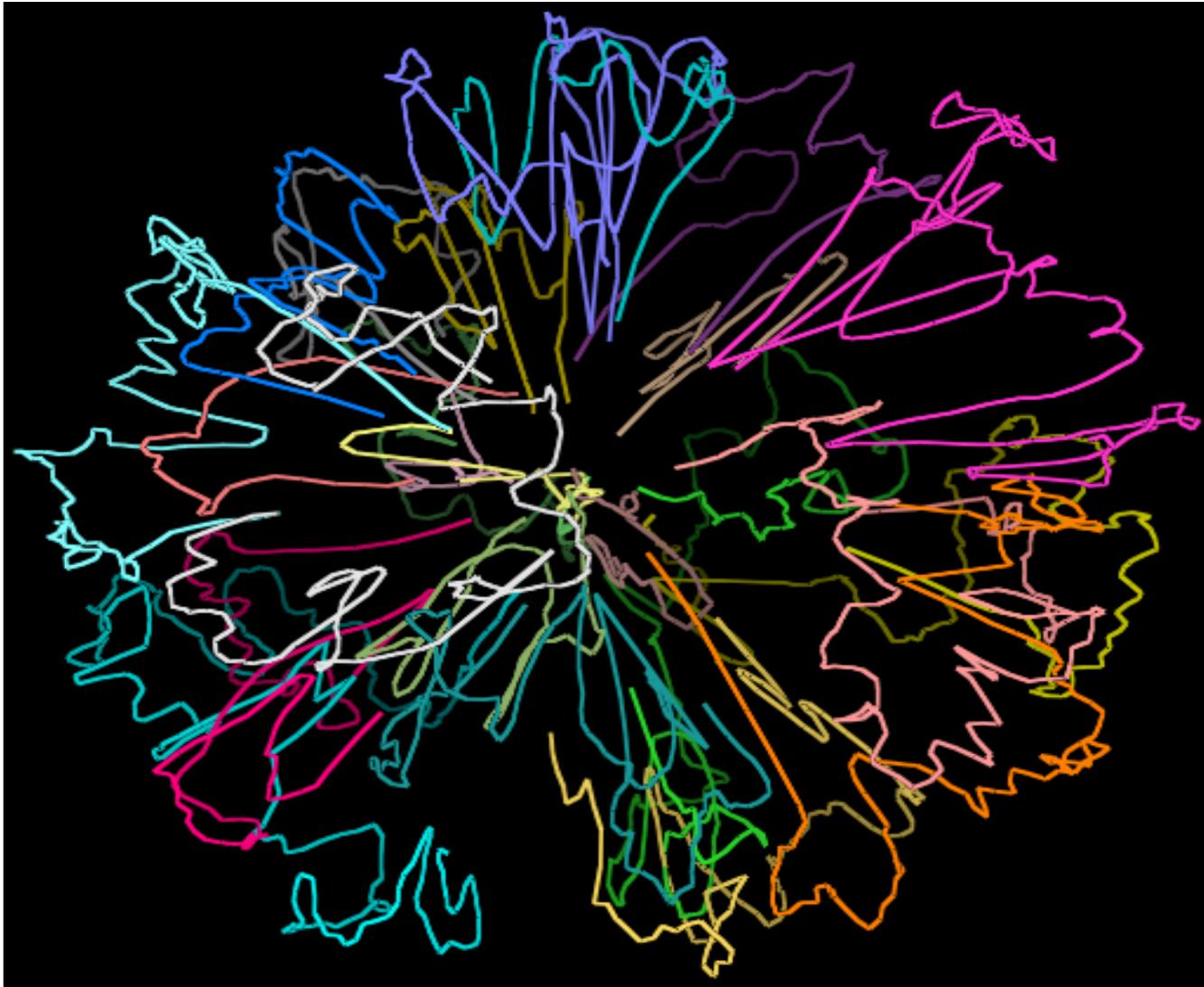
For Educational Use Only



# Correlation Between Reconstructed Distances and True Distances



# 3D Genome Model for Human



# A Open Source Tool - LorDG

- Github: <https://github.com/BDM-Lab/LorDG>

BDM-Lab / LorDG

Watch 2 Star 1 Fork 0

Code Issues 0 Pull requests 0 Projects 0 Insights

Join GitHub today  
GitHub is home to over 20 million developers working together to host and review code, manage projects, and build software together.  
Sign up

3D genome reconstruction with Lorentzian objective function

7 commits 1 branch 1 release 0 contributors

Branch: master New pull request Find file Clone or download

File/Folder	Commit Message	Time
bin	run 1 thread with the main thread	a year ago
examples	initial commit	a year ago
lib	add dependency library	a year ago
make_input	add jar file to make contact matrices from raw HiC data in 2009	4 months ago
src	add dependency library	a year ago
README.md	add readme file	a year ago
paramaters_settings.txt	initial commit	a year ago

README.md

3D Genome Structure Modeling by Lorentzian Objective Function

# Contact-Driven Modeling at Chromosome Scale

- **Input:** initial representation of chromosome, contact map, and physical distance restraints
- **Objective:** find 3D chromosome structures that satisfy the contact map and physical contact restraints as much as possible.
- **Scoring Function**
- **Optimization**
- **Output:** an ensemble of 3D shapes

# Data Preparation

- **Data Sets:** Normal B-Cell, ALL B-Cell
- **Unit Size:** 1Mb
- **Unit Number:** Chr. 1, 248 – Chr. 22, 50
- **Contact map normalization:**  
 $C_{ij}' = C_{ij} / \text{expected IF}$
- **Remove noisy contacts** with low interaction frequency

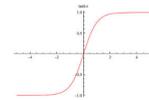
**Normalized Interaction Frequencies (IF) on the normal B-Cell**

Chromosome	Max IF	Min IF	Average IF
1	86.34	0.035	1.35
2	83.04	0.037	1.53
3	59.81	0.037	1.74
4	70.76	0.035	1.95
5	171.92	0.040	1.84
6	59.68	0.031	1.86
7	130.00	0.035	1.99
8	122.22	0.056	2.29
9	190.85	0.043	2.86
10	76.12	0.085	2.31
11	58.51	0.040	2.19
12	85.86	0.042	2.23
13	151.09	0.212	3.62
14	112.48	0.116	3.42
15	103.02	0.110	3.14
16	100.09	0.046	3.31
17	86.08	0.081	3.12
18	108.33	0.173	3.80
19	96.56	0.061	4.82
20	106.89	0.089	4.05
21	63.01	0.174	5.98
22	79.05	0.057	5.83

# Scoring Function for Optimization

$$S = \sum_{\substack{\text{contacts} \\ |i-j| \neq 1}} \left( \tanh(d_c - d(i,j)) * \frac{IF_{ij}}{T} + W1 * \frac{\tanh(d(i,j) - d_{min})}{T} \right) \\ + \sum_{\substack{\text{non-contacts} \\ |i-j| \neq 1}} \left( W2 * \frac{\tanh(d_{max} - d(i,j))}{T} + W3 * \frac{\tanh(d(i,j) - d_c)}{T} \right) \\ + \sum_{|i-j|=1} \left( IF_{max} * \frac{\tanh(da_{max} - d(i,j))}{T} + W1 * \frac{\tanh(d(i,j) - d_{min})}{T} \right)$$

***tanh***: hyperbolic tangent function



***IF<sub>ij</sub>***: interaction frequency between units *i* and *j*

***T***: total interaction frequencies

***d<sub>c</sub>***: contact distance threshold

***d<sub>min</sub>***: minimum distance between two units

***d<sub>max</sub>***: maximum distance between two units

***da<sub>max</sub>***: maximum distance between two adjacent units

***W1, W2, W3, W4***: weight parameters in order to maximize % satisfied contacts + % satisfied non-contacts

**Table S1** The weight parameters and the percent of contacts after removing nose for 23 pairs of chromosomes at 1MB resolution.

Chromosome	Percentage of contact pairs	$W_1$	$W_2$	$W_3$	$W_4$
1	49%	2.0	1.0	1.0	1.0
2	55%	2.0	1.0	1.0	1.0
3	66%	1.2	1.0	1.0	1.0
4	65%	1.0	1.0	1.0	1.0
5	64%	1.0	1.1	1.1	1.1
6	70%	1.0	1.2	1.2	1.2
7	74%	1.0	1.5	1.5	1.5
8	77%	1.0	1.6	1.6	1.6
9	75%	1.0	2.2	2.2	2.2
10	78%	1.0	1.8	1.8	1.8
11	77%	1.0	1.8	1.8	1.8
12	75%	1.0	1.5	1.5	1.5
13	94%	1.0	4.6	4.6	4.6
14	86%	1.0	3.5	3.5	3.5
15	89%	1.0	3.0	3.0	3.0
16	86%	1.0	3.9	3.9	3.9
17	85%	1.0	3.0	3.0	3.0
18	97%	1.0	9.5	9.5	9.5
19	92%	1.0	4.0	4.0	4.0
20	93%	1.0	5.0	5.0	5.0
21	94%	1.0	6.0	6.0	6.0
22	96%	1.0	6.0	6.0	6.0
23	83%	1.0	2.0	2.0	2.0

# Estimating Parameters

- FISH data (*Mateos-Langerak et al., 2009*) for physical distances of Chr. 1 and 11 at various genomic distances
- $d_{\min}$ ,  $d_{\max}$ : min and max distance between pairs at 1Mb away. (0.2  $\mu\text{m}$  , 1.8  $\mu\text{m}$ )
- $d_{\max}$ : max distance between all pairs
- $d_c$ : a threshold resulting in the same percent of contacts in our data
- $d_c$  and  $d_{\max}$  are chromosome length dependent (1.73 – 2.24  $\mu\text{m}$ , 2.45 – 3.32  $\mu\text{m}$ )

# Steepest Gradient Ascent with Backtracking Line Search

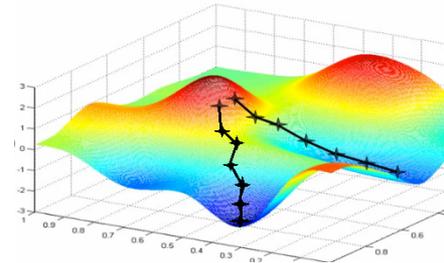
- **Random Initialization:**  $(x_1^0, y_1^0, z_1^0), (x_2^0, y_2^0, z_2^0), \dots, (x_N^0, y_N^0, z_N^0)$  in  $[-0.5, 0.5]$
- **Update:**

$$X_1^{t+1} = X_1^t + \eta * \Delta X \quad Y_1^{t+1} = Y_1^t + \eta * \Delta Y \quad Z_1^{t+1} = Z_1^t + \eta * \Delta Z$$

$$X_2^{t+1} = X_2^t + \eta * \Delta X \quad Y_2^{t+1} = Y_2^t + \eta * \Delta Y \quad Z_2^{t+1} = Z_2^t + \eta * \Delta Z$$

⋮

$$X_N^{t+1} = X_N^t + \eta * \Delta X \quad Y_N^{t+1} = Y_N^t + \eta * \Delta Y \quad Z_N^{t+1} = Z_N^t + \eta * \Delta Z$$



*Step size ( $\eta$ ) is adjusted dynamically during iterations to avoid too big moves.*

# Structure Modeling Movie

# Structure Modeling Movie

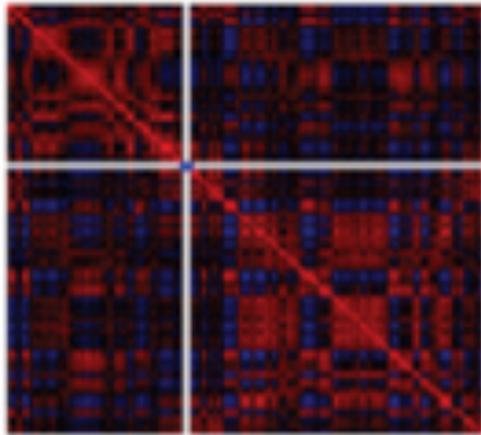
At YouTube without music: <http://www.youtube.com/watch?v=C03R7A9kYc8>

At NSF CAREER project web site with music:

[http://people.cs.missouri.edu/~chengji/genome\\_modeling\\_movie.mp4](http://people.cs.missouri.edu/~chengji/genome_modeling_movie.mp4)

**J. Cheng, NSF CAREER Project Plan, 2011, 2012, 2013. T. Tuan made the movie.**

# Two Compartment Validation on Normal B-Cell



Chromosome 7

Purple and green denote regions in two different components using principle component analysis on contact correlation map

# Model Selection

- **Use TM-Score** (Zhang & Skolnick, 2004) **to superpose every pair of models in an ensemble of models**
- **Calculate GDT-HA score: percent of unit pairs within specific distance thresholds**
- **Choose centroid model as representative**

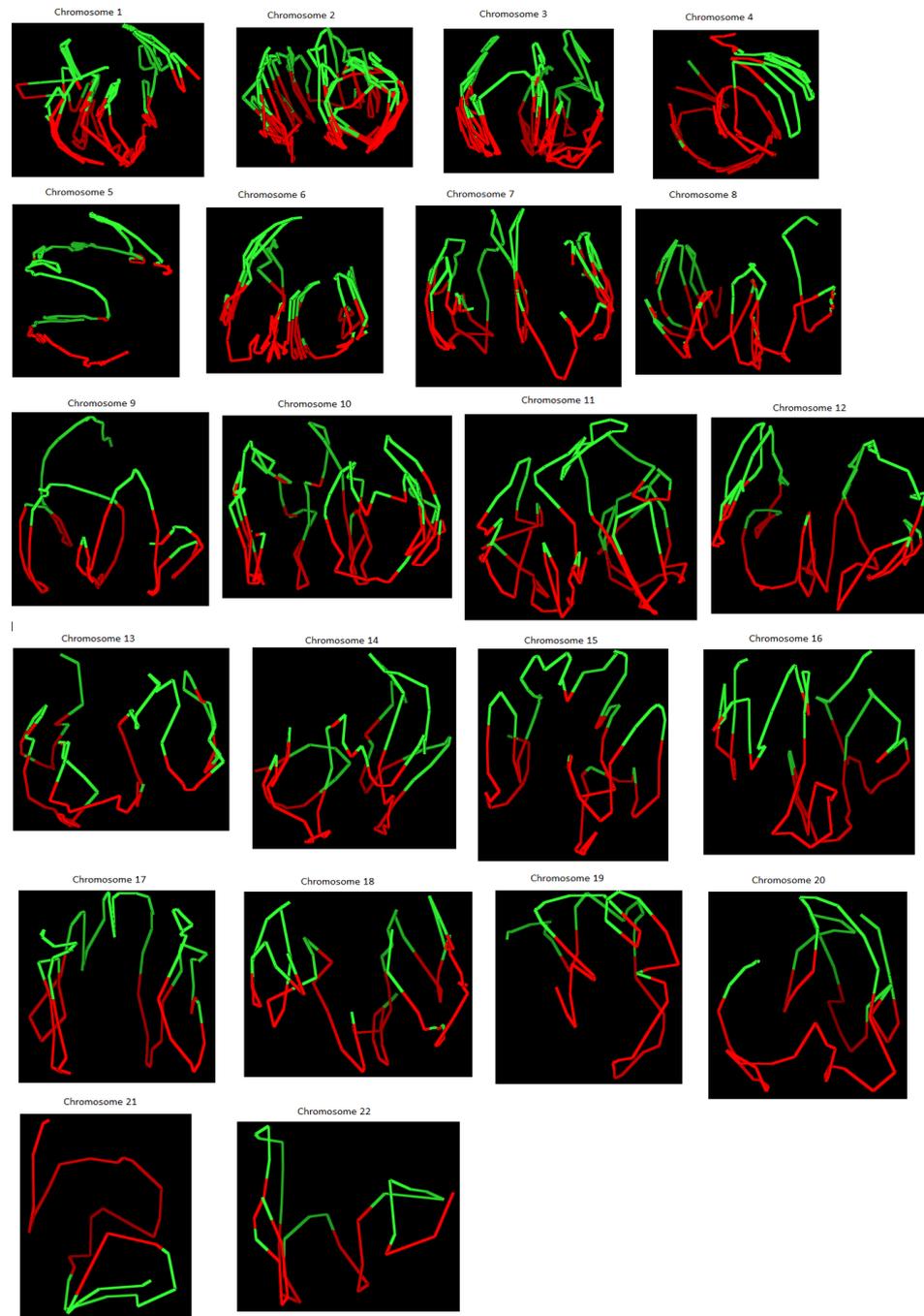
# Satisfaction of Contacts in Representative Models of Normal B-Cell

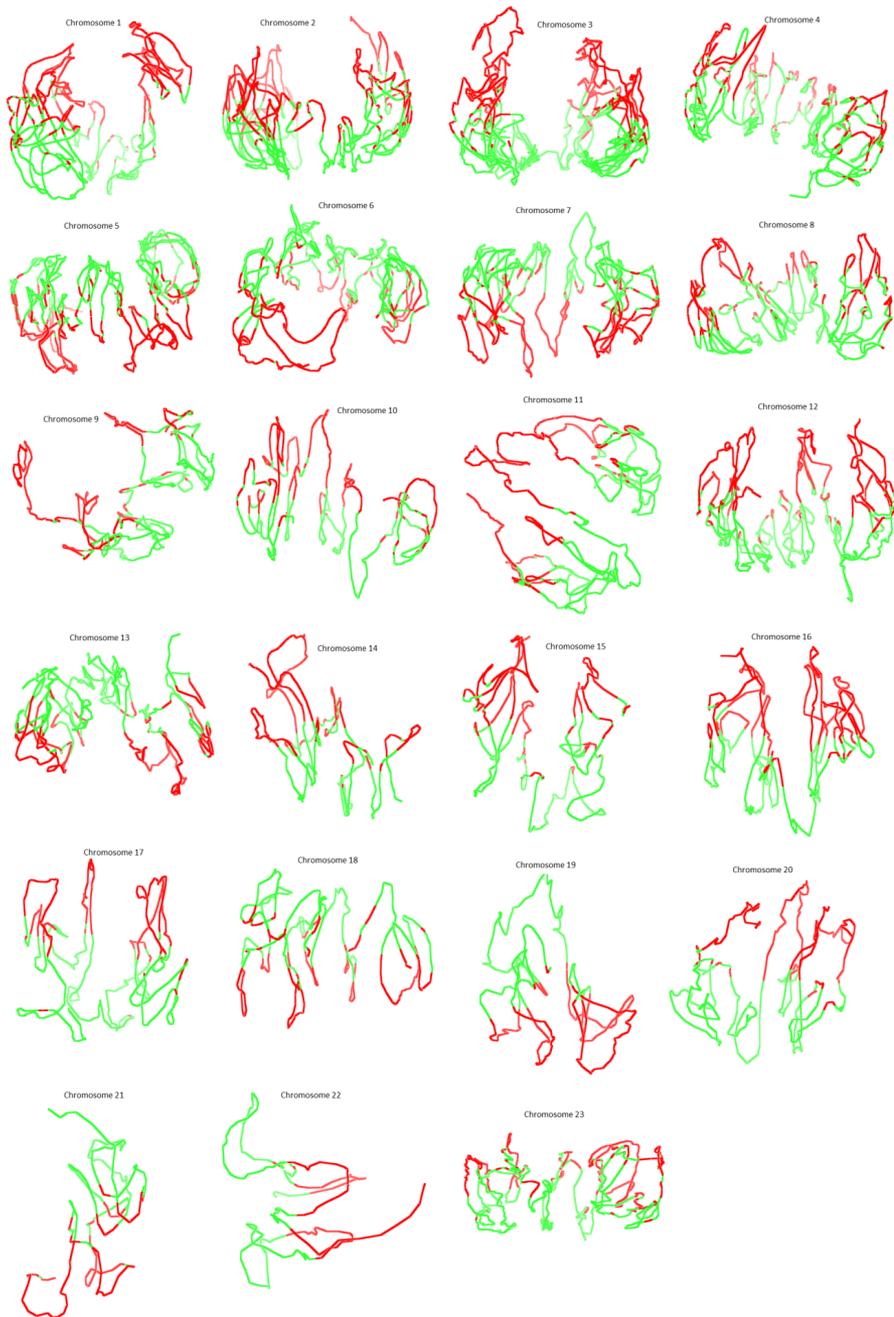
Chromosome	Satisfied contact pairs (%)	Satisfied Non-contact pairs (%)
1	82	80
2	81	80
3	82	84
4	85	85
5	83	86
6	83	83
7	86	82
8	89	83
9	88	90
10	87	88
11	86	86
12	88	88
13	90	89
14	93	84
15	91	90
16	89	88
17	90	92
18	89	94
19	92	96
20	90	90
21	92	99
22	91	92

# Violations of Contacts on Normal B-Cell

Chromosome	Average distance (unsatisfied contact pairs)	Average IF (unsatisfied contact pairs)	Average distance (unsatisfied non-contact pairs)	$d_c$	Average IF
1	1.98	0.44	1.43	1.73	1.6
2	1.97	0.52	1.48	1.73	1.84
3	1.94	0.60	1.51	1.73	2.09
4	2.05	0.57	1.31	1.73	2.36
5	2.03	0.55	1.39	1.73	2.22
6	1.97	0.60	1.42	1.73	2.23
7	2.37	0.62	1.68	2.00	2.39
8	2.34	0.69	1.60	2.00	2.76
9	2.42	0.72	1.72	2.12	3.49
10	2.58	0.71	1.92	2.23	2.78
11	2.56	0.71	1.80	2.23	2.6
12	2.54	0.67	1.89	2.23	2.69
13	2.55	1.32	1.75	2.23	4.34
14	2.47	1.15	1.69	2.23	4.12
15	2.46	1.02	1.85	2.23	3.78
16	2.58	0.91	1.79	2.34	4.02
17	2.56	0.93	2.06	2.34	3.78
18	2.66	1.36	1.72	2.34	4.54
19	2.57	1.40	1.84	2.34	5.85
20	2.56	1.35	1.78	2.34	4.88
21	2.62	1.60	1.93	2.34	7.28
22	2.58	1.78	1.64	2.34	7.05

# 3D Models for 22 Chromosomes of Normal B- Cell



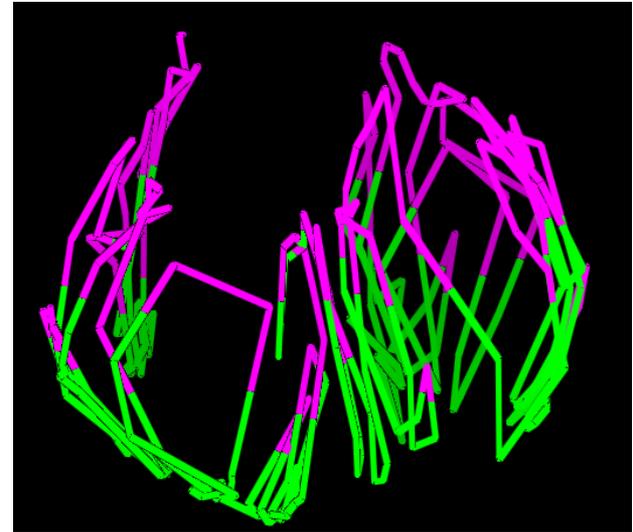


Models:  
200KB resolution

# Models of Chr. 2 for Normal and Malignant Cells

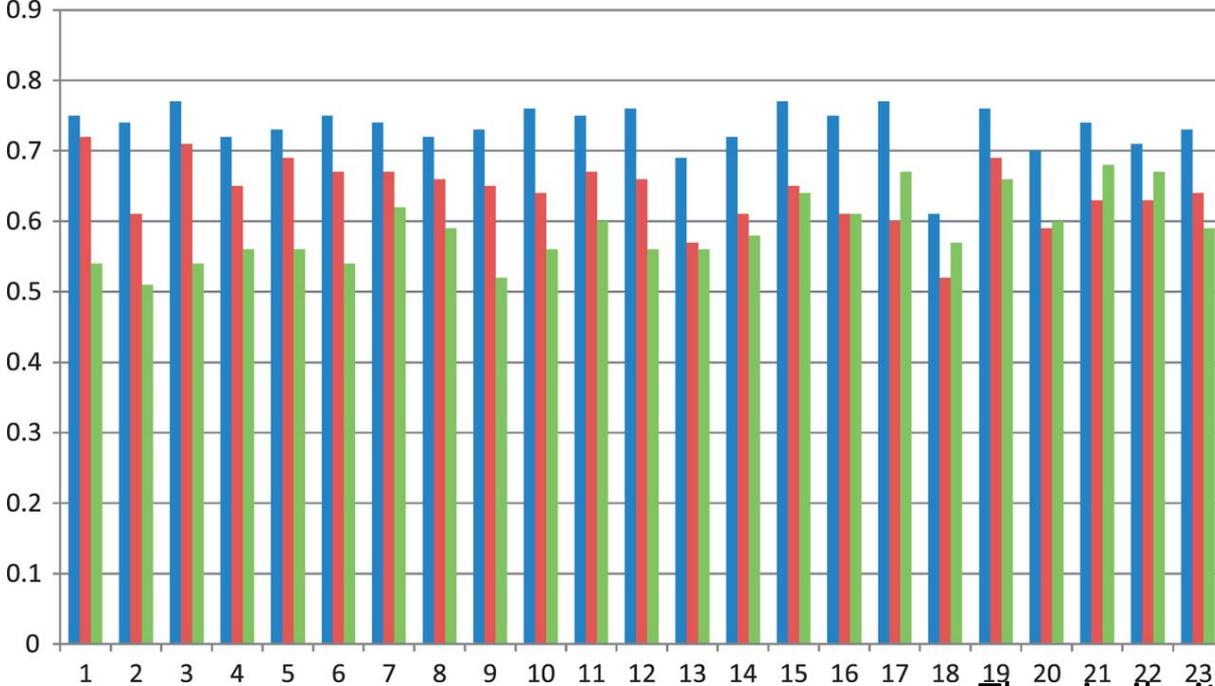


**Primary Leukemia B-Cell**



**Normal B-Cell**

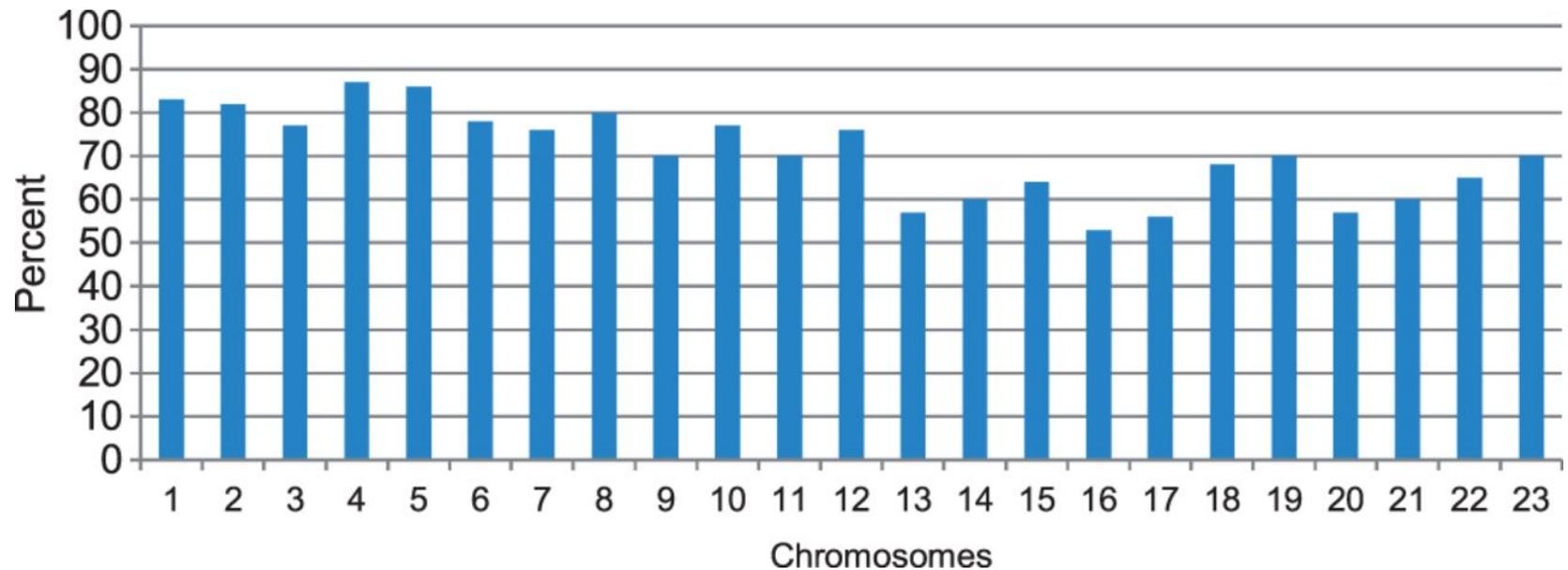
# Model Consistency



The similarity scores measured as average GDT-HA scores. Y-axis denotes the similarity scores and X-axis the indices of chromosomes. 'Blue bars' represent the average GDT-HA scores of models within the same model ensemble constructed from the whole normalized data sets of the normal B-cell for each chromosome, 'red bars' the average GDT-HA scores between models constructed from sampled data sets with those constructed from the whole data sets and 'green bars' the average GDT-HA scores between models of the leukemia B-cell and those of the normal B-cell.

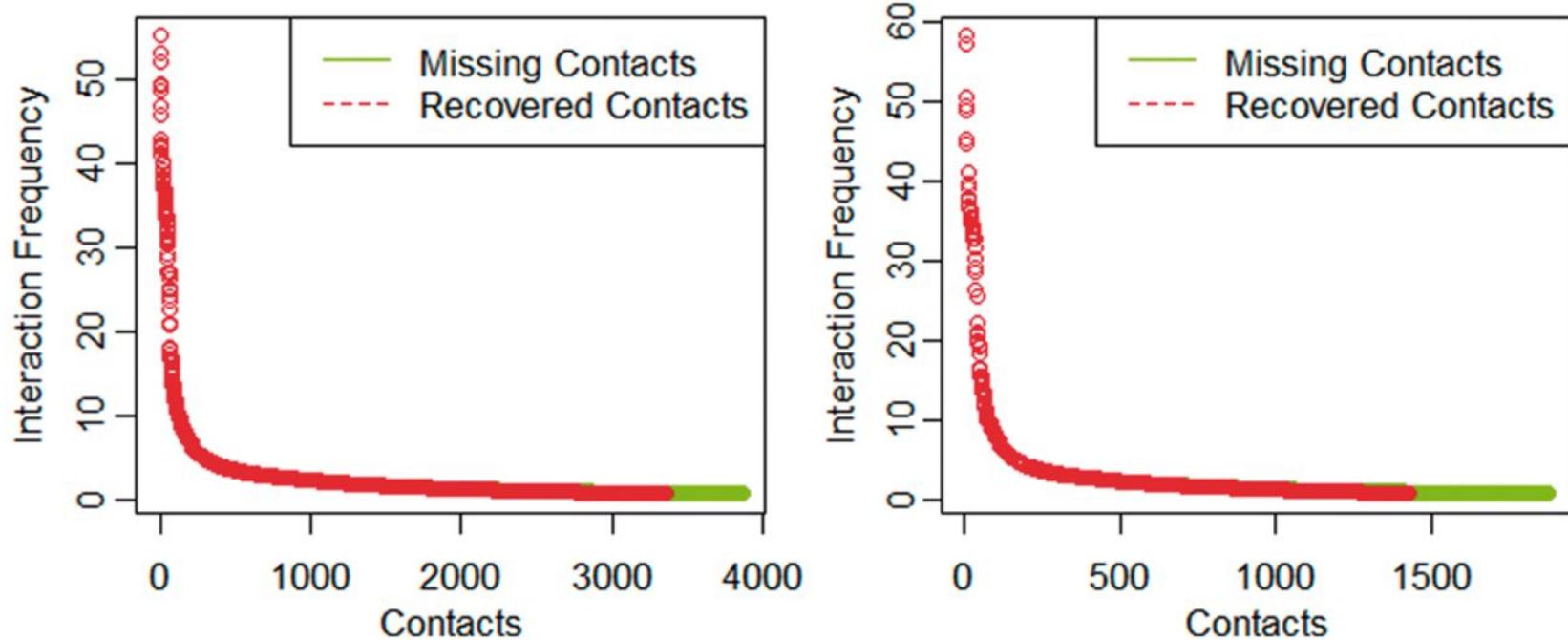
Trieu T , and Cheng J Nucl. Acids Res. 2014;nar.gkt1411

# The percentage of recovered contacts in all chromosomes in validation



Trieu T , and Cheng J Nucl. Acids Res. 2014;nar.gkt1411

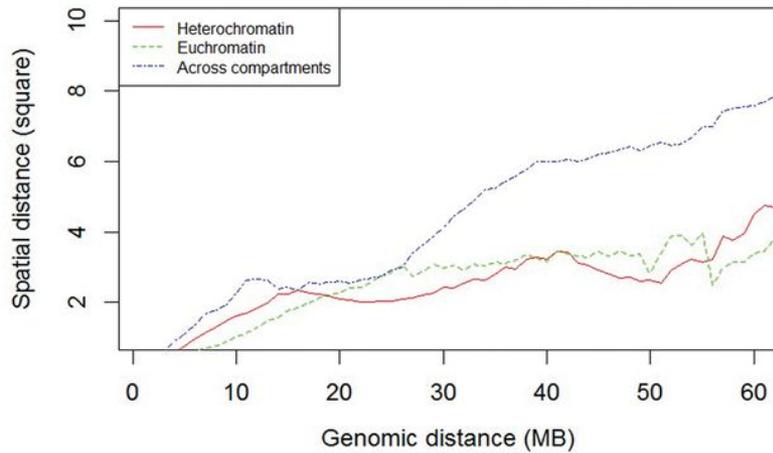
The superimposition of the plots of IFs of all missing contacts and the plots of IFs of recovered contacts for chromosome 1 (left) and chromosome 11 (right).



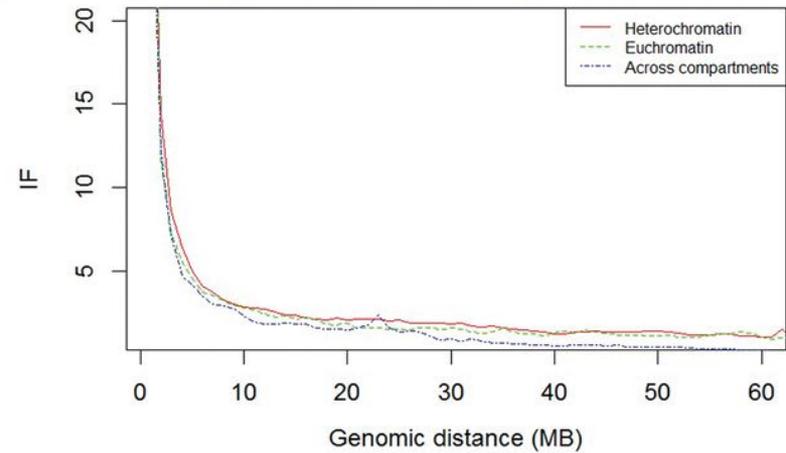
The superimposition of the plots of IFs of all missing contacts and the plots of IFs of recovered contacts for chromosome 1 (left) and chromosome 11 (right). Y-axis denotes interaction frequencies and X-axis the indices of contacts. The green tails visualized the contacts that were not recovered and had lower IFs.

# Average of spatial distances and IFs of region pairs within and across (between) compartments in chromosomes 1 and 11.

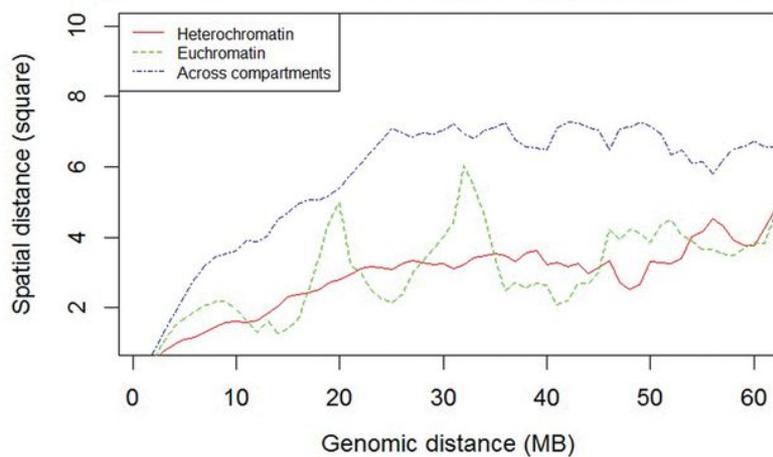
### Distances within and across two compartments in Chr.1



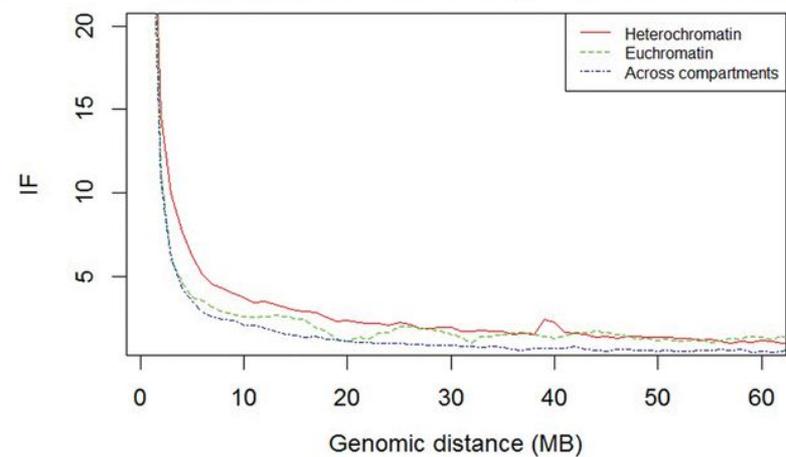
### IF within and across two compartments in Chr.1



### Distances within and across two compartments in Chr.11

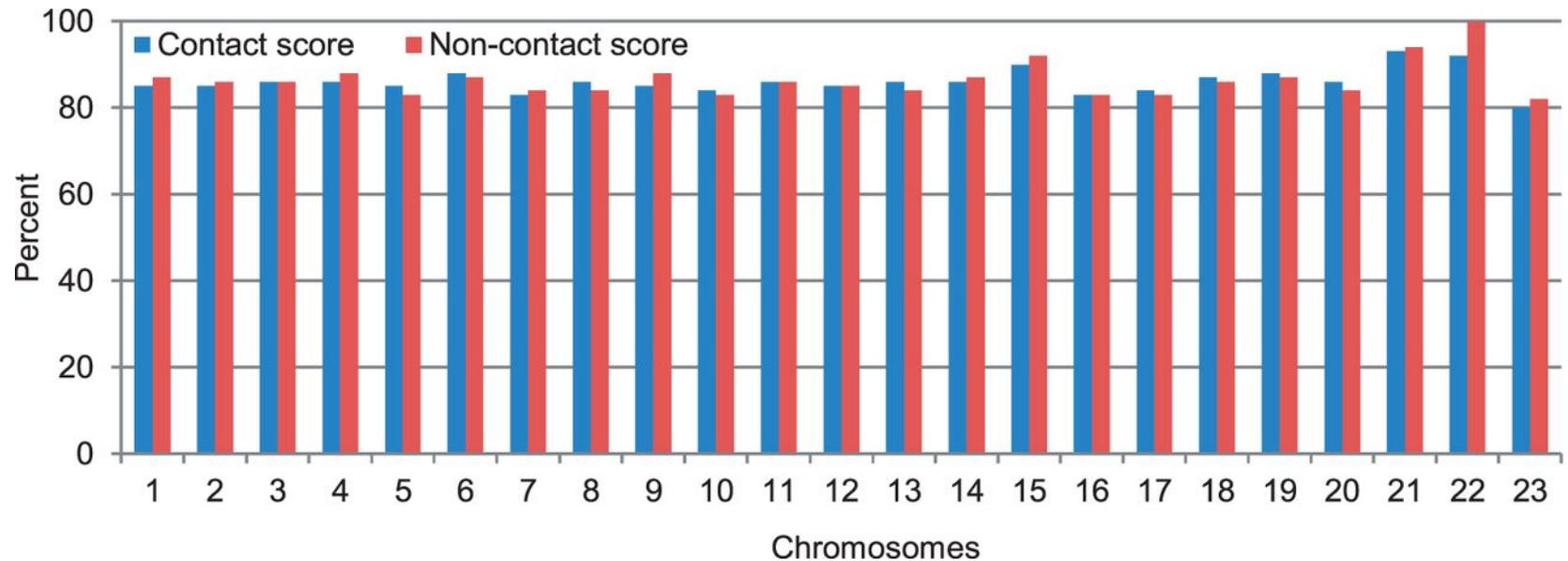


### IF within and across two compartments in Chr.11



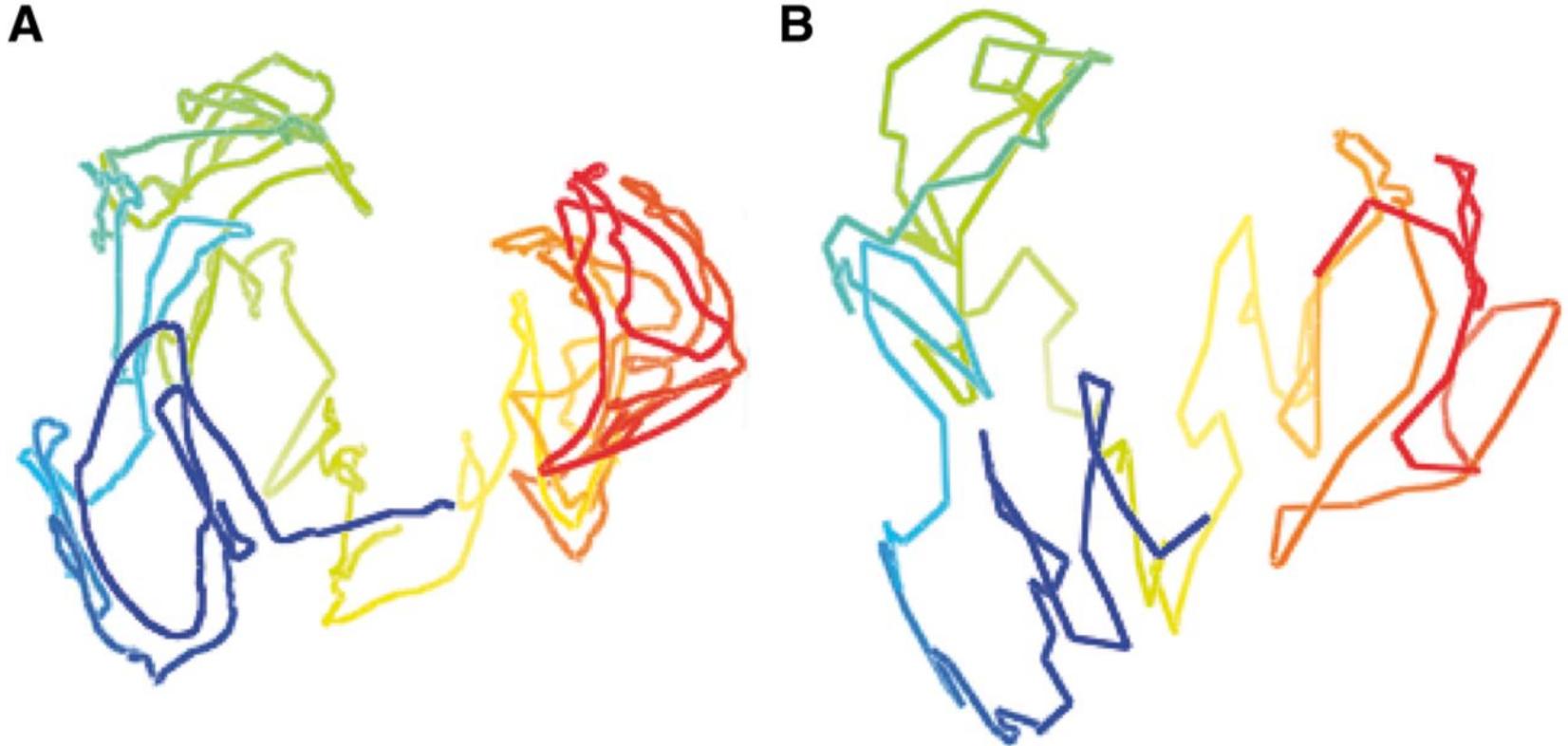
Trieu T , and Cheng J Nucl. Acids Res. 2014;nar.gkt1411

# The contact scores and non-contact scores of the chromosomal models of the leukemia B-cell.

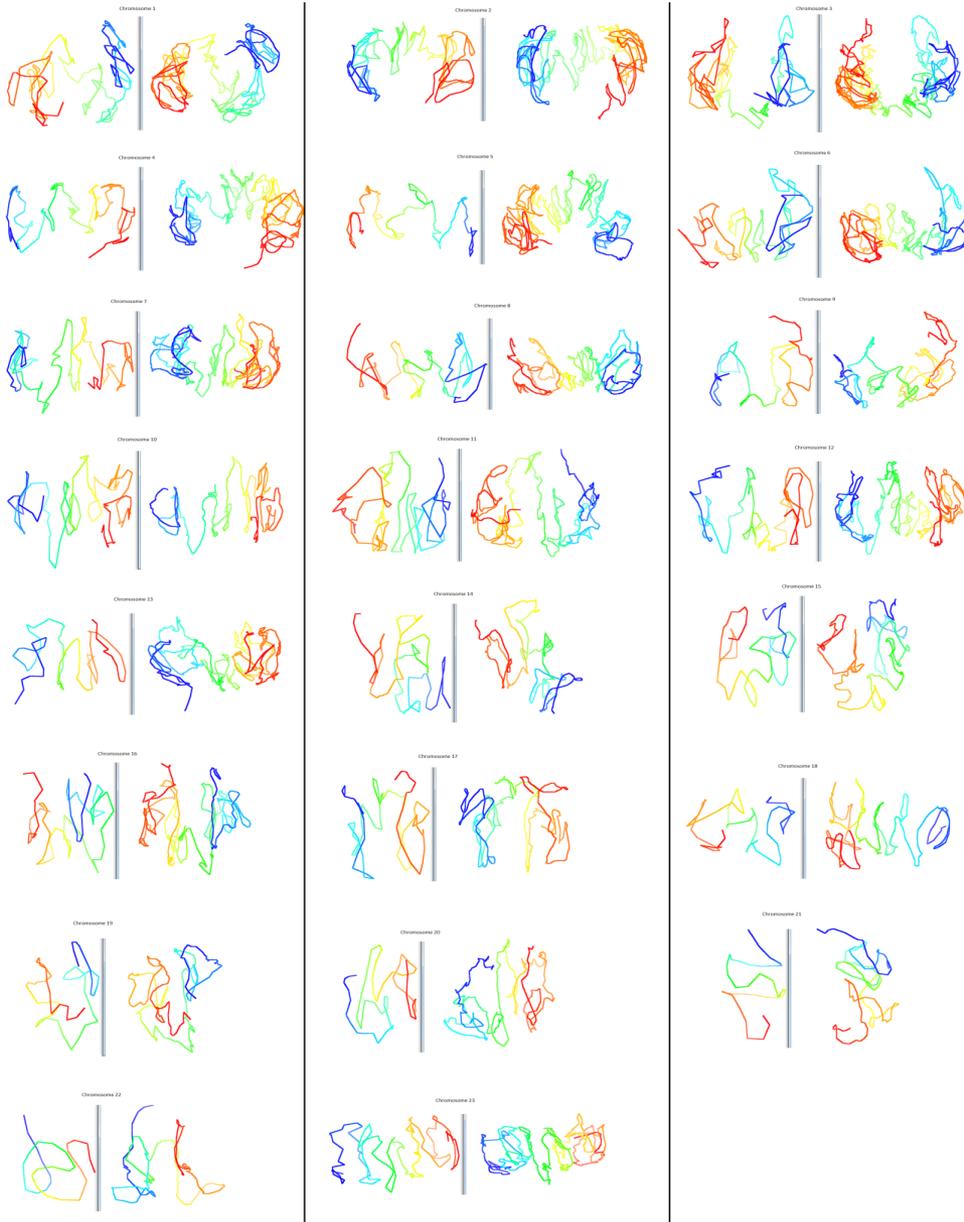


Trieu T , and Cheng J Nucl. Acids Res. 2014;nar.gkt1411

The models of chromosome 1 at resolution of 200 K (A) and at resolution of 1 MB (B).



Trieu T , and Cheng J Nucl. Acids Res. 2014;nar.gkt1411



Pairwise comparison:

1Mb versus 200Kb

# Scoring Function for Genome Optimization

$$\begin{aligned}
 Fn = & \sum_{\{(i,j): |i-j| \neq 1\}} \text{contacts} \left( W_1[\text{chr1}, \text{chr2}] * \tanh(d_c^2 - d_{ij}^2) * \frac{F_{ij}}{\text{totalIF}} + \right. \\
 & \left. W_2[\text{chr1}, \text{chr2}] * \frac{\tanh(d_{ij}^2 - d_{min}^2)}{\text{totalIF}} \right) + \\
 & \sum_{\{(i,j): |i-j|=1 \ \& \ \text{chr1}=\text{chr2}\}} \left( W_1[\text{chr1}, \text{chr2}] * IF_{max} * \frac{\tanh(da_{max}^2 - d_{ij}^2)}{\text{totalIF}} + \right. \\
 & \left. W_2[\text{chr1}, \text{chr2}] * \frac{\tanh(d_{ij}^2 - d_{min}^2)}{\text{totalIF}} \right) + \\
 & \sum_{\{(i,j): |i-j| \neq 1, \text{chr1}=\text{chr2}\}} \text{non-contacts} \left( W_3[\text{chr1}, \text{chr1}] * \frac{\tanh(d_{max\_intra}^2 - d_{ij}^2)}{\text{totalIF}} + \right. \\
 & \left. W_4[\text{chr1}, \text{chr1}] * \frac{\tanh(d_{ij}^2 - d_c^2)}{\text{totalIF}} \right) + \\
 & \sum_{\{(i,j): |i-j| \neq 1, \text{chr1} \neq \text{chr2}\}} \text{non-contacts} \left( W_3[\text{chr1}, \text{chr2}] * \frac{\tanh(d_{max\_inter}^2 - d_{ij}^2)}{\text{totalIF}} + \right. \\
 & \left. W_4[\text{chr1}, \text{chr2}] * \frac{\tanh(d_{ij}^2 - d_c^2)}{\text{totalIF}} \right)
 \end{aligned}$$

**Contact Pairs**

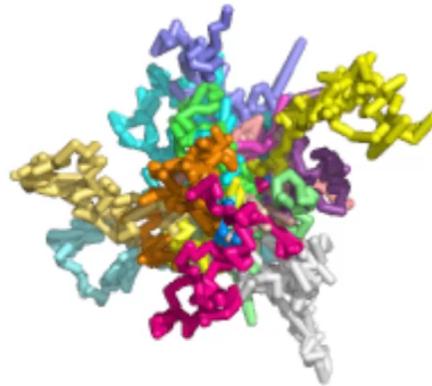
**Adjacent Pairs**

**Non-contacts on same chromosome**

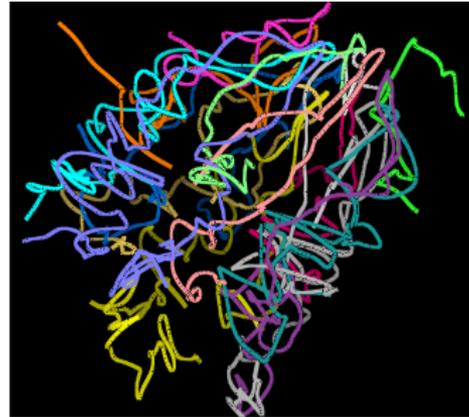
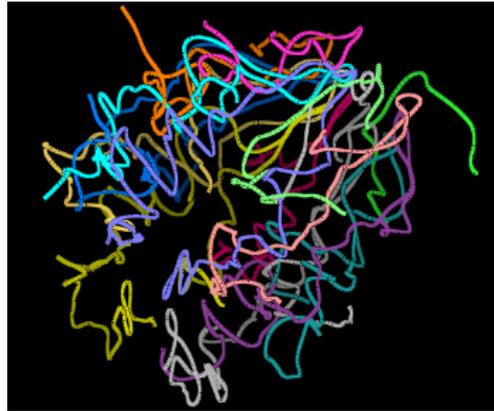
**Non-contacts on different chromosomes**

$$\begin{aligned}
 d_{min}^2 &= 0.2 \mu m^2, \quad da_{max}^2 = 1.8(\mu m^2) \\
 d_c^2 &= 6\mu m^2, \quad d_{max\_intra}^2 = d_{max}^2 = 20\mu m^2 \\
 d_{max\_inter} &= 13\mu m
 \end{aligned}$$

# Modeling of 3D Genome

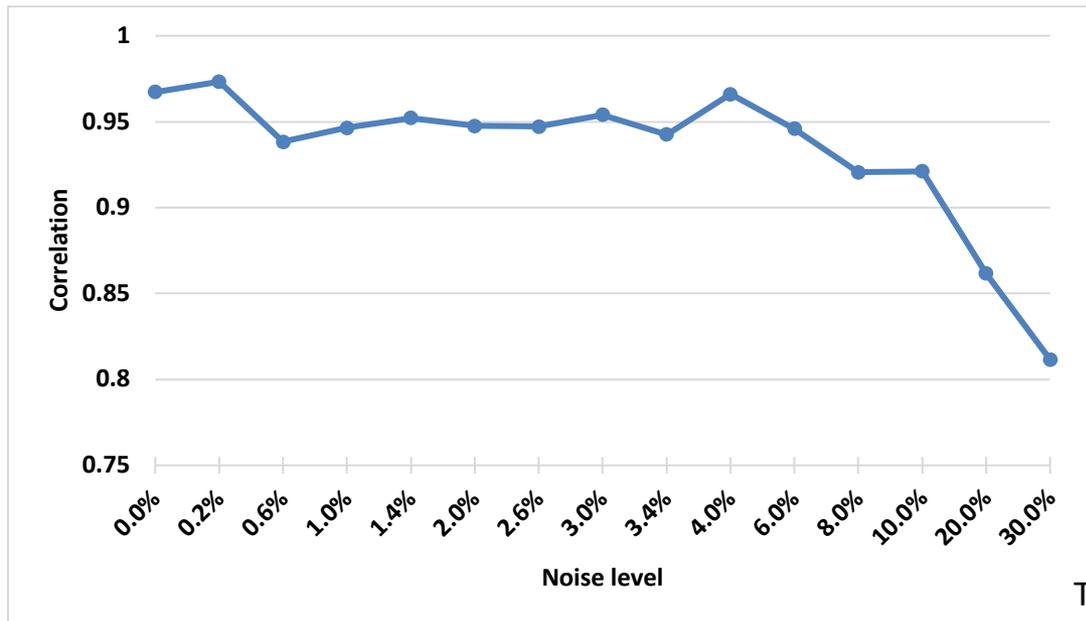


# Test on Simulated Data of the Yeast Genome Structure

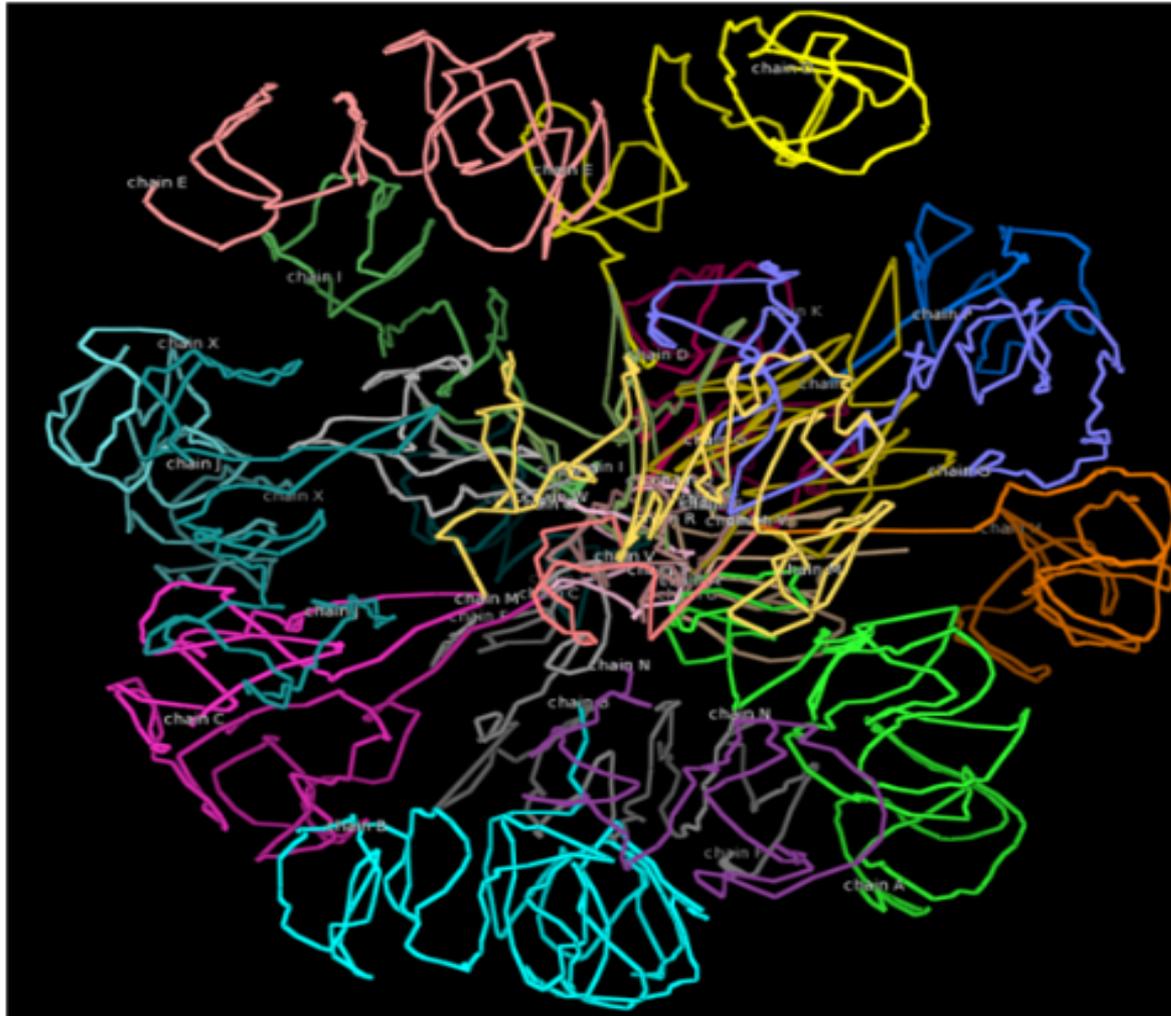


Known Structure (Duan et al., 2010)

Reconstructed Model



# A 3D Model of Genome of Human B-Cell in an Ensemble



# Largely Conserved Chromosomal Structure, Dynamic Genome Structure

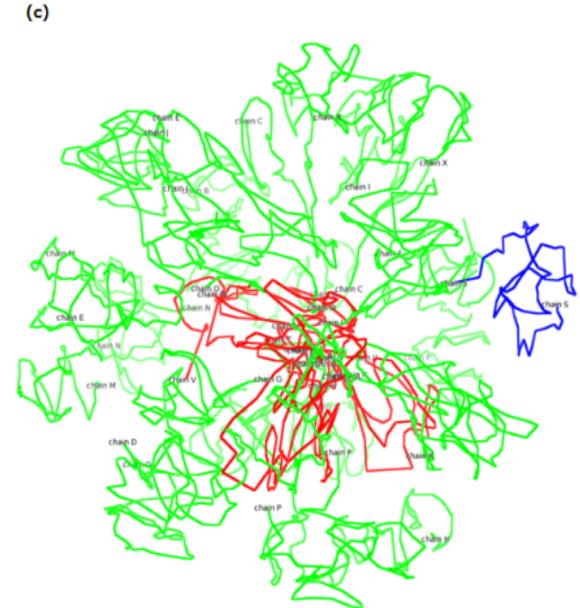
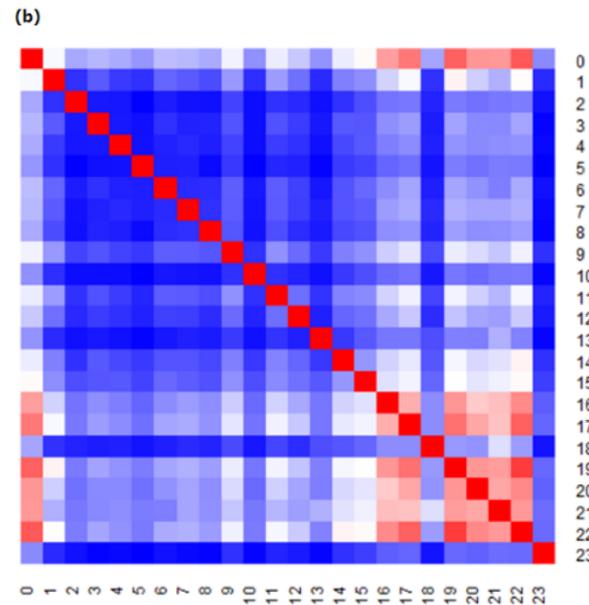
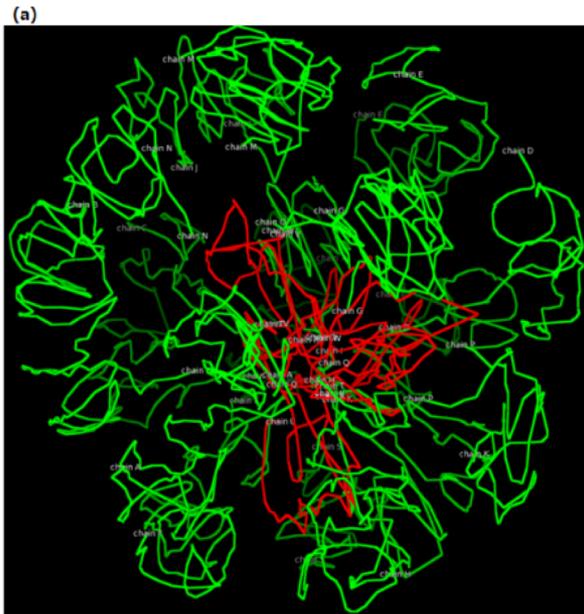


**Chromosome 11**

**Pairwise structural similarity between chromosomal models: 0.71**

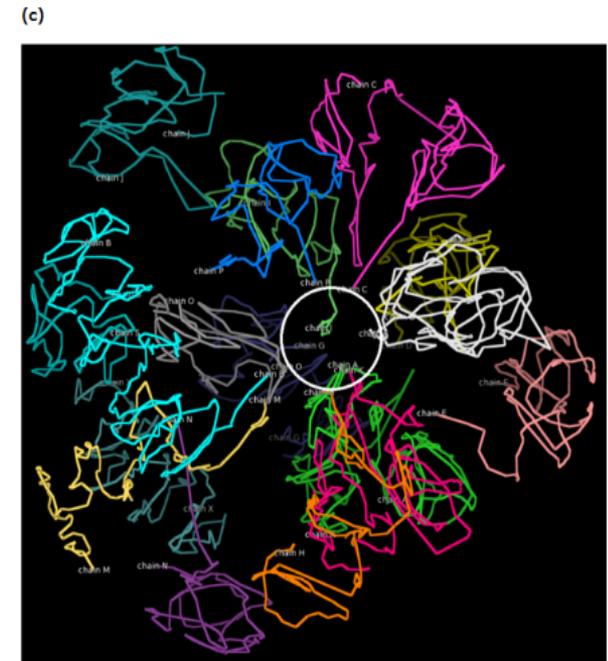
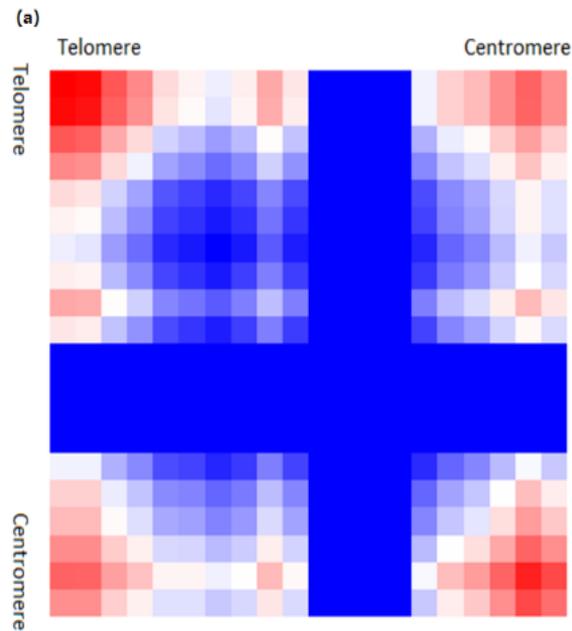
**Pairwise structural similarity between genome models: 0.16**

# Center – Periphery Architecture and Dynamic Inter-chromosomal Interactions



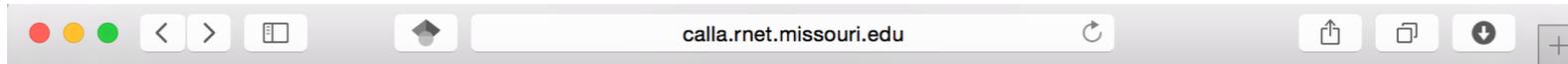
- Small chromosomes in center
- Large chromosomes in periphery
- Genome structure ensemble (dynamics)
- Correlation of inter-chromosomal interactions (0.48), p-value <  $e-16$

# Enriched Telomere and Centromere Inter-Chromosomal Interactions



- **Telomere-telomere interactions > telomere/centromere interactions > other interactions**

# MOGEN: Model of Genome (Tool)



1. [Reconstruction process video: http://calla.rnet.missouri.edu/mogen/video/3DGenome\\_Movie.mp4](http://calla.rnet.missouri.edu/mogen/video/3DGenome_Movie.mp4)



2. [Genome structures of healthy cells: http://calla.rnet.missouri.edu/mogen/normal\\_cell/](http://calla.rnet.missouri.edu/mogen/normal_cell/)
3. [Genome structures of malignant cells: http://calla.rnet.missouri.edu/mogen/cancer\\_cell/](http://calla.rnet.missouri.edu/mogen/cancer_cell/)
4. [Experimental data on synthetic datasets: https://missouri.box.com/s/5cnwys9e1qgt8o20pdh8544meb5myay0](https://missouri.box.com/s/5cnwys9e1qgt8o20pdh8544meb5myay0)
5. [MOGEN: Executable program - https://missouri.box.com/s/a8nynitxow4aixj6fd0k2xf6n4wof772](https://missouri.box.com/s/a8nynitxow4aixj6fd0k2xf6n4wof772)

<http://calla.rnet.missouri.edu/mogen/>

# MOGEN at GitHub

- <https://github.com/BDM-Lab/MOGEN>
- Trieu, Tuan, and Jianlin Cheng. "MOGEN: a tool for reconstructing 3D models of genomes from chromosomal conformation capturing data." *Bioinformatics* 32.9 (2015): 1286-1292.

Secure | <https://github.com/BDM-Lab/MOGEN>

The software for modeling the 3D structure of a genome using Hi-C chromosome conformation capturing data

19 commits   1 branch   1 release   1 contributor   GPL-3.0

Branch: master ▾   New pull request   Find file   Clone or download ▾

**tuan** add buildfile   Latest commit 07e31c6 on Mar 10, 2017

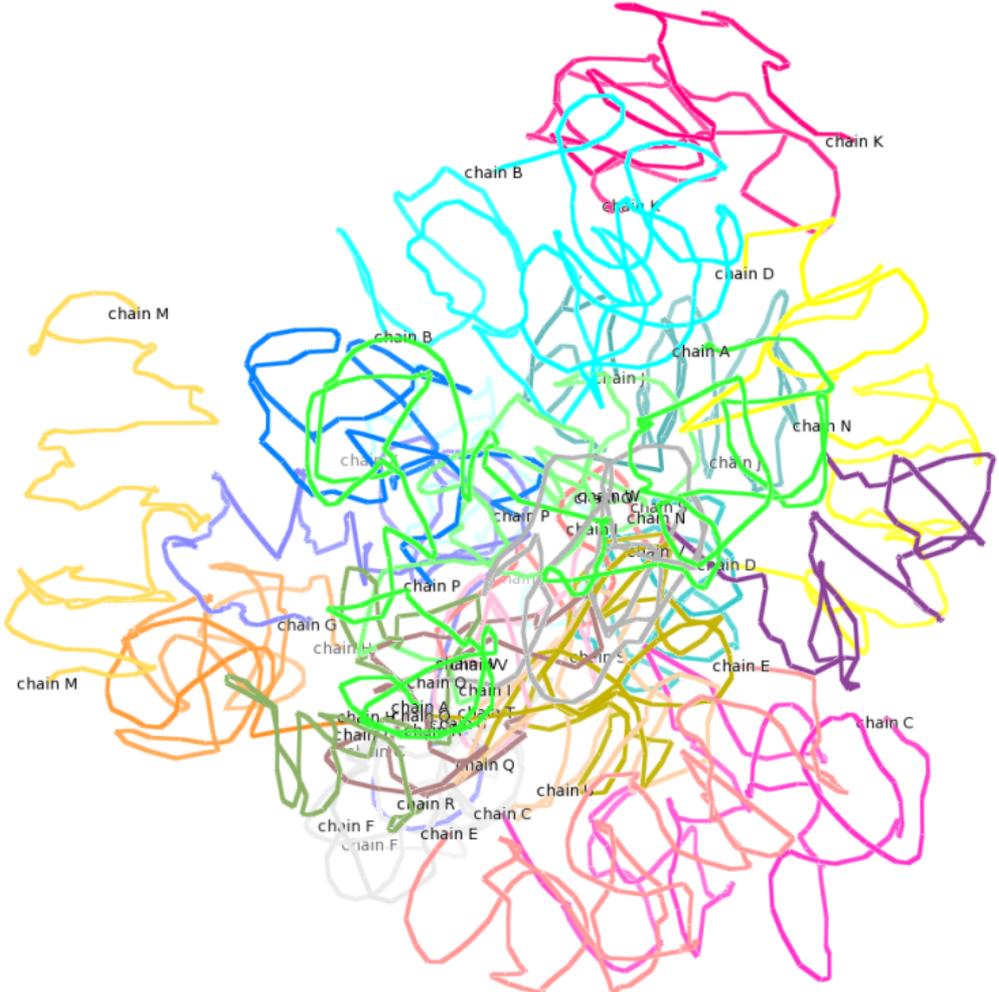
bin	MOGEN source code	2 years ago
documents	add supplementary document	a year ago
examples/hiC	remove file > 100 mb	2 years ago
src	add supplementary document	a year ago
LICENSE	Initial commit	2 years ago
README.md	add question 3	a year ago
build.xml	add buildfile	a year ago

README.md

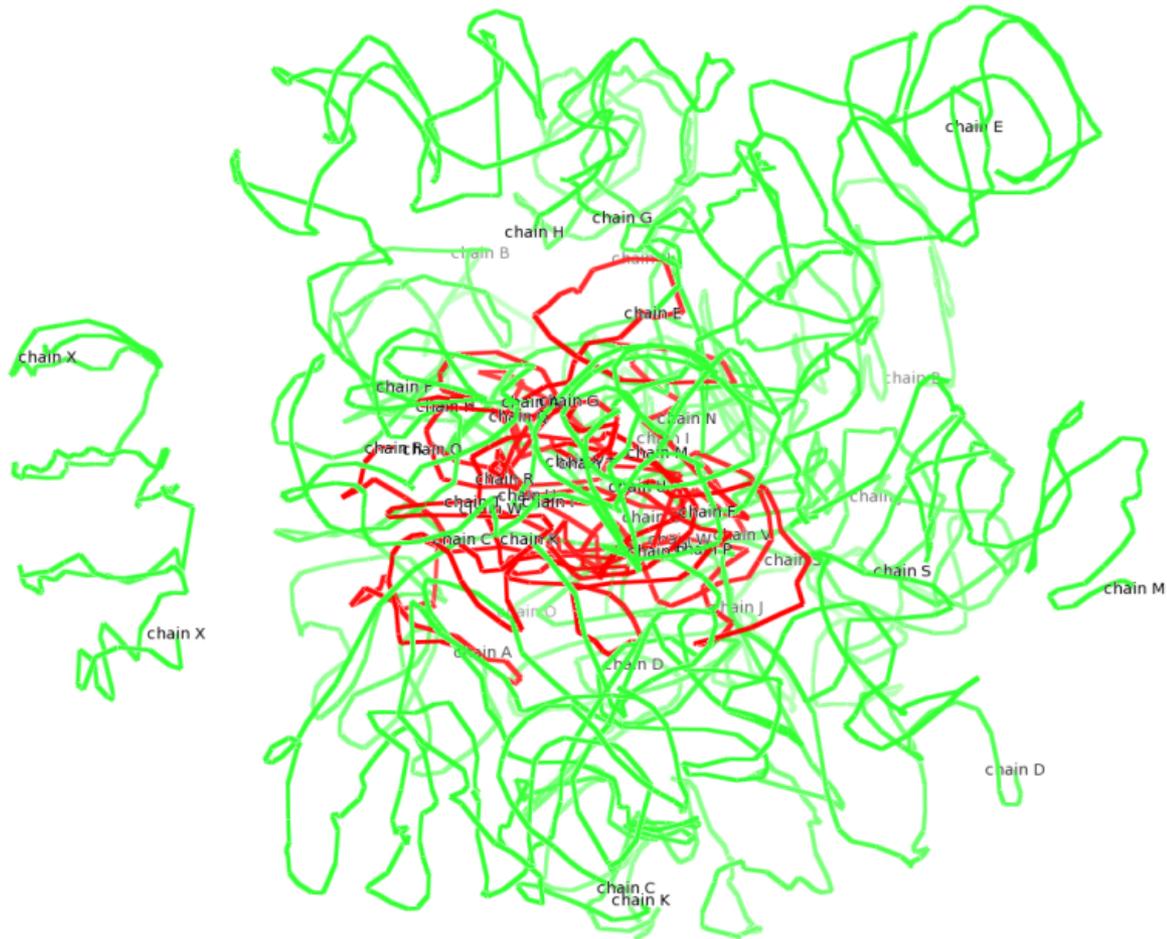
---

**MOGEN: The software for modeling the 3D structure of a genome using Hi-C chromosome conformation capturing data**

# 3D Structure of the Entire Genome



# Smaller Chromosomes are Closer to the Center



★ [Add a Review](#)

↓ [1 Download](#) (This Week)

📅 Last Update: 2014-02-28

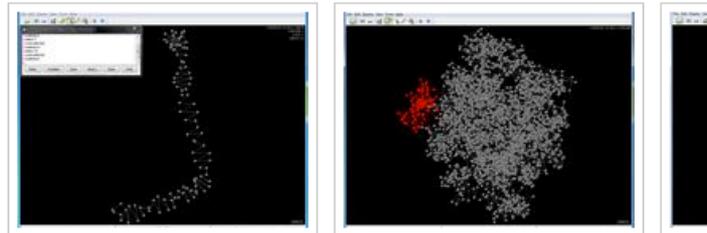
[Tweet](#) 1

[g+](#) 0

[Like](#) 2



[Browse All Files](#)



## Description

GMOL is an application designed to visualize genome structure in 3D. It allows users to view the genome structure at multiple scales, including: global, chromosome, loci, fiber, nucleosome, and nucleotide. This software was built upon the pre-existing Jmol package.

[GMOL Web Site](#) >

### Categories

[Bio-Informatics](#), [Visualization](#)

### License

[GNU Library or Lesser General Public License version 3.0 \(LGPLv3\)](#)

## Features

- Interactively visualize genome structures in 3D
- Supports multiple scales/resolutions: global, chromosome, loci, fiber, nucleosome, nucleotide
- Measure distances and angles between points in the structure
- Rotate and scale models to analyze them even more
- Select parts of the structure based on index, scale, or sequence information
- Get DNA sequence for selected structures or portions of structures
- Includes existing Jmol functions

### Featured Download

#### Why attend IBM Edge 2014?

750 technical sessions & over 5500 of your peers. Sponsored by Intel@



[Register](#)



**PROGRESS**  
**GOT A PAIN IN THE SaaS?**  
We've got your SaaS data access cure.  
**FREE TRIAL**

## Recommended Projects

[Jmol](#)

[JSpeView Project](#)

[Open Virtual Machine Tools](#)

## Latest Tech Jobs

Powered by

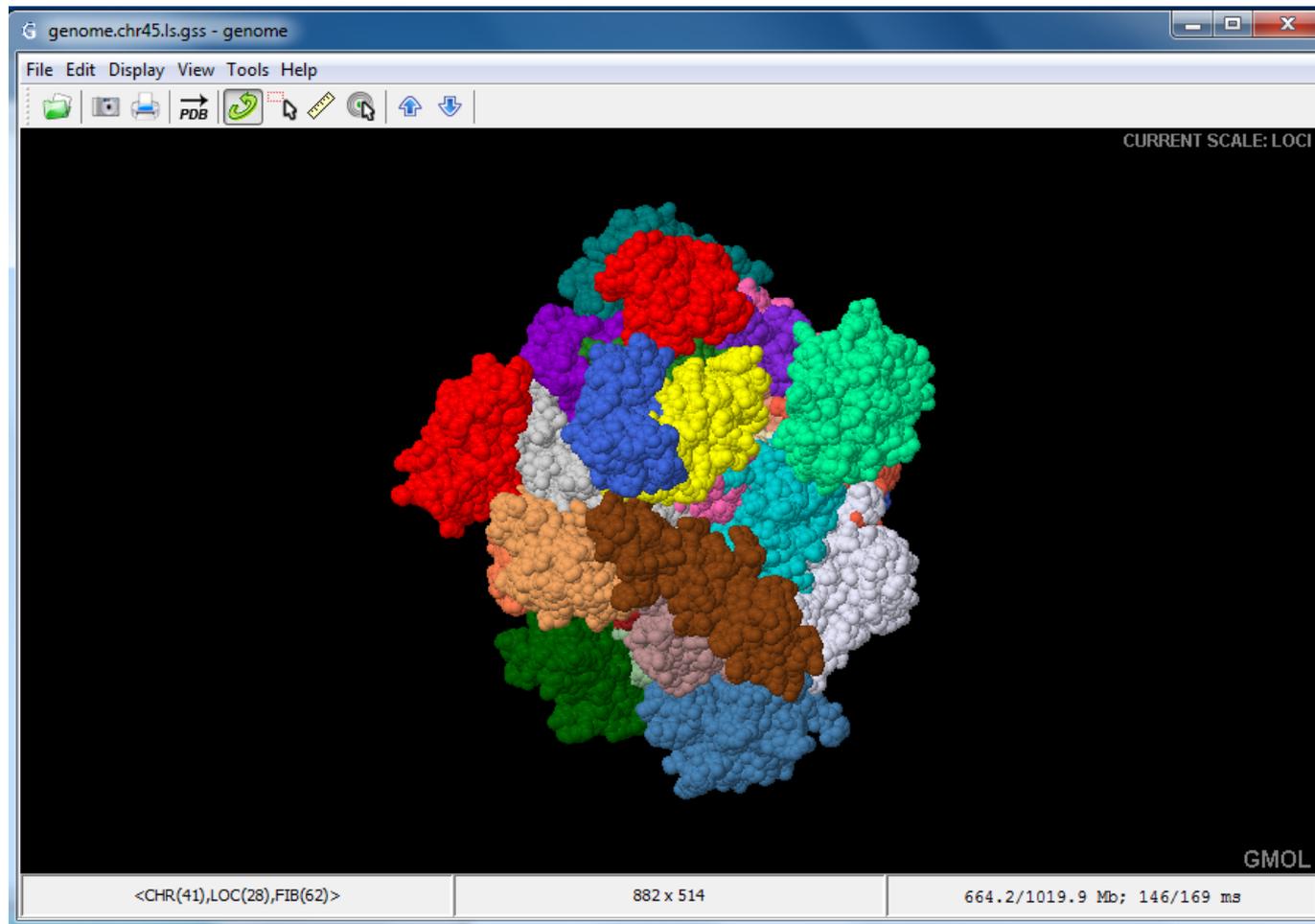
[Find](#)

[Senior PeopleSoft DBA](#)

SA Technologies Inc - St. Louis, MO

[National Account Partner](#)

# GMOL - Multi-Scale Visualization of 3D Genome Structure



# Project 4

- Apply a gradient descent method LorDG / MOGEN to build 3D models for human chromosome 7.
- References: (1) T. Trieu, J. Cheng. **3D Genome Structure Modeling by Lorentzian Objective Function**. *Nucleic Acids Research*, accepted, 2016; (2) Trieu, Cheng. **Large-scale reconstruction of 3D structures of human chromosomes from chromosomal contact data**. *Nucleic Acids Research*, 2014; (3) T. Tuan, J. Cheng. **MOGEN: a tool for reconstructing 3D models of genomes from chromosomal conformation capturing data**. *Bioinformatics*, accepted, doi: 10.1093/bioinformatics/btv754.
- How to assess the models (check the paper)?
- Visualization of models
- Reference models (constructed from somewhat different contact data of Chr. 7: <http://calla.rnet.missouri.edu/mogen/>)
- LorDG at GitHub: <https://github.com/BDM-Lab/LorDG>
- MOGEN at GitHub: <https://github.com/BDM-Lab/MOGEN>

# Data II: Hi-C data of Chromosome 7

HWI-EAS313\_0025:7:78:7863:19096|7 40786044 7 45502707|case 0  
HWI-EAS313\_0025:7:63:14188:6924|7 128562332 7 128562745|case 0  
HWI-EAS313\_0025:7:94:3739:17610|7 154672365 7 154672125|case 0  
HWI-EAS313\_0025:7:78:17921:8167|7 90166635 7 90166267|case 0  
HWI-EAS313\_0025:7:98:5753:2860|7 146851262 7 82792917|case 0  
WI-EAS313\_0025:7:104:3993:13804|7 101472788 7 101472557|case 0  
HWI-EAS313\_0025:7:4:13123:19420|7 63615352 7 62522230|case 0  
HWI-EAS313\_0025:7:54:7610:3364|7 24688078 7 24687498|case 0  
HWI-EAS313\_0025:7:6:8245:1169|7 47788402 7 47788122|case 0

[http://sysbio.rnet.missouri.edu/3dgenome/contact\\_ALL\\_chrom7](http://sysbio.rnet.missouri.edu/3dgenome/contact_ALL_chrom7)

The LorDG /MOGEN software package includes some sample data and examples of how to run the program.

# Timeline

- April 30: discussion of plan
- May 2: presentation of the plan
- May 9: Presentation of your results

# Acknowledgements

- Tuan Trieu, Oluwatosin Oluwadare, Chenfeng He, Sharif Ahmed, Lingfei Xu
- Zheng Wang, Renzhi Cao, Avery Wells
- Charles Caldwell, Kristen Taylor, Aaron Briley

