

3DGenomics

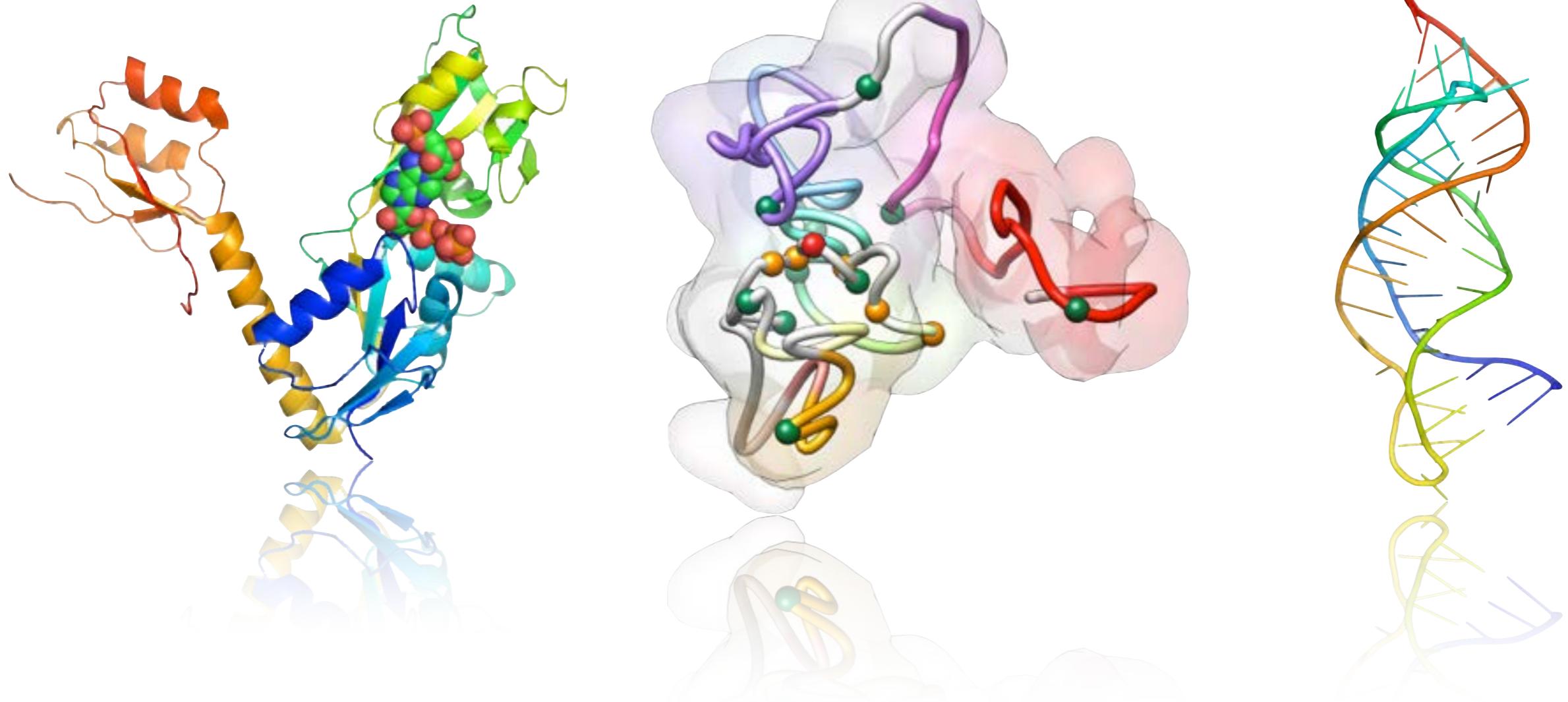
Marc A. Martí-Renom
Genome Biology Group (CNAG)
Structural Genomics Group (CRG)



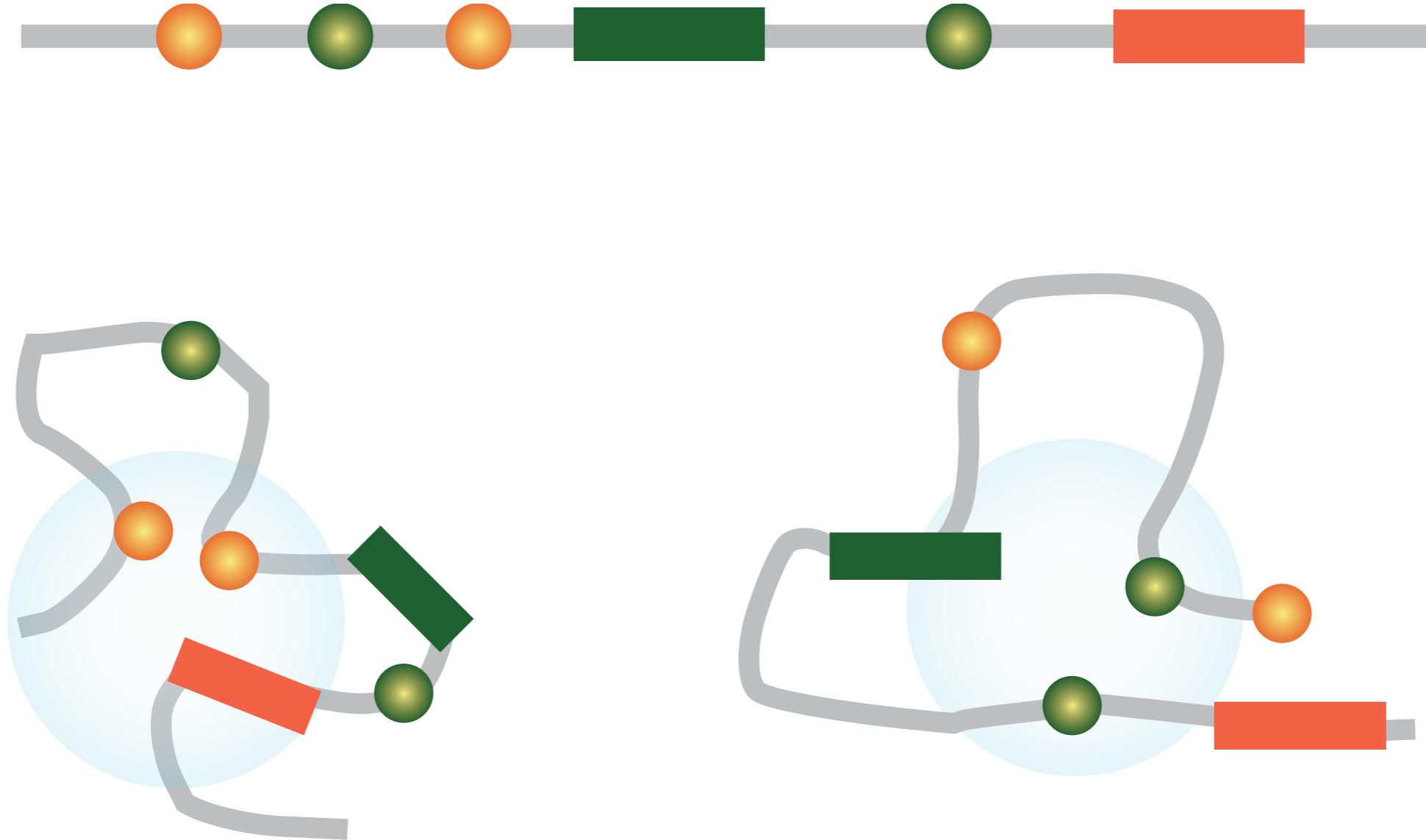


Structural Genomics Group

<http://www.marciuslab.org>

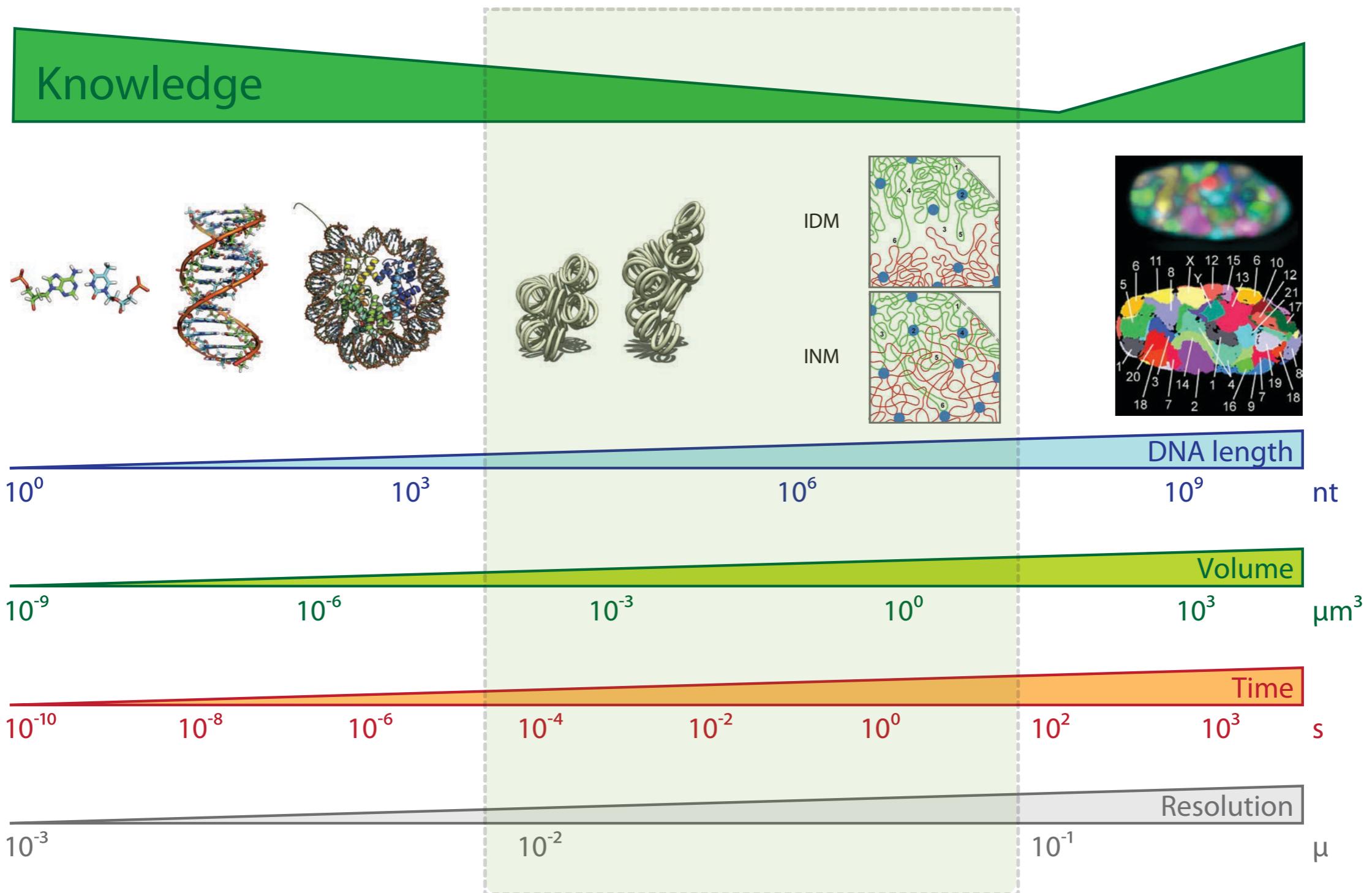


Complex genome organization



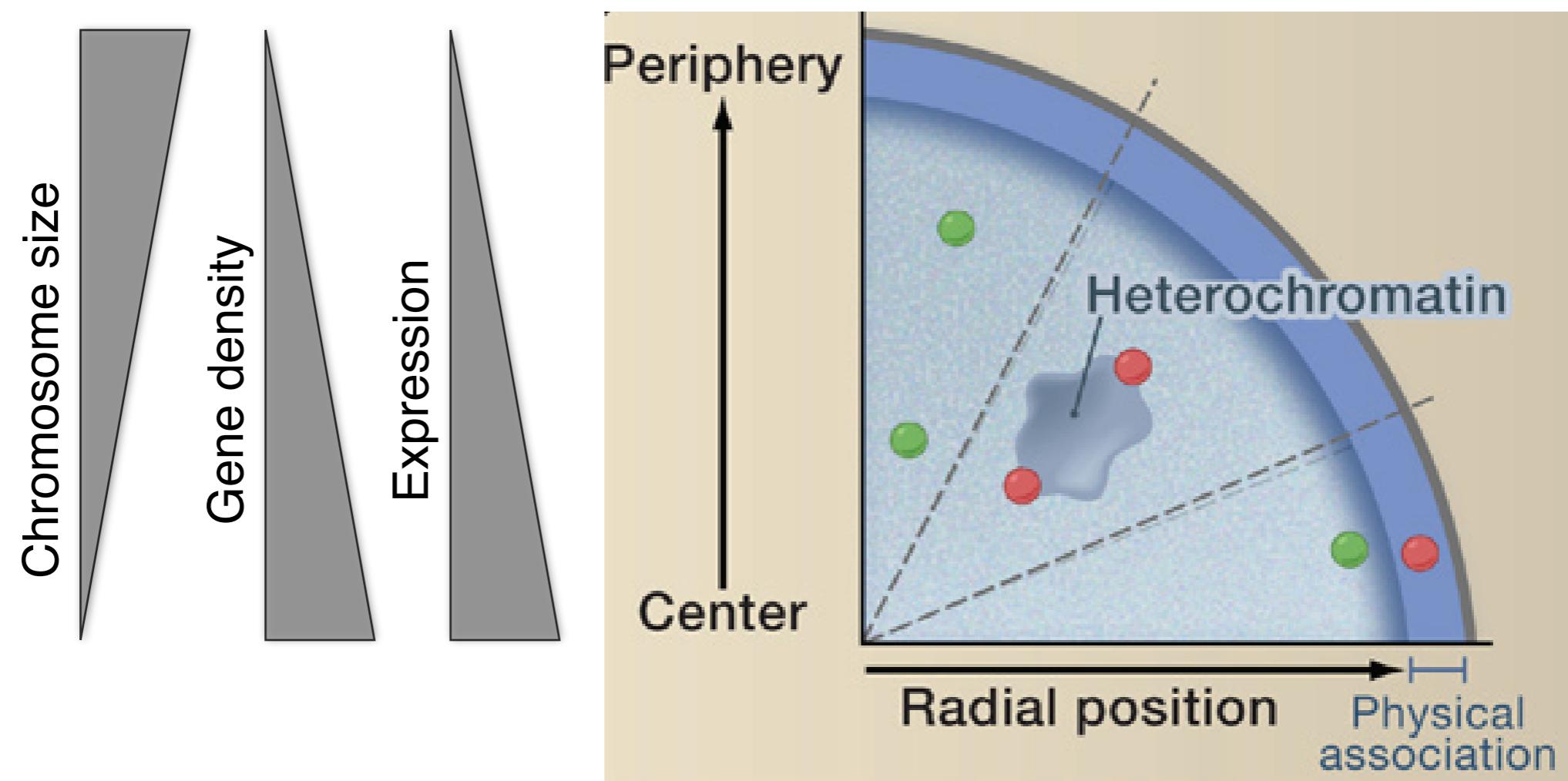
Resolution Gap

Marti-Renom, M. A. & Mirny, L. A. PLoS Comput Biol 7, e1002125 (2011)



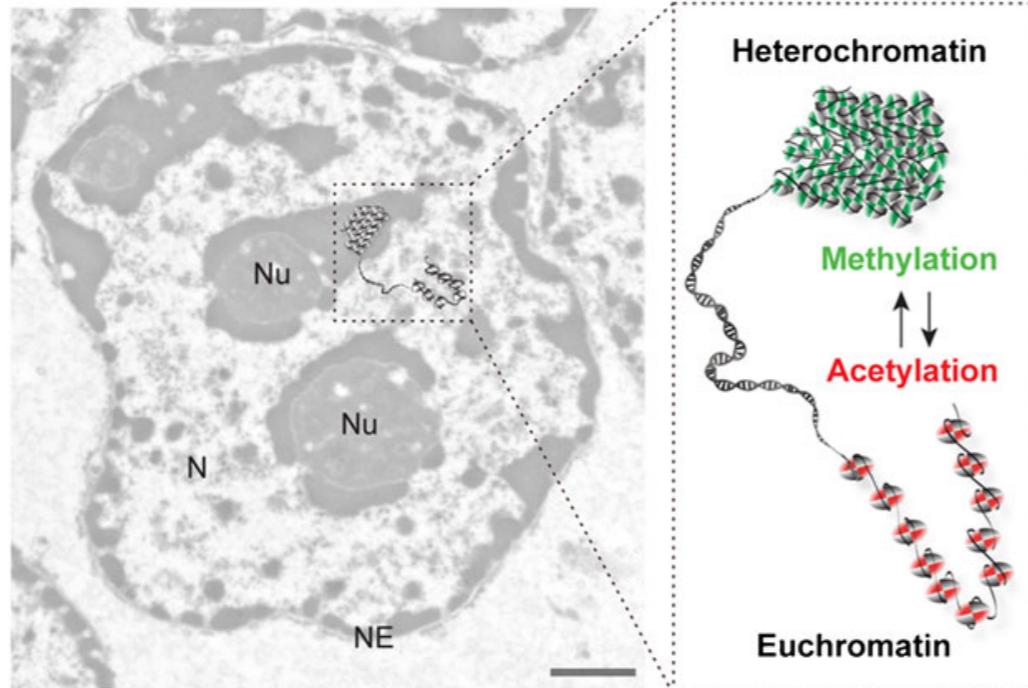
Level I: Radial genome organization

Takizawa, T., Meaburn, K. J. & Misteli, T. The meaning of gene positioning. *Cell* 135, 9–13 (2008).

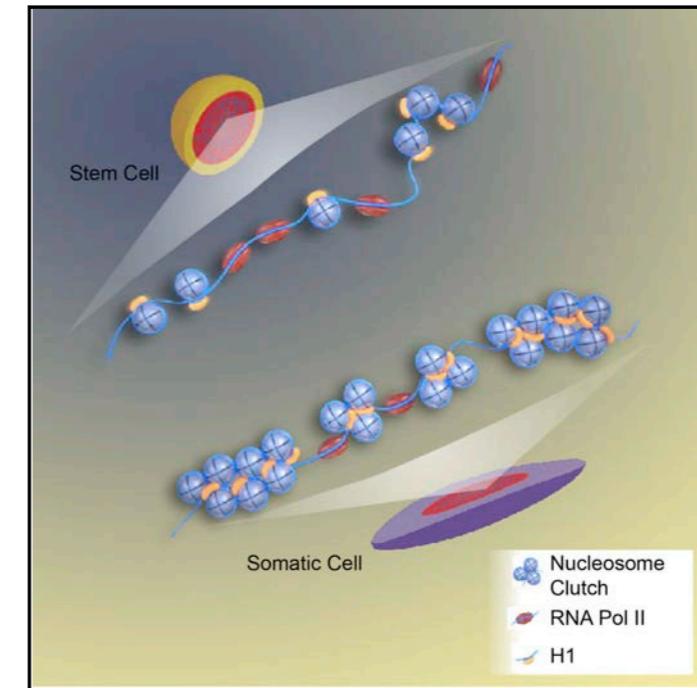


Level II: Euchromatin vs heterochromatin

Electron microscopy



Nanoscropy



Euchromatin:

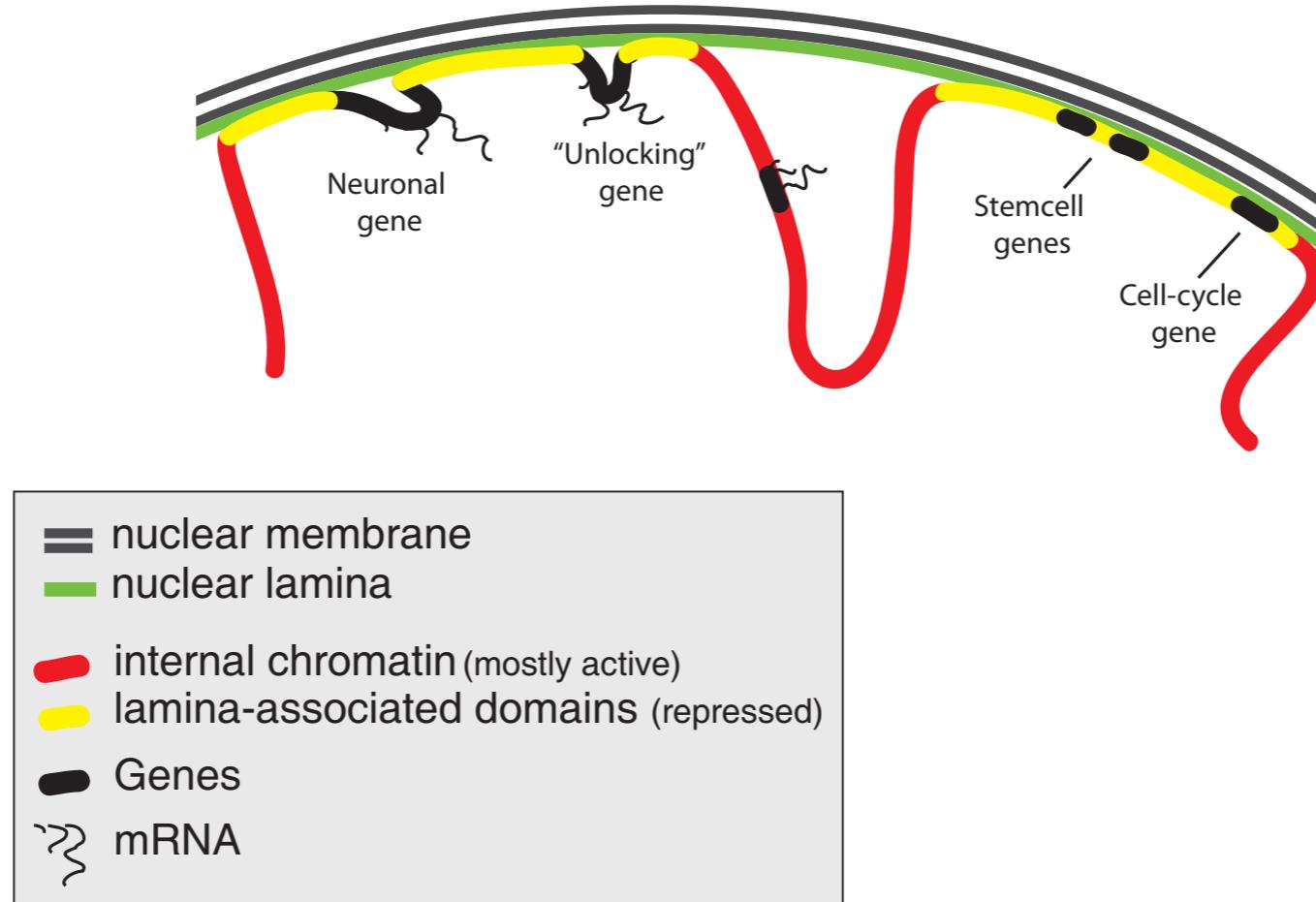
chromatin that is located away from the nuclear lamina, is generally less densely packed, and contains actively transcribed genes

Heterochromatin:

chromatin that is near the nuclear lamina, tightly condensed, and transcriptionally silent

Adapted from Cell, 160(6), 1145-1158. 2015

Level III: Lamina-genome interactions

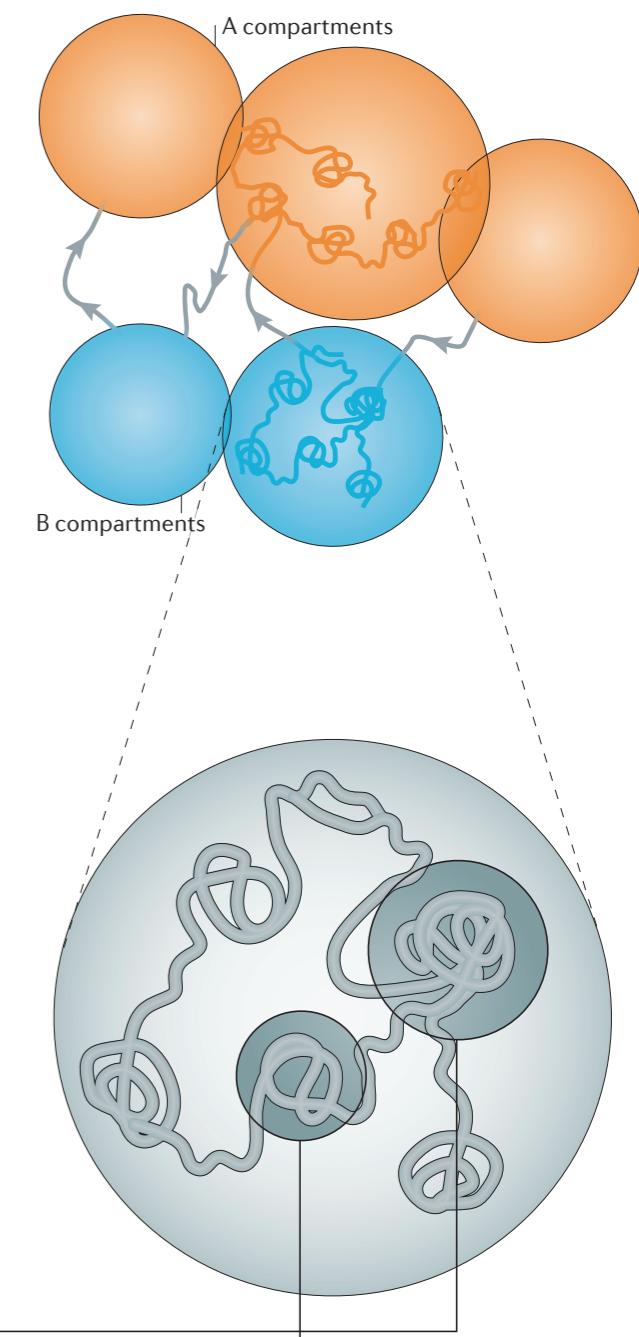
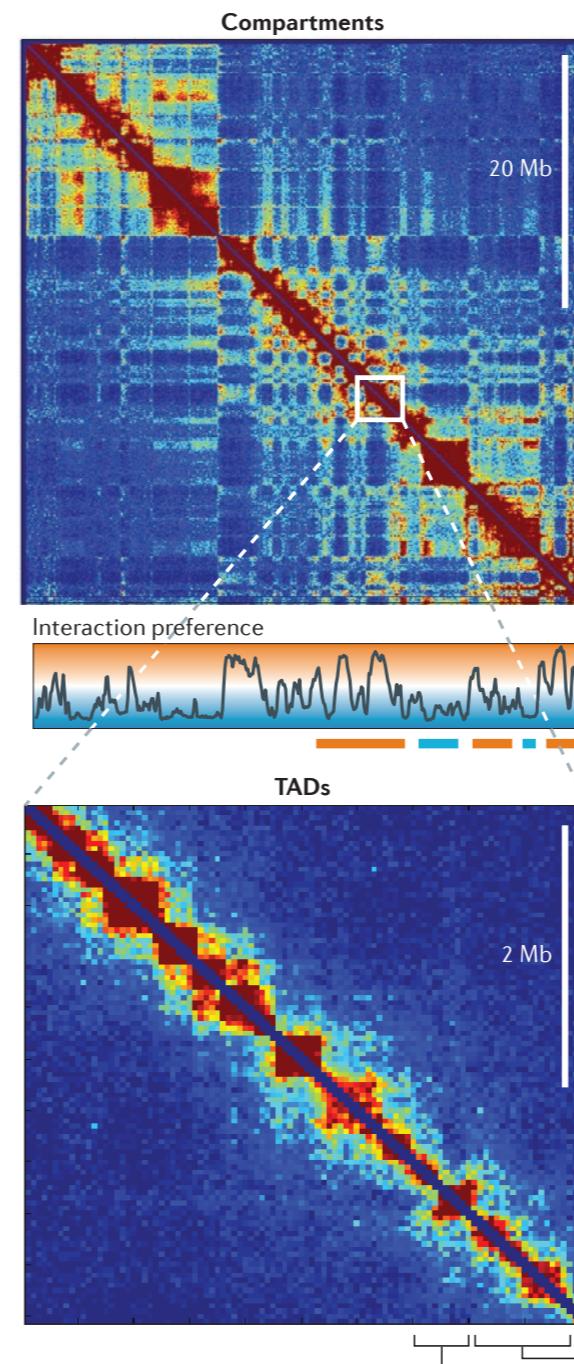
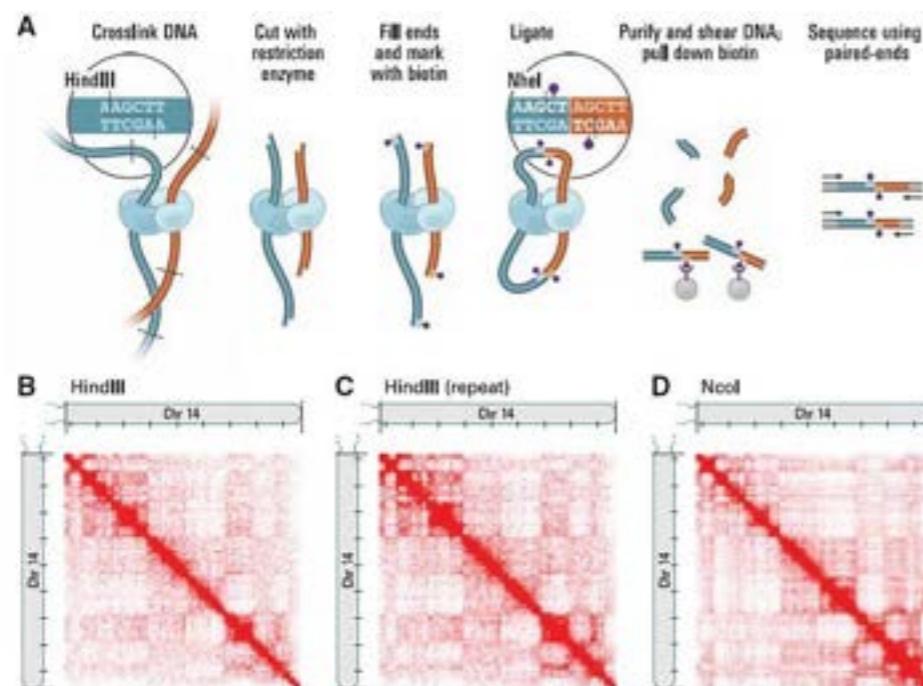


Most genes in Lamina Associated Domains are transcriptionally silent, suggesting that **lamina-genome interactions** are widely involved in the control of **gene expression**

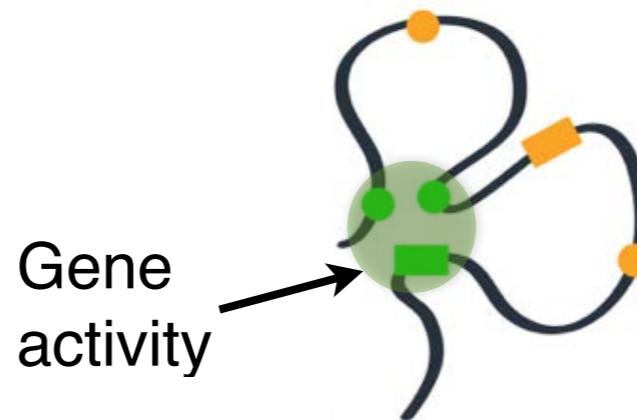
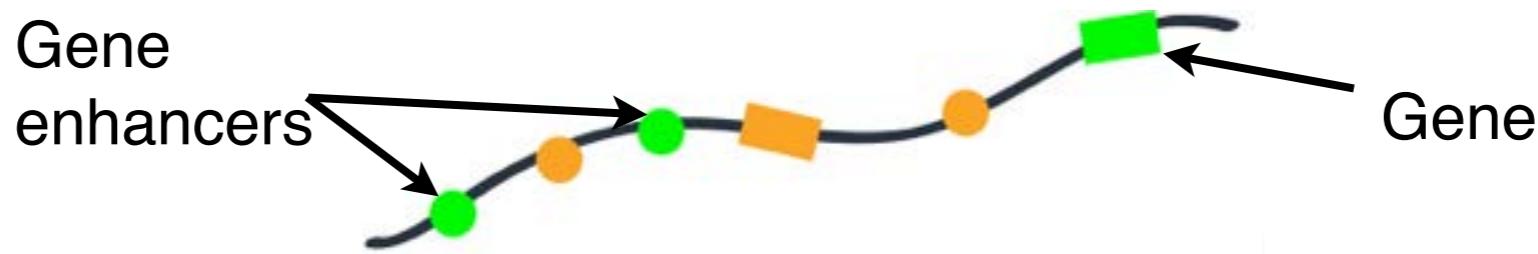
Adapted from Molecular Cell 38, 603-613, 2010

Level IV: Higher-order organization

Dekker, J., Marti-Renom, M. A. & Mirny, L. A. Exploring the three-dimensional organization of genomes: interpreting chromatin interaction data. Nat Rev Genet 14, 390–403 (2013).



Level V: Chromatin loops



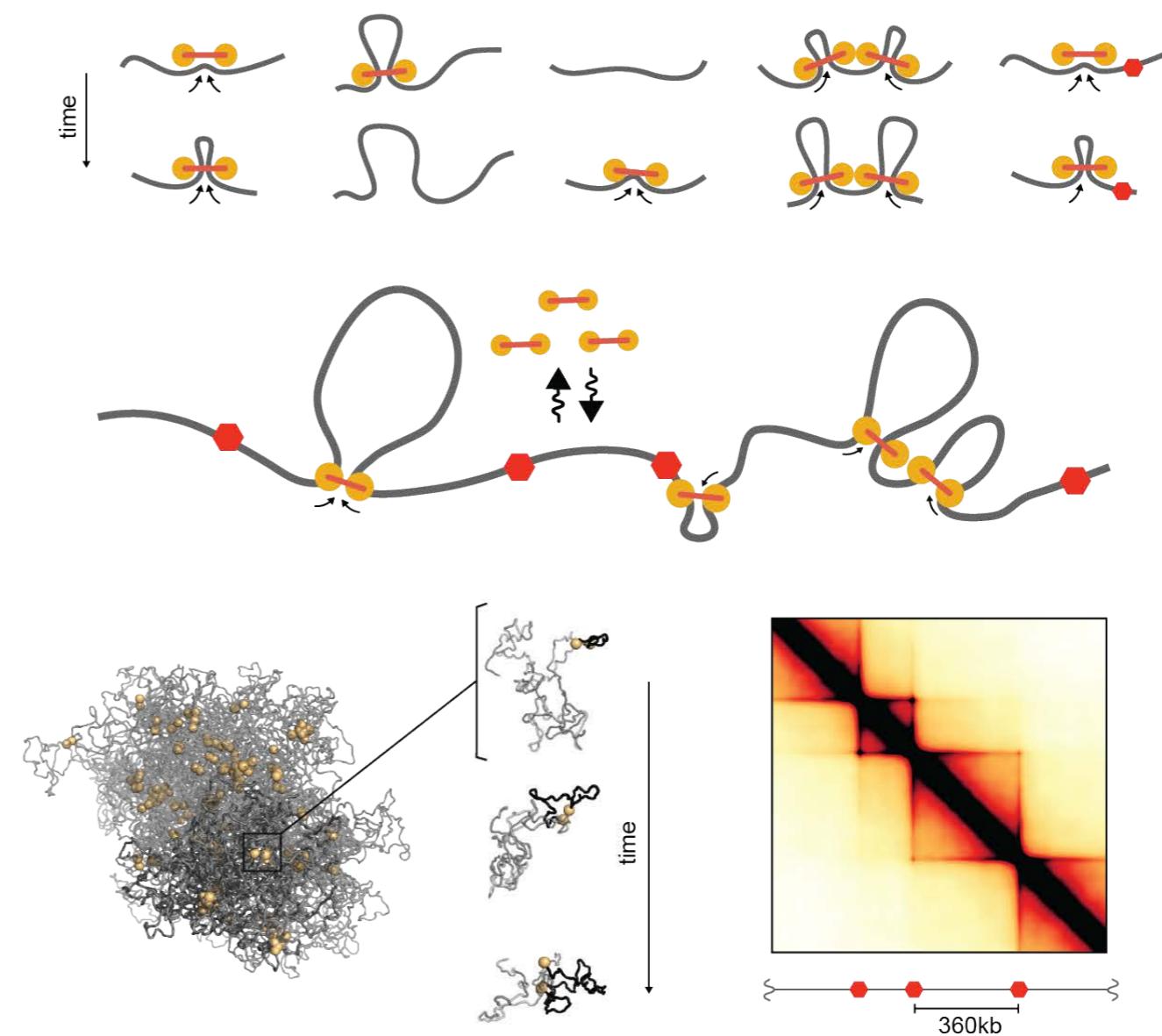
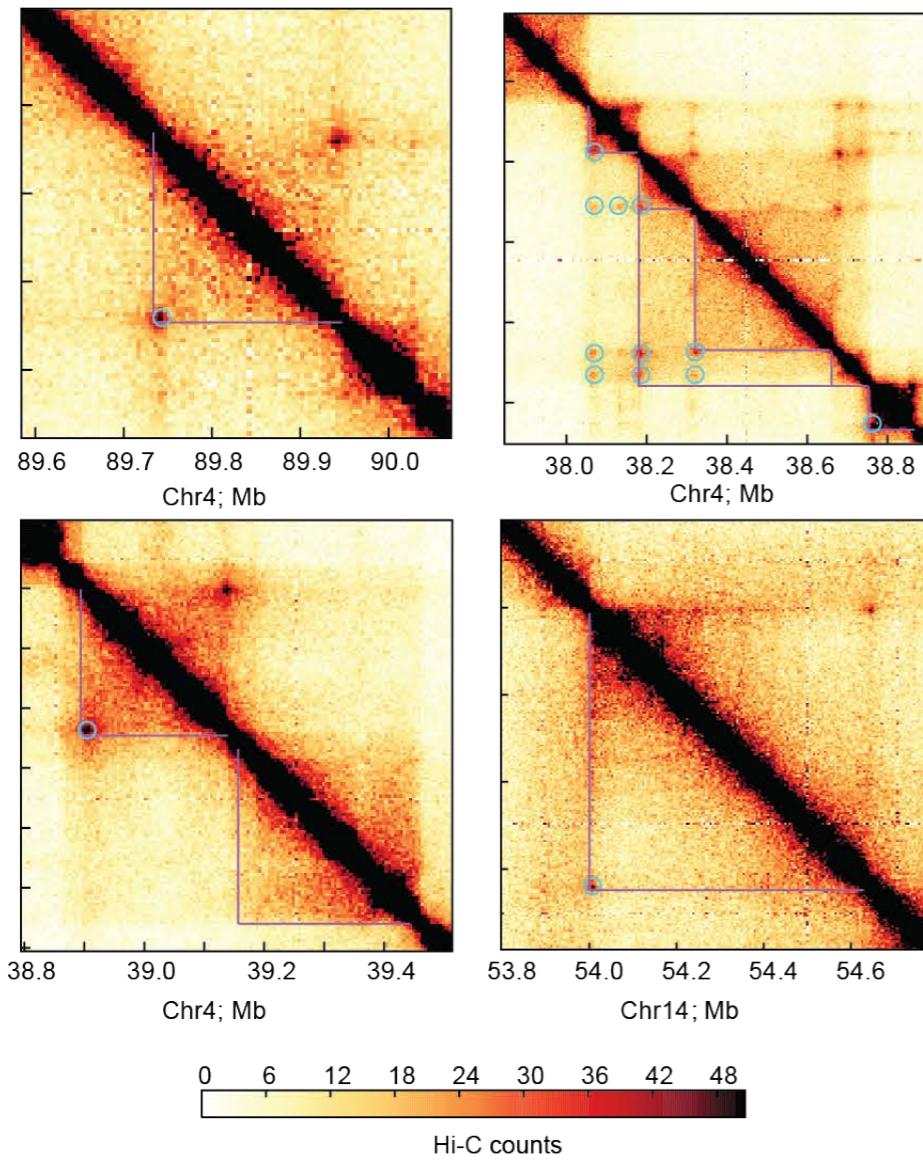
Loops bring distal genomic regions in close proximity to one another

This in turn can have profound effects on gene transcription

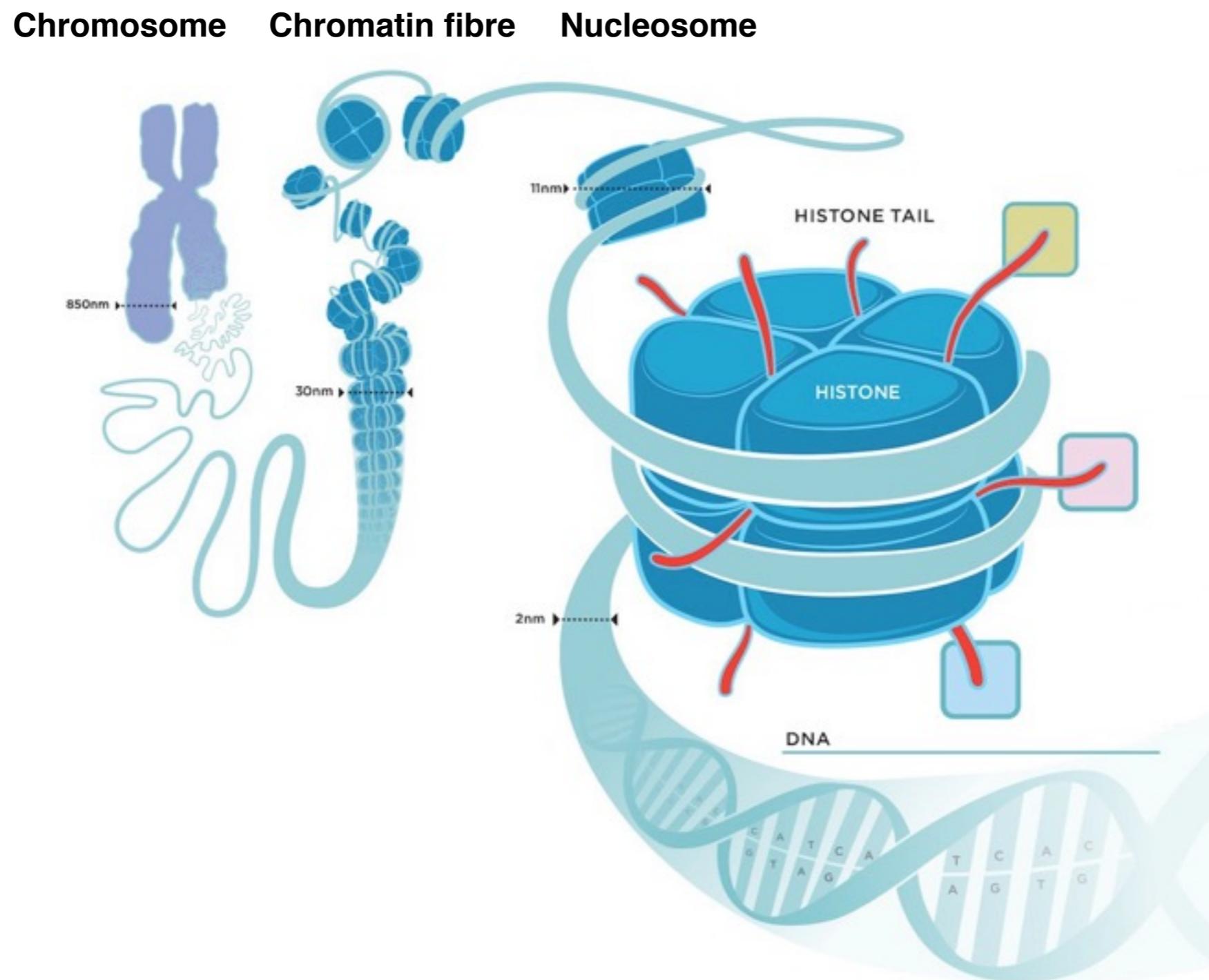
Enhancers can be thousands of kilobases away from their target genes in any direction (or even on a separate chromosome)

Level V: Loop-extrusion as a driving force

Fudenberg, G., Imakaev, M., Lu, C., Goloborodko, A., Abdennur, N., & Mirny, L. A. (2015).
Formation of Chromosomal Domains by Loop Extrusion. *bioRxiv*.



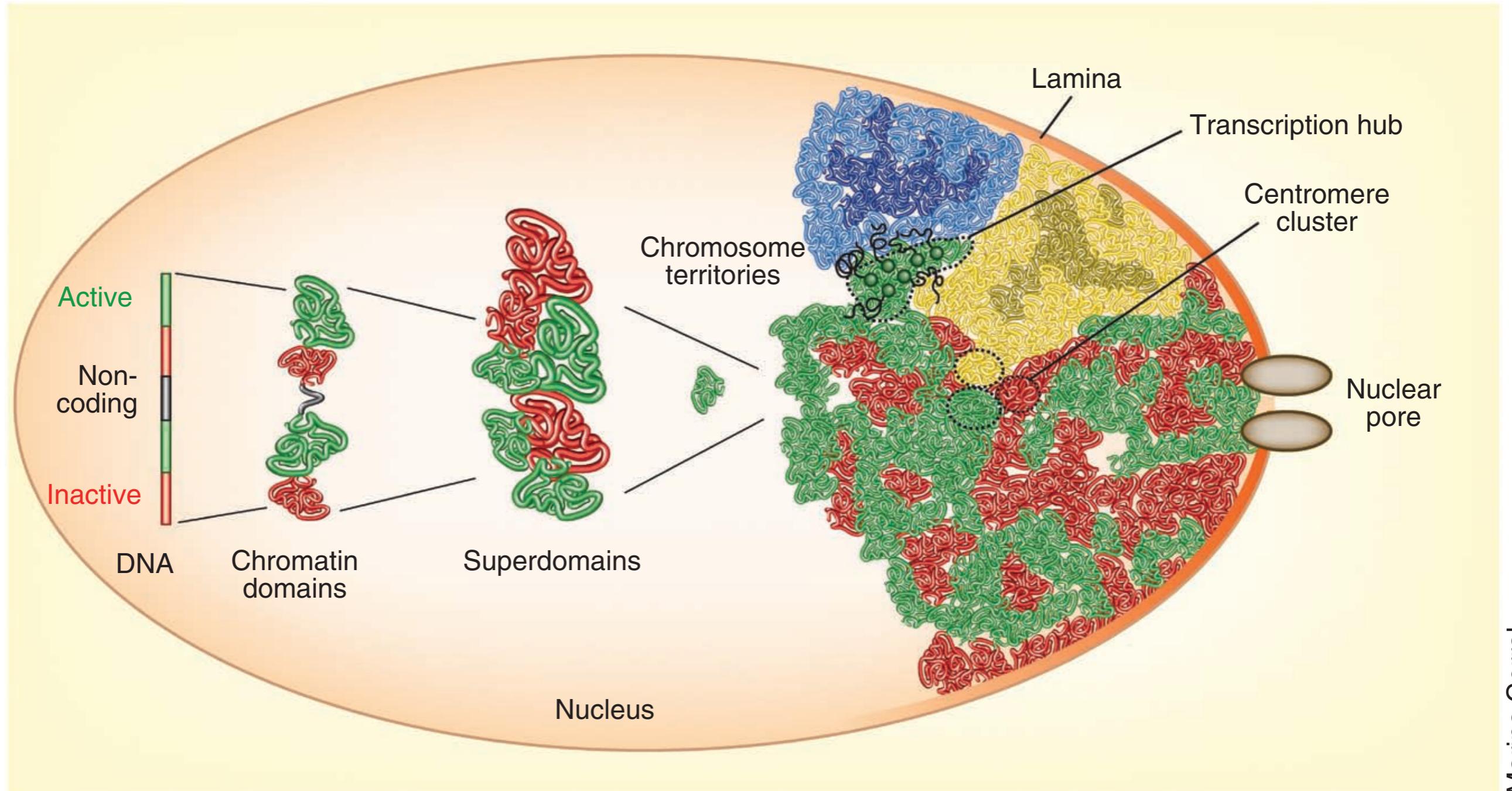
Level VI: Nucleosome



Adapted from Richard E. Ballermann, 2012

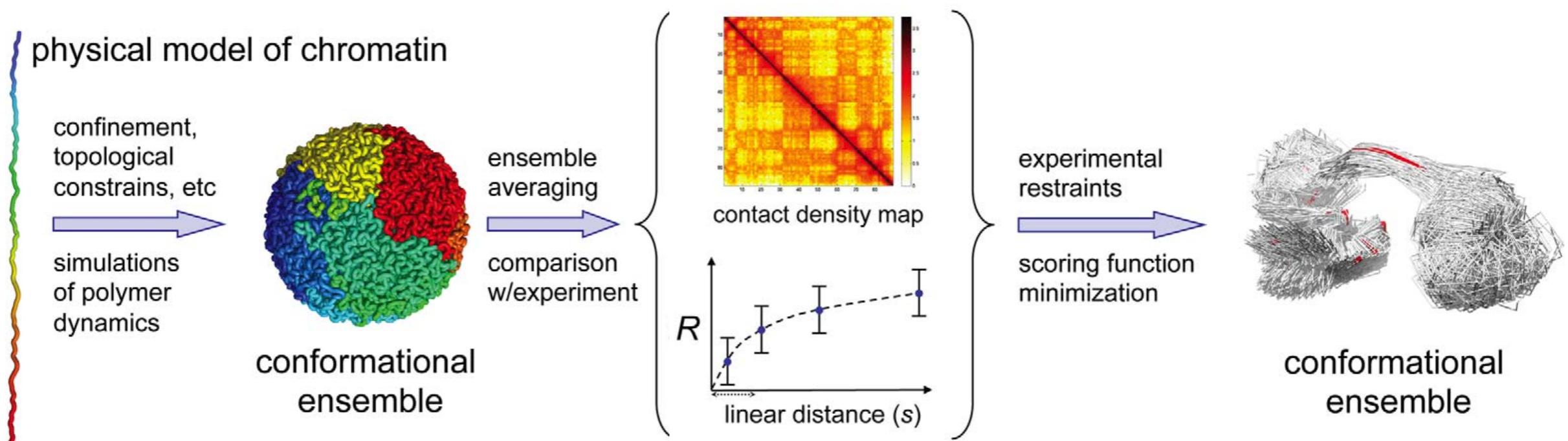
Complex genome organization

Cavalli, G. & Misteli, T. Functional implications of genome topology. *Nat Struct Mol Biol* 20, 290–299 (2013).

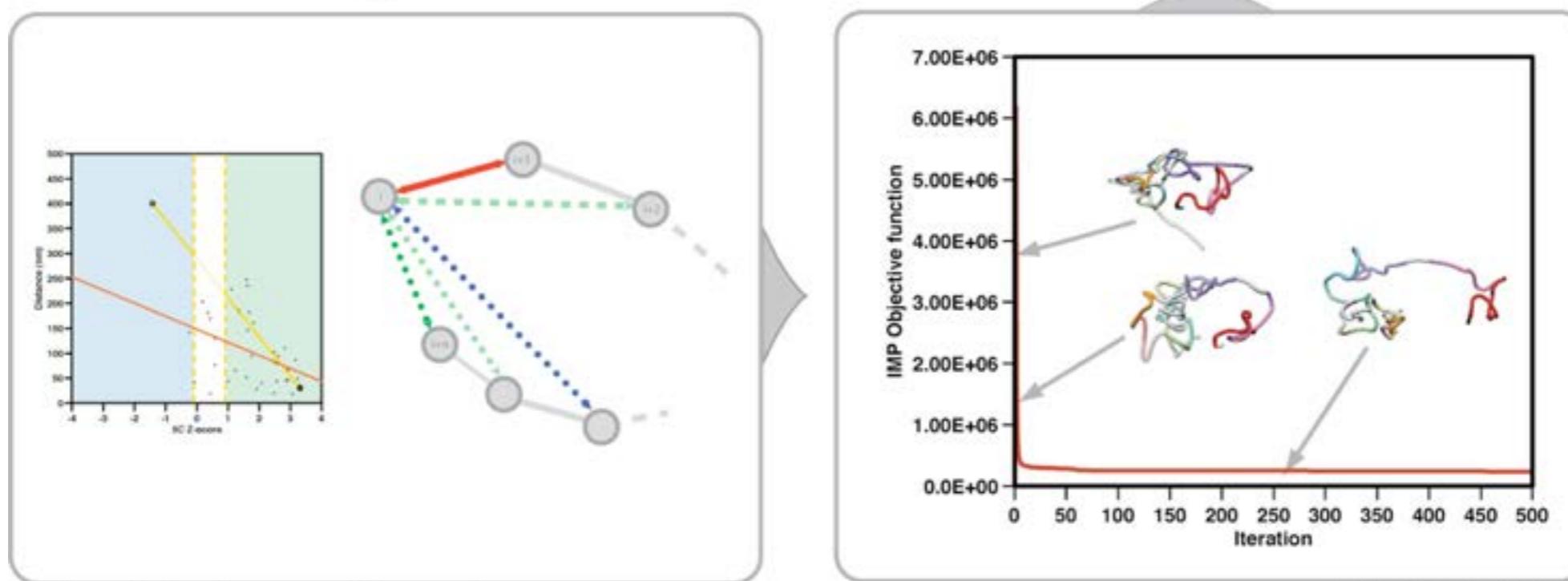
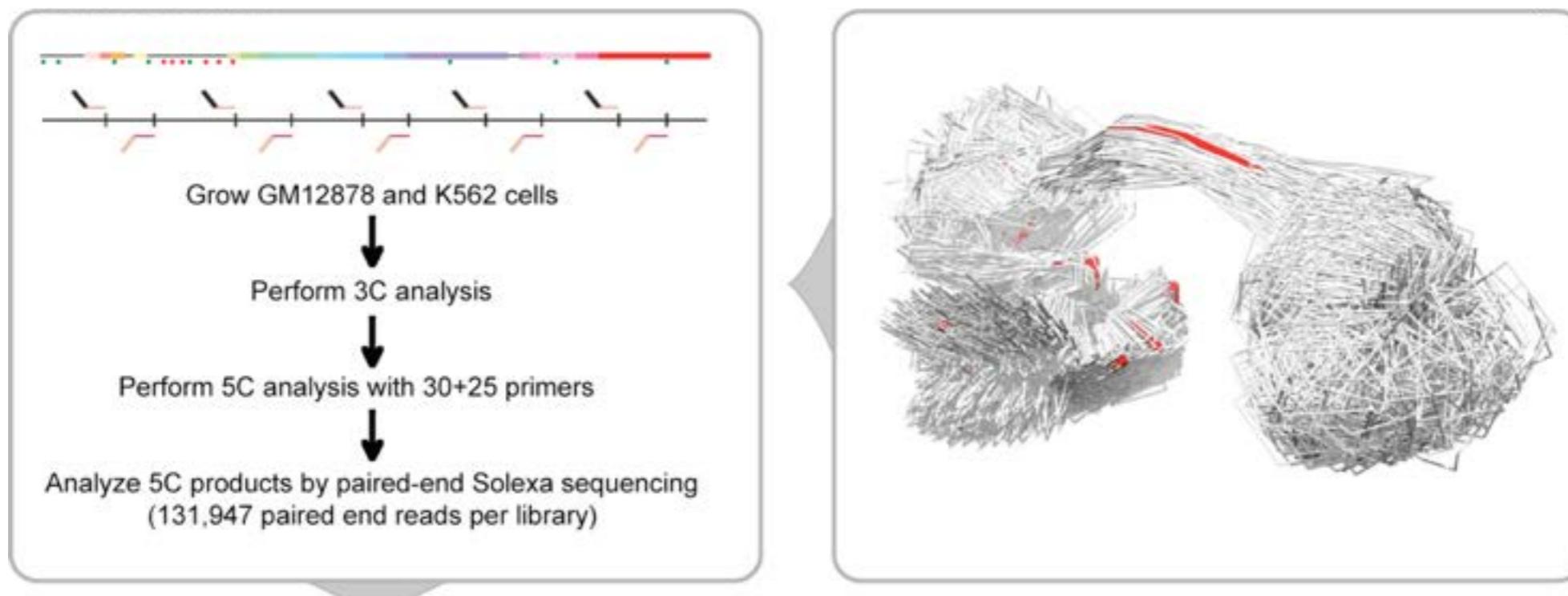


Modeling Genomes

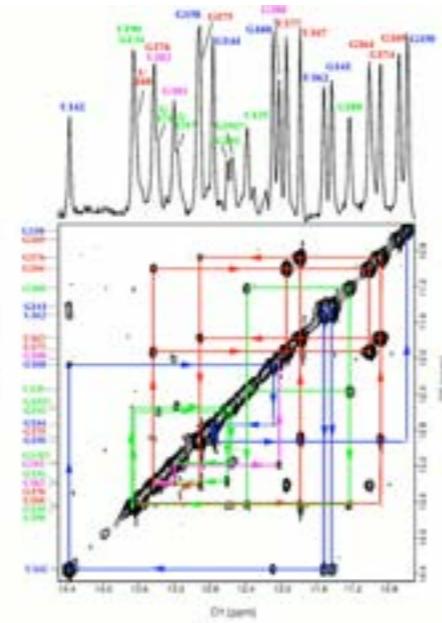
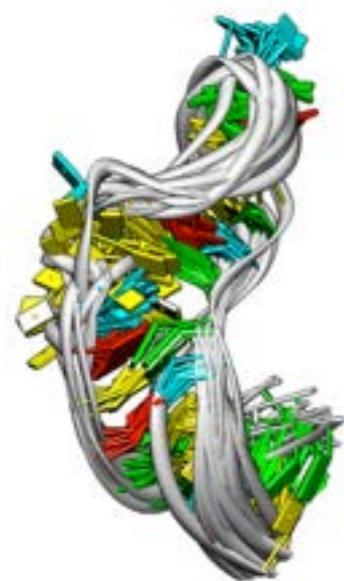
Marti-Renom, M. A. & Mirny, L. A. PLoS Comput Biol 7, e1002125 (2011)



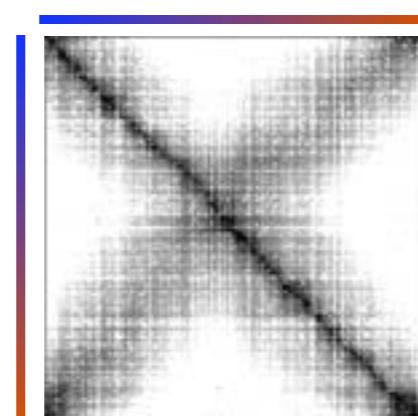
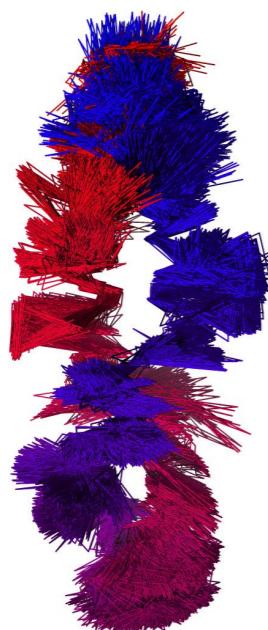
Experiments



Computation

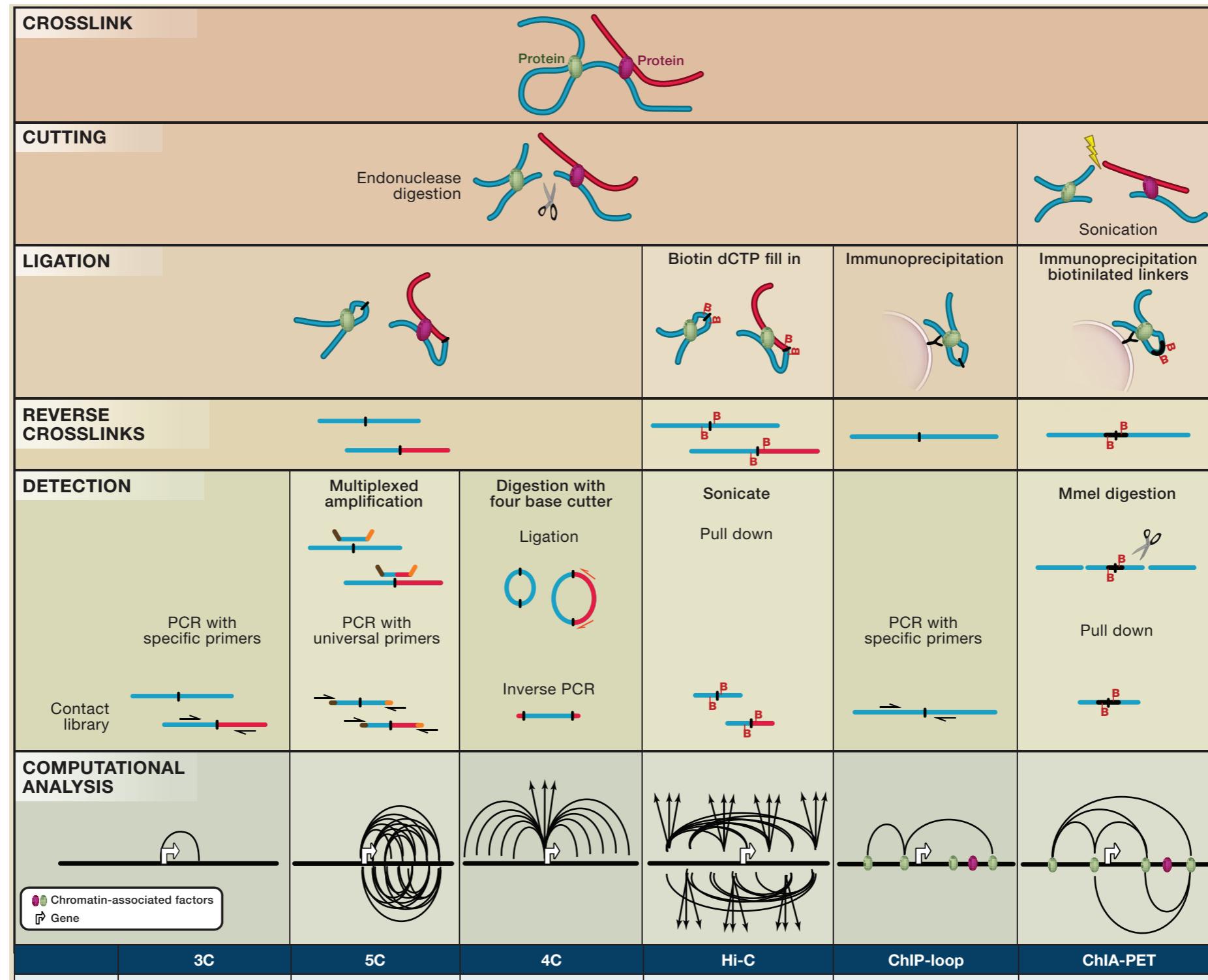


Biomolecular structure determination
2D-NOESY data



Chromosome structure determination
5C data

Chromosome Conformation Capture



Hakim, O., & Misteli, T. (2012). SnapShot: Chromosome Confirmation Capture. *Cell*, 148(5), 1068–1068.e2.

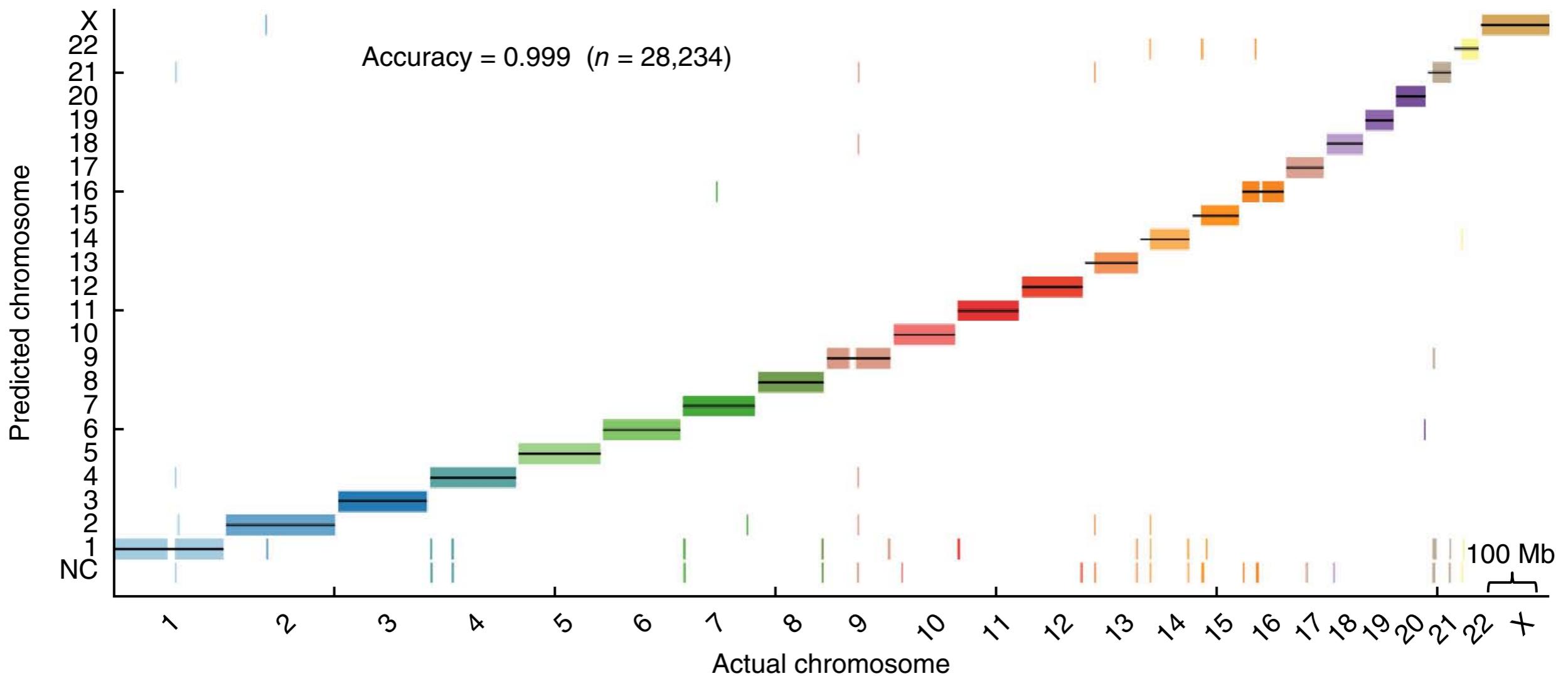
Chromosome Conformation Capture

	3C	5C	4C	Hi-C	ChIP-loop	ChIA-PET
Principle	Contacts between two defined regions ^{3,17}	All against all ^{4,18}	All contacts with a point of interest ¹⁴	All against all ¹⁰	Contacts between two defined regions associated with a given protein ⁸	All contacts associated with a given protein ⁶
Coverage	Commonly < 1Mb	Commonly < 1Mb	Genome-wide	Genome-wide	Commonly < 1Mb	Genome-wide
Detection	Locus-specific PCR	HT-sequencing	HT-sequencing	HT-sequencing	Locus-specific qPCR	HT-sequencing
Limitations	Low throughput and coverage	Limited coverage	Limited to one viewpoint		Rely on one chromatin-associated factor, disregarding other contacts	
Examples	Determine interaction between a known promoter and enhancer	Determine comprehensively higher-order chromosome structure in a defined region	All genes and genomic elements associated with a known LCR	All intra- and interchromosomal associations	Determine the role of specific transcription factors in the interaction between a known promoter and enhancer	Map chromatin interaction network of a known transcription factor
Derivatives	PCR with TaqMan probes ⁷ or melting curve analysis ¹		Circular chromosome conformation capture ²⁰ , open-ended chromosome conformation capture ¹⁹ , inverse 3C ¹² , associated chromosome trap (ACT) ¹¹ , affinity enrichment of bait-ligated junctions ²	Yeast ^{5,15} , tethered conformation capture ⁹		ChIA-PET combined 3C-ChIP-cloning (6C) ¹⁶ , enhanced 4C (e4C) ¹³

Hakim, O., & Misteli, T. (2012). SnapShot: Chromosome Confirmation Capture. *Cell*, 148(5), 1068–1068.e2.

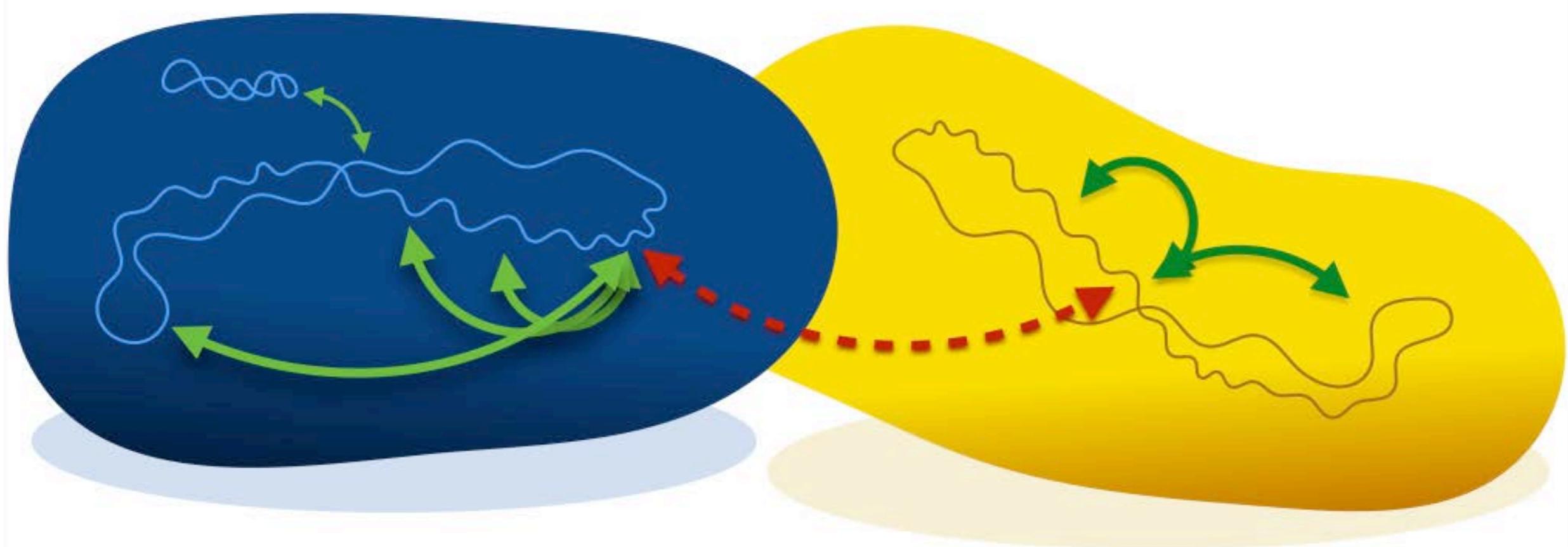
... and one more thing

Chromosome Conformation Capture for de-novo assembly



Kaplan, N., & Dekker, J. (2013). High-throughput genome scaffolding from *in vivo* DNA interaction frequency. *Nature Biotechnology*, 31(12), 1143–1147.

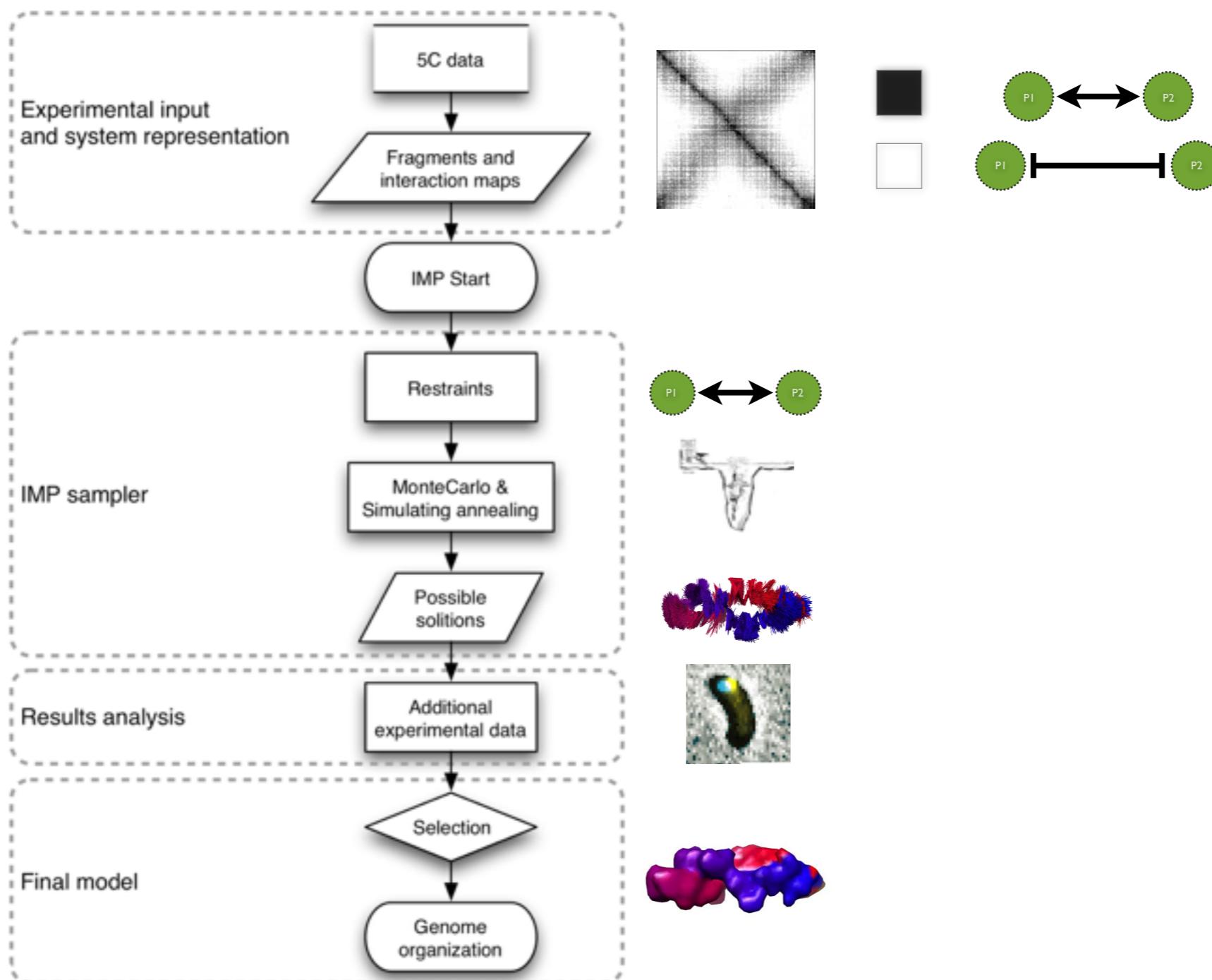
Chromosome Conformation Capture for meta genomics



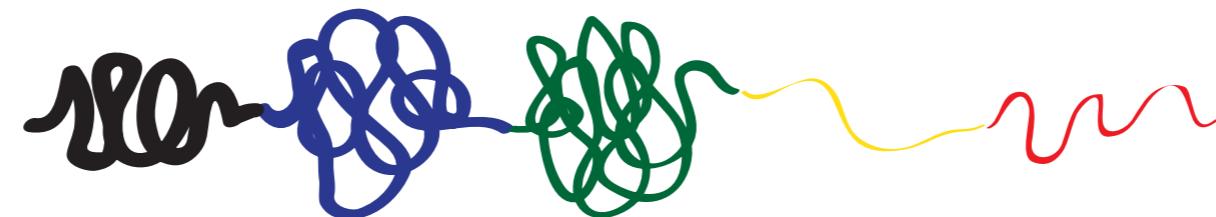
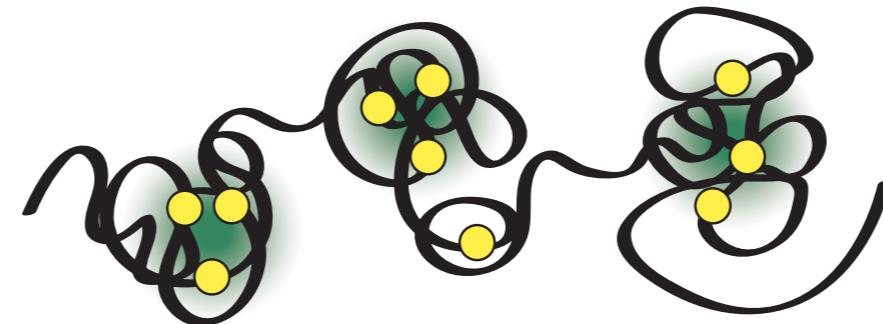
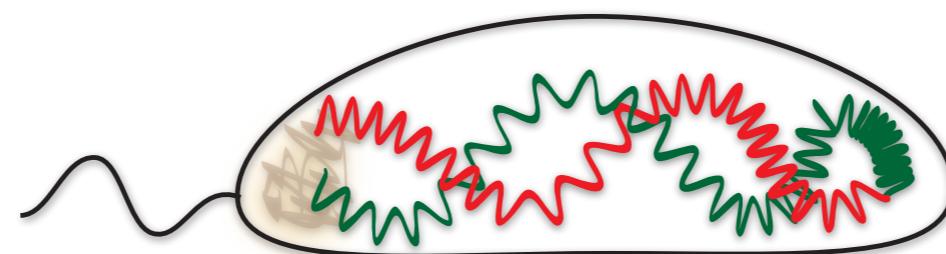
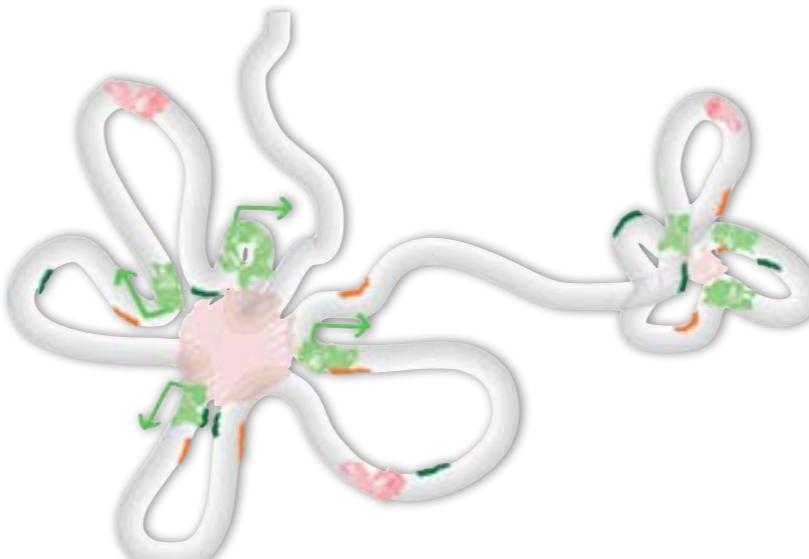
Beitel, C. W., Froenicke, L., Lang, J. M., Korf, I. F., Michelmore, R. W., Eisen, J. A., & Darling, A. E. (2014). Strain- and plasmid-level deconvolution of a synthetic metagenome by sequencing proximity ligation products. doi:10.7287/peerj.preprints.260v1

Modeling 3D Genomes

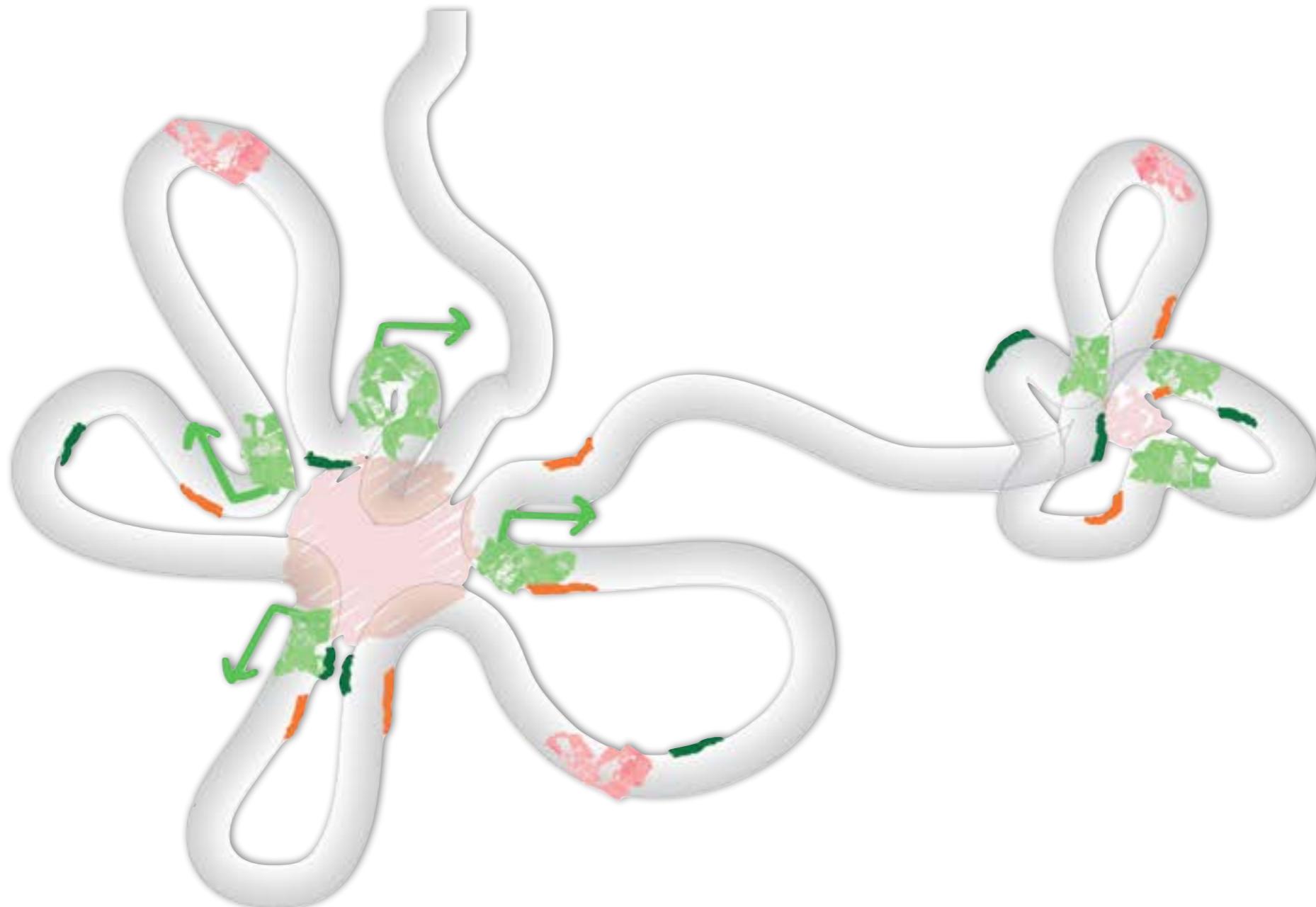
Baù, D. & Martí-Renom, M. A. Methods 58, 300–306 (2012).



Examples...



Human α -globin domain



Human α -globin domain

ENm008 genomic structure and environment

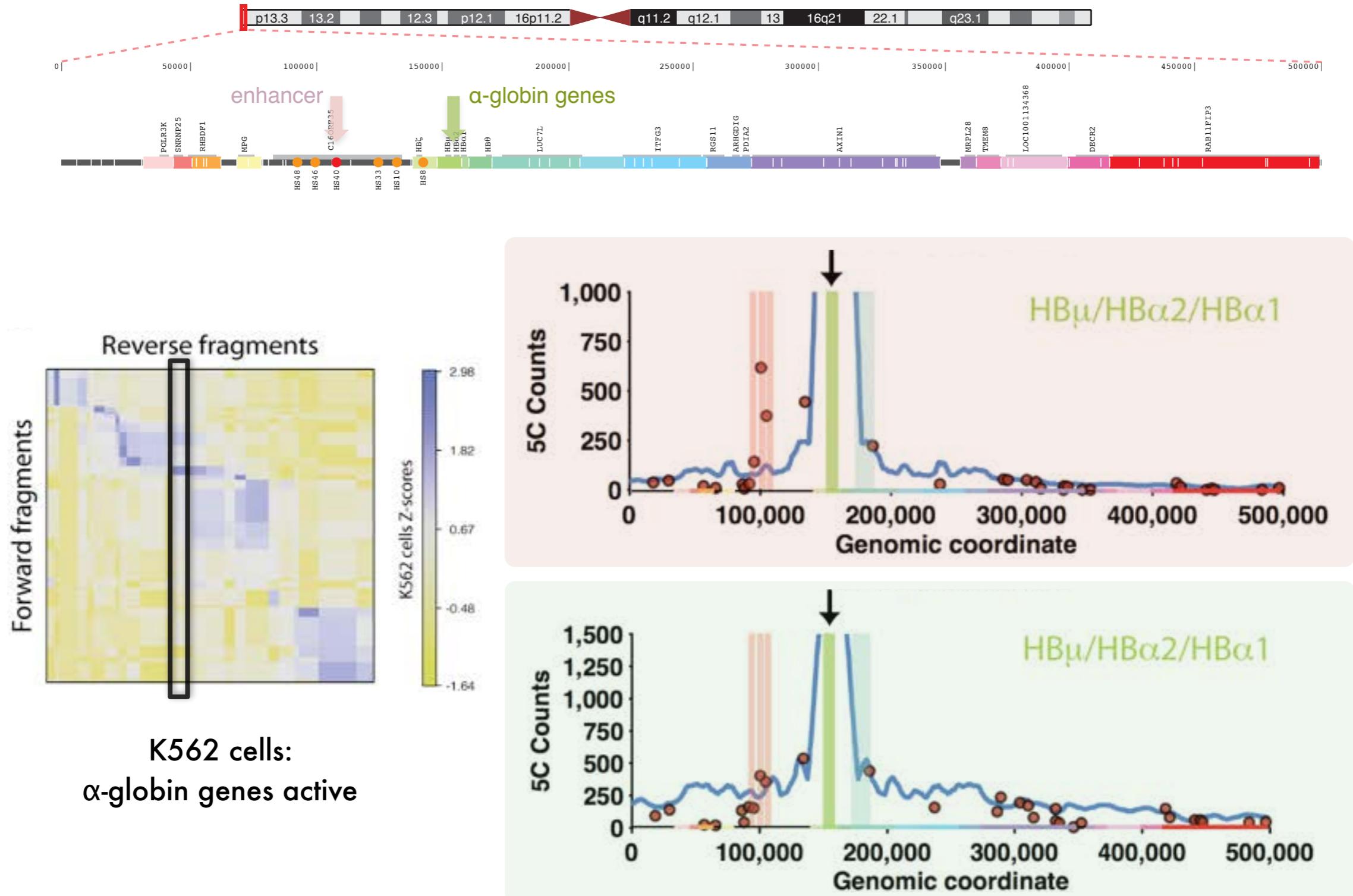


The ENCODE data for ENm008 region was obtained from the UCSC Genome Browser tracks for: RefSeq annotated genes, Affymetrix/CSHL expression data (Gingeras Group at Cold Spring Harbor), Duke/NHGRI DNaseI Hypersensitivity data (Crawford Group at Duke University), and Histone Modifications by Broad Institute ChIP-seq (Bernstein Group at Broad Institute of Harvard and MIT).

ENCODE Consortium. Nature (2007) vol. 447 (7146) pp. 799-816

Human α -globin domain

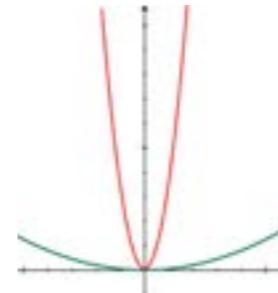
ENm008 genomic structure and environment



Representation

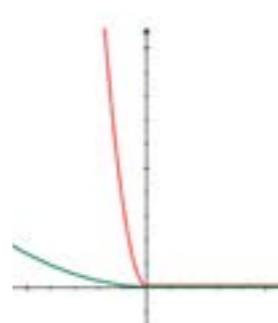
Harmonic

$$H_{i,j} = k(d_{i,j} - d_{i,j}^0)^2$$



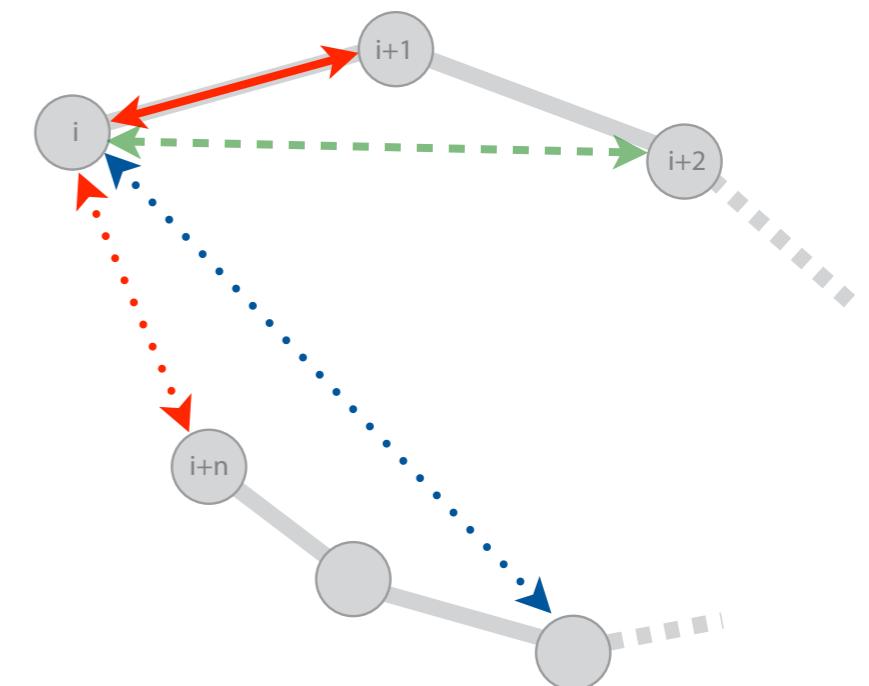
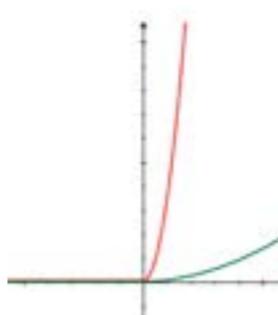
Harmonic Lower Bound

$$\begin{cases} \text{if } d_{i,j} \leq d_{i,j}^0; & lbH_{i,j} = k(d_{i,j} - d_{i,j}^0)^2 \\ \text{if } d_{i,j} > d_{i,j}^0; & lbH_{i,j} = 0 \end{cases}$$

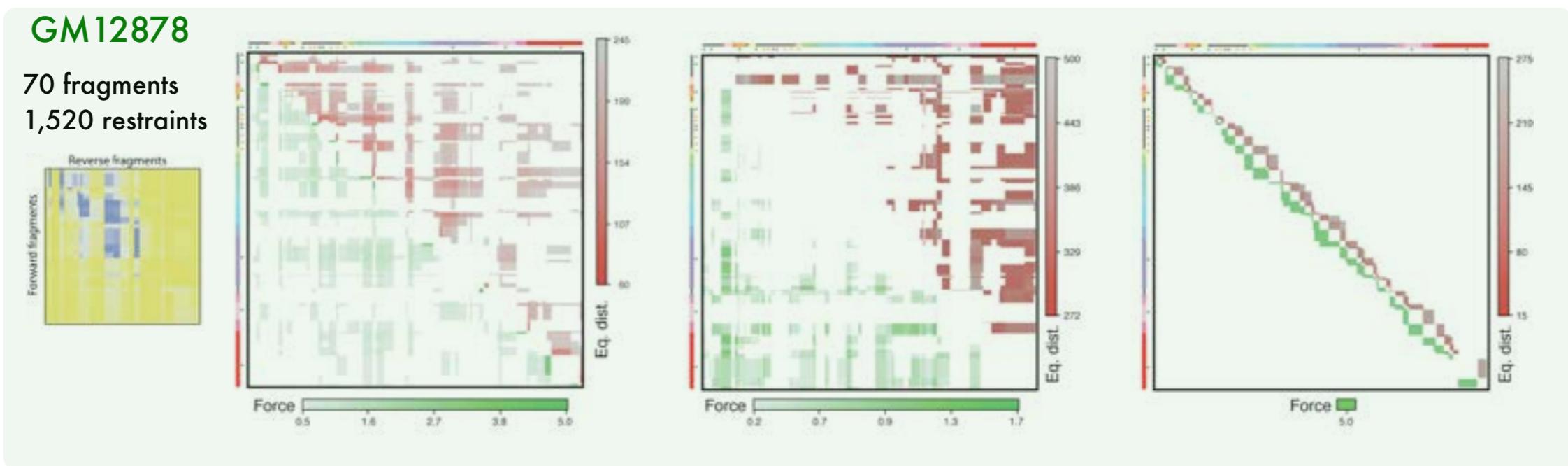


Harmonic Upper Bound

$$\begin{cases} \text{if } d_{i,j} \geq d_{i,j}^0; & ubH_{i,j} = k(d_{i,j} - d_{i,j}^0)^2 \\ \text{if } d_{i,j} < d_{i,j}^0; & ubH_{i,j} = 0 \end{cases}$$



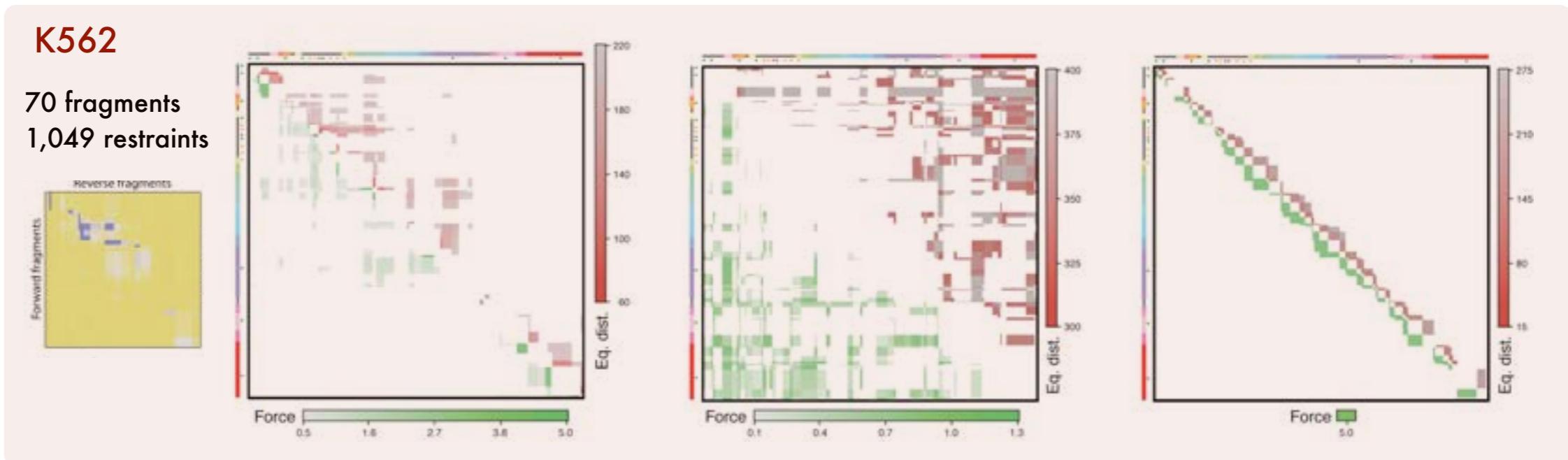
Scoring



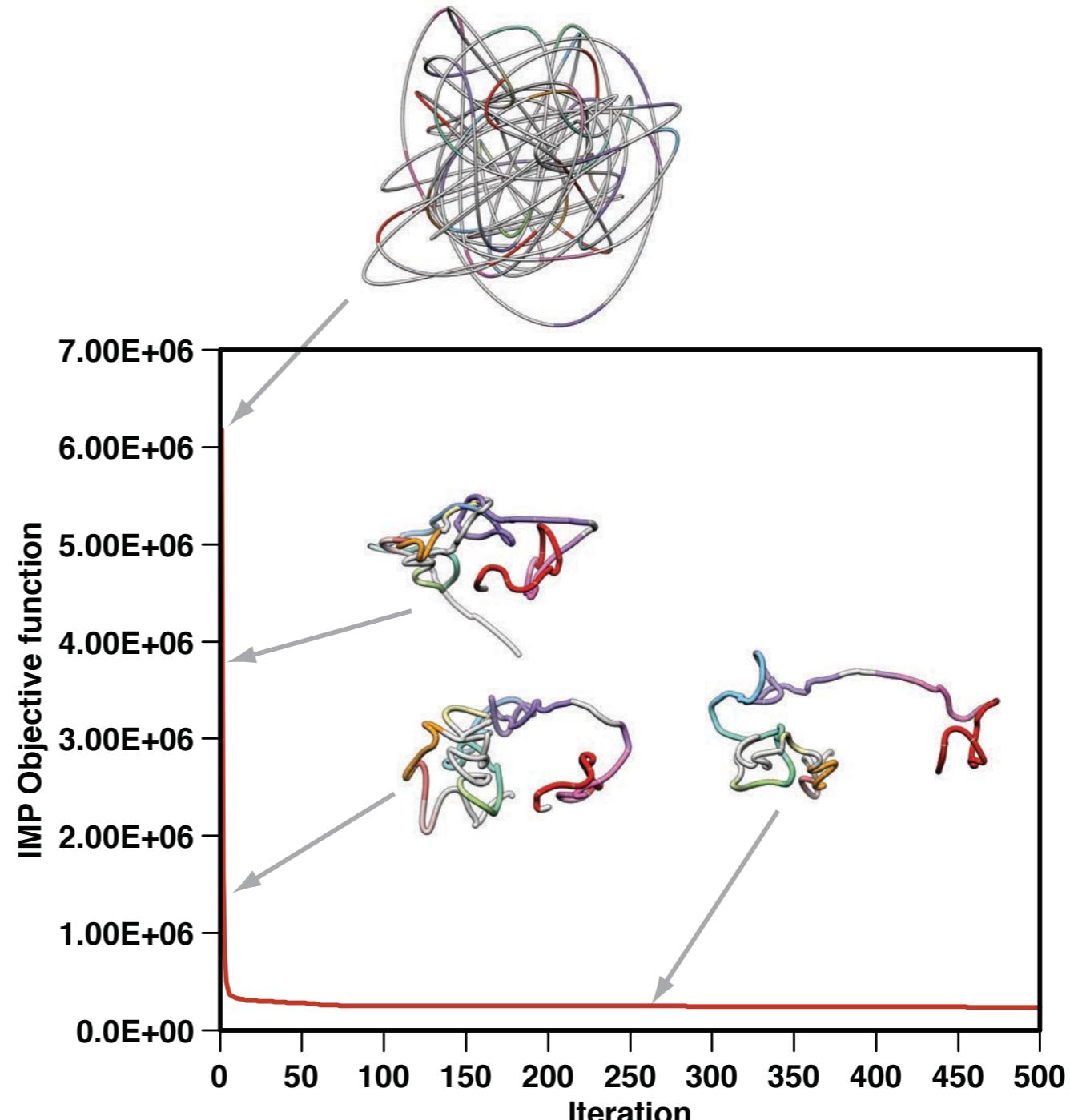
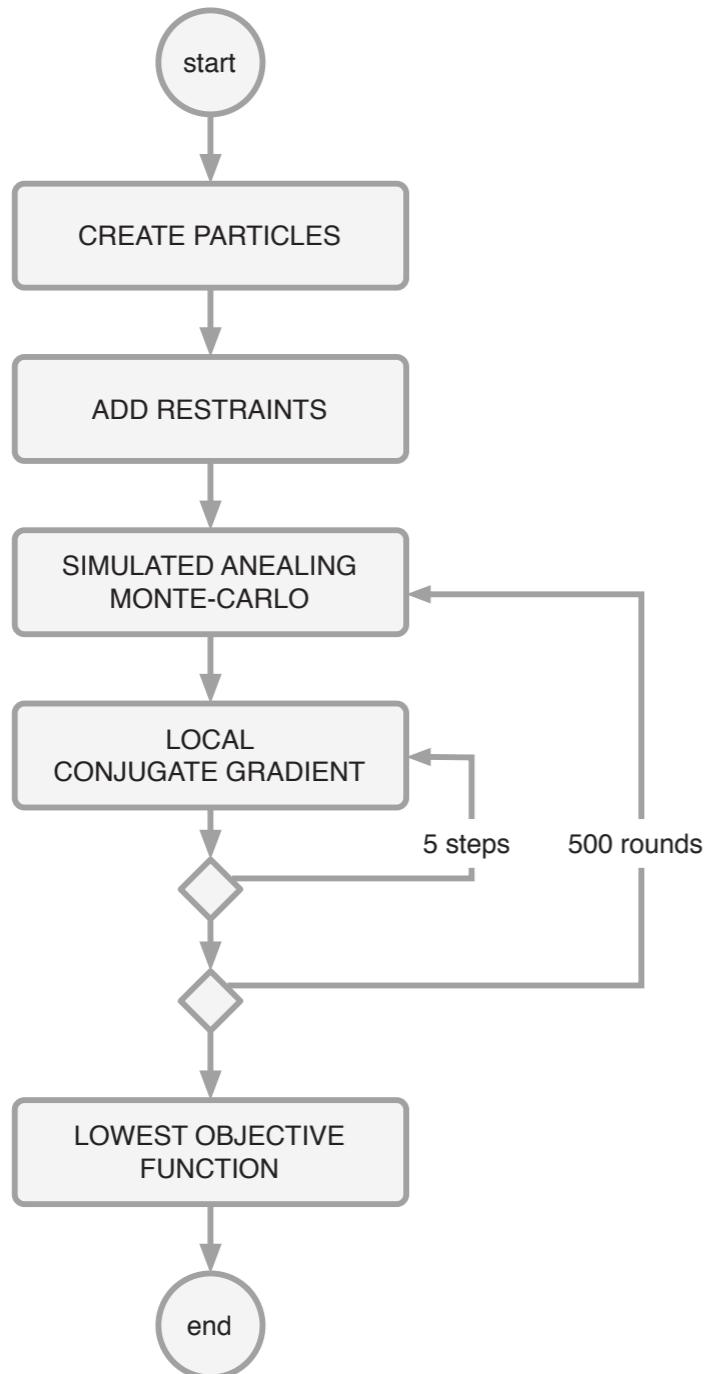
Harmonic

Harmonic Lower Bound

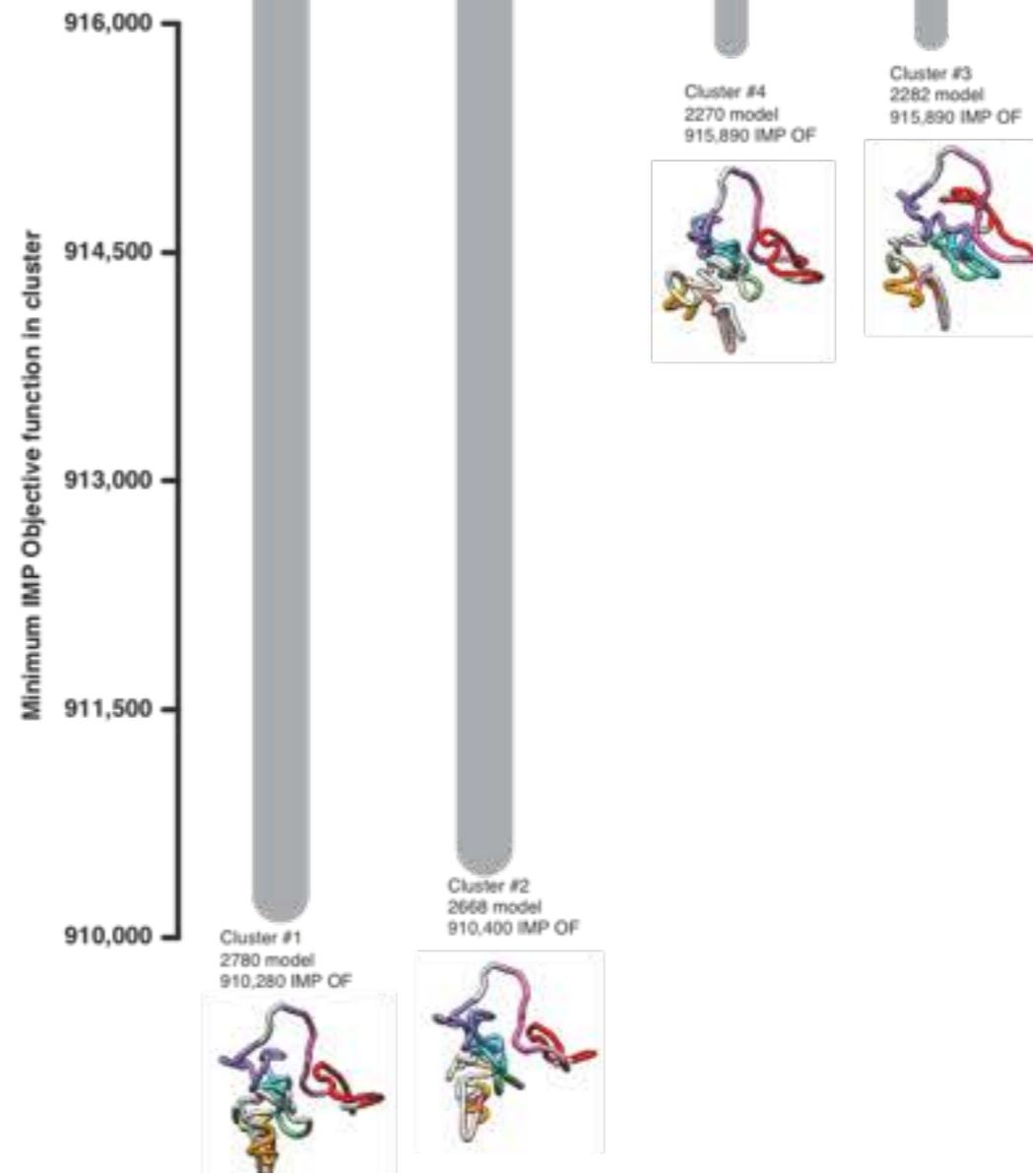
Harmonic Upper Bound



Optimization

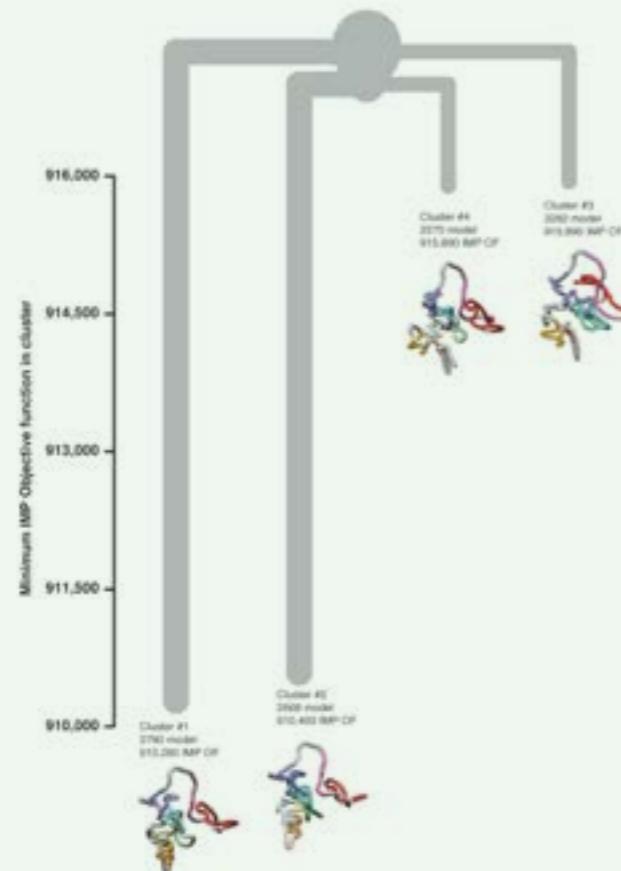
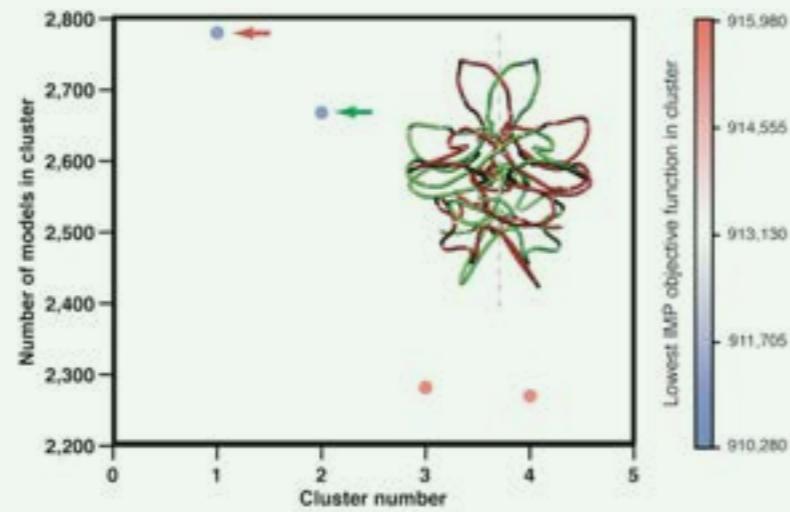


Clustering

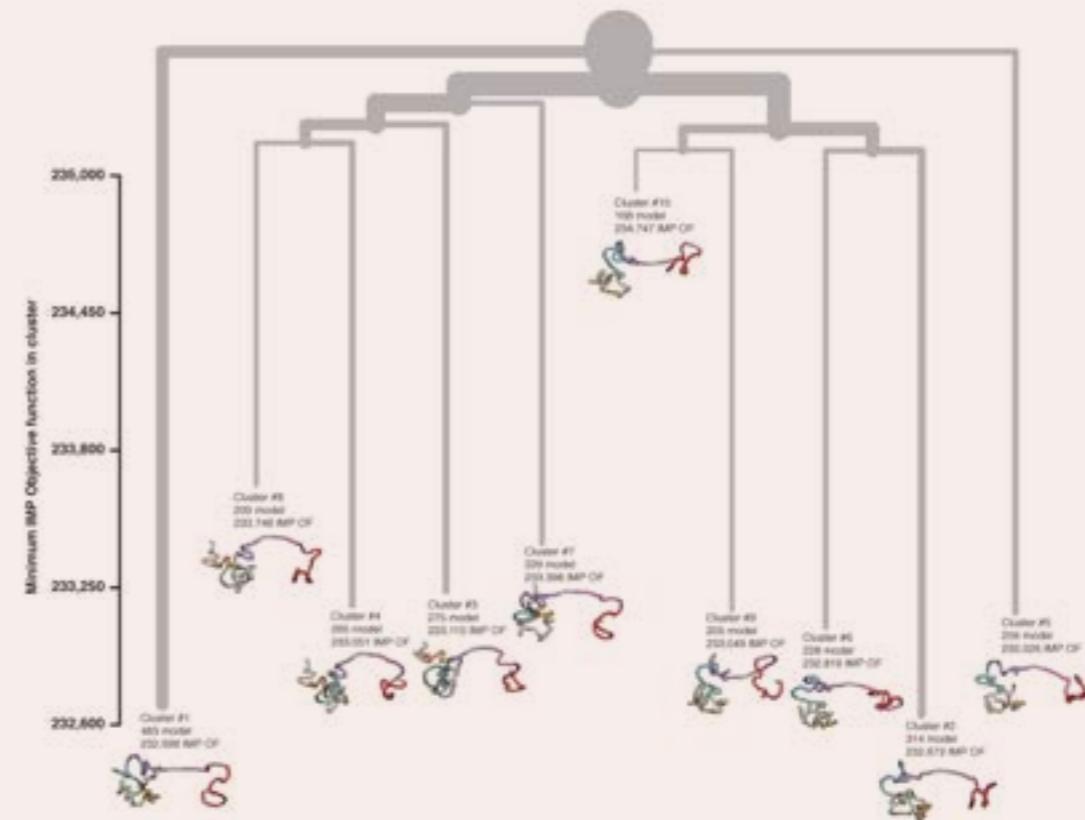
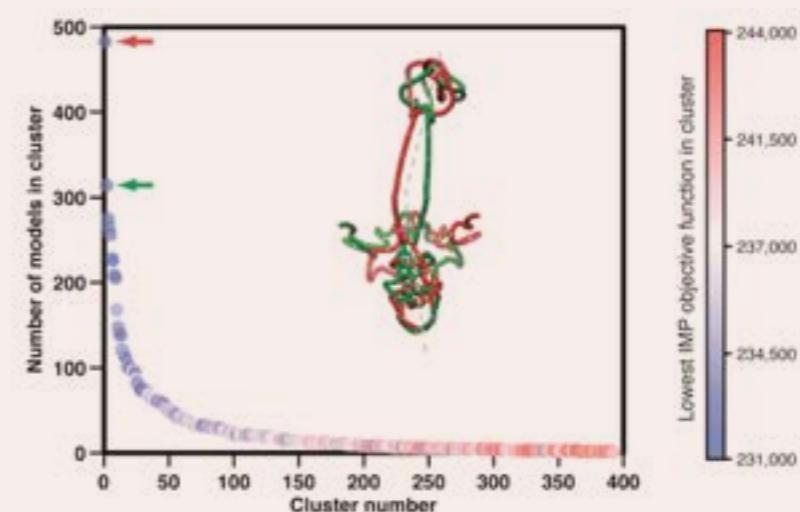


Not just one solution

GM12878



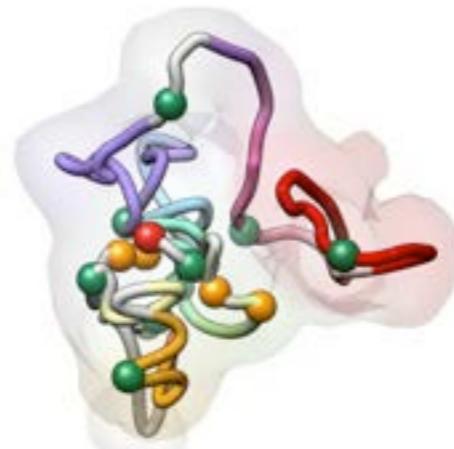
K562



Regulatory compartmentalization

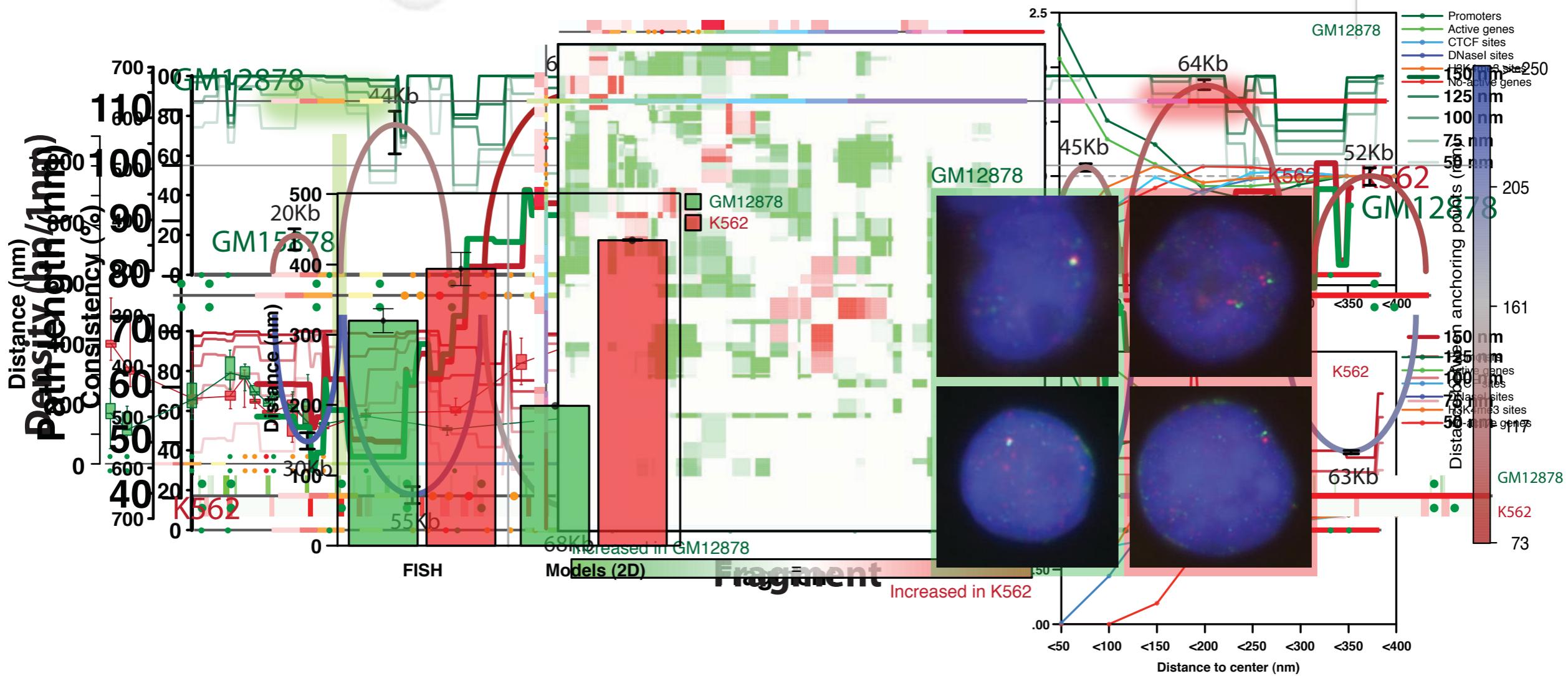
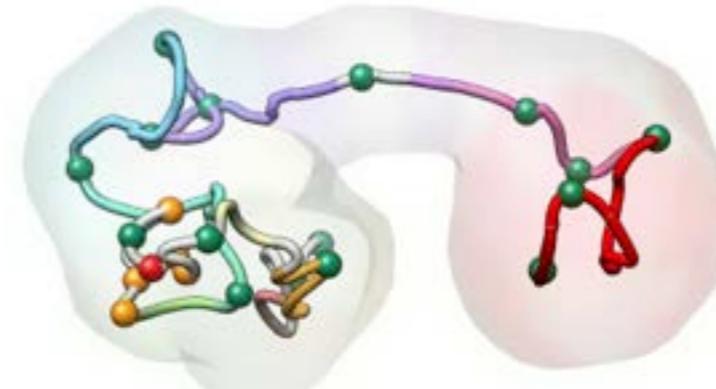
GM12878

Cluster #1
2780 model

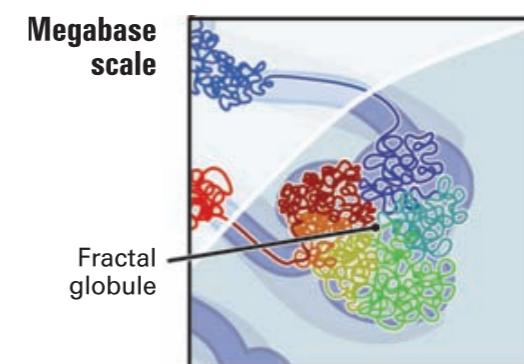
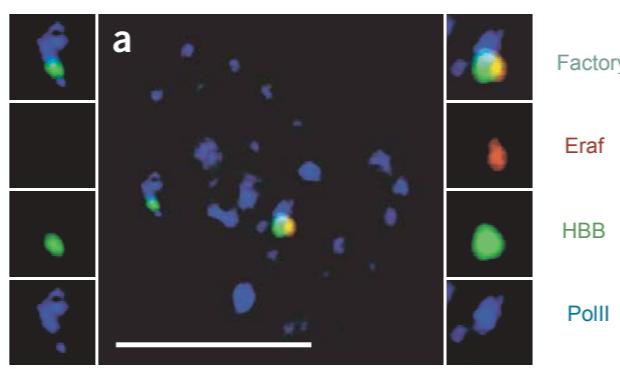
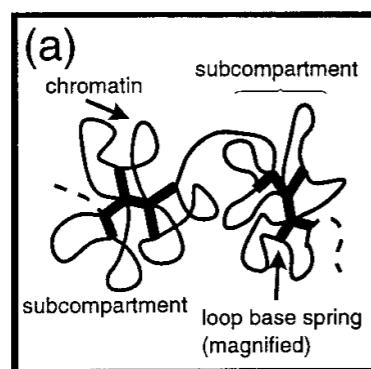
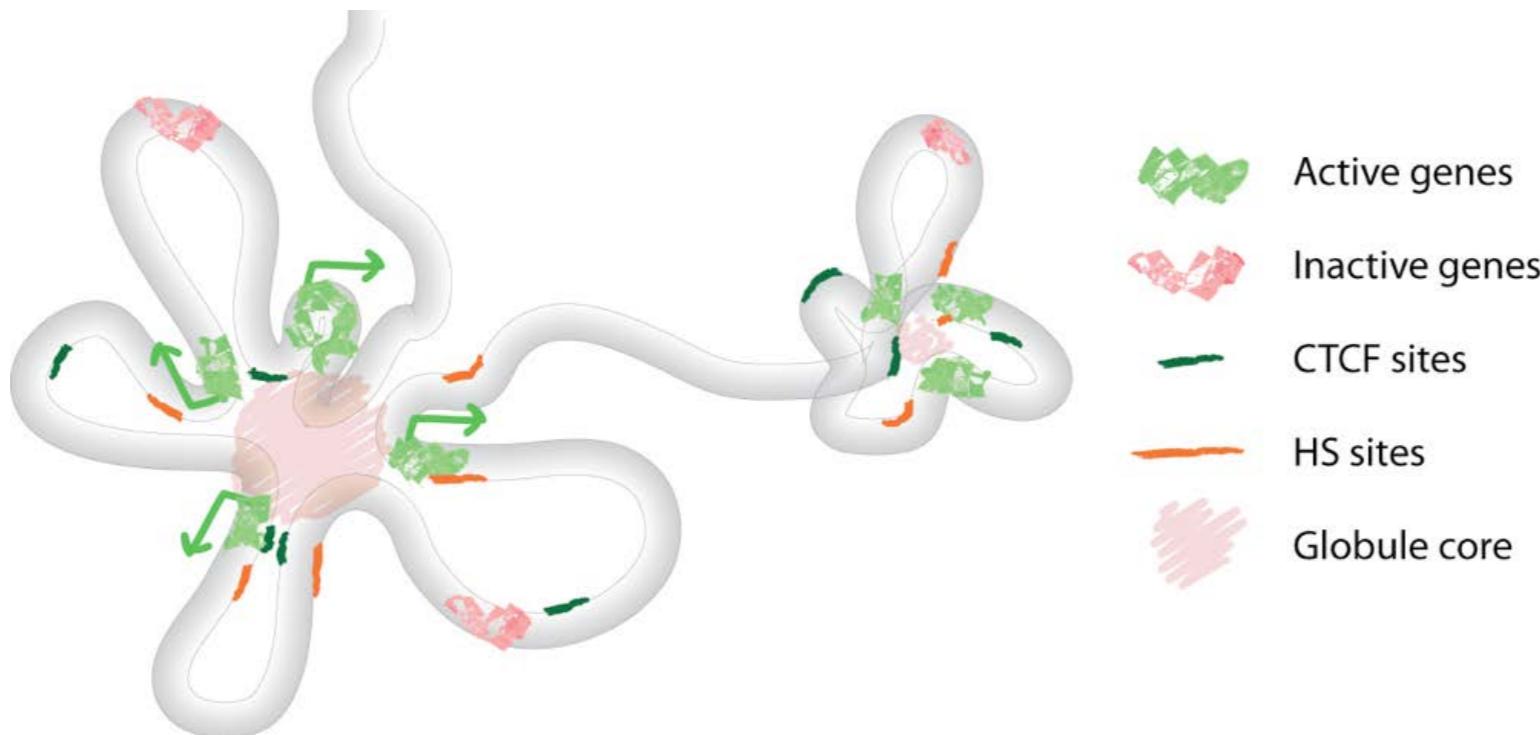


K562

Cluster #2
314 model

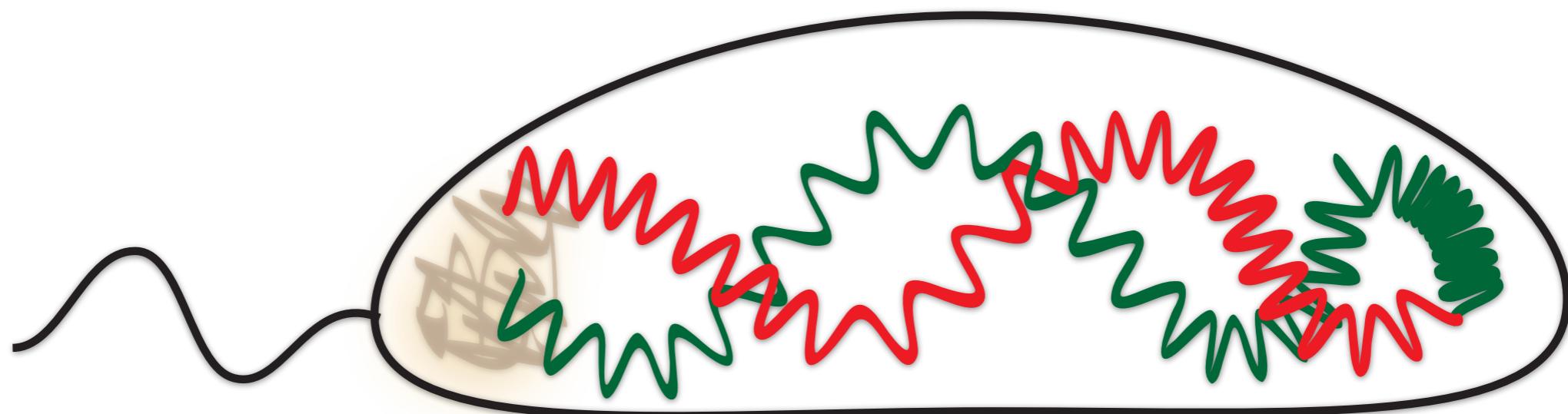


The “Chromatin Globule” model



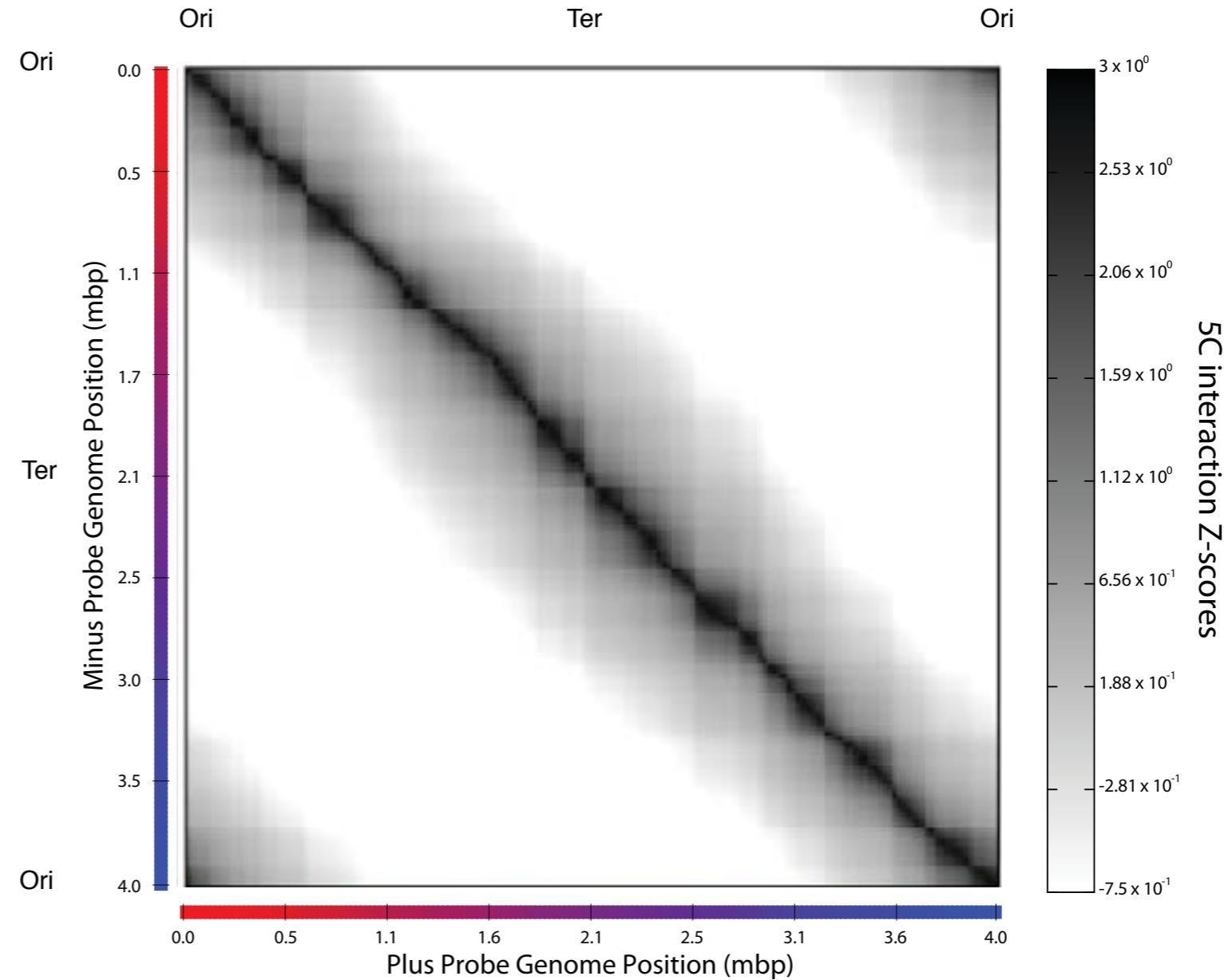
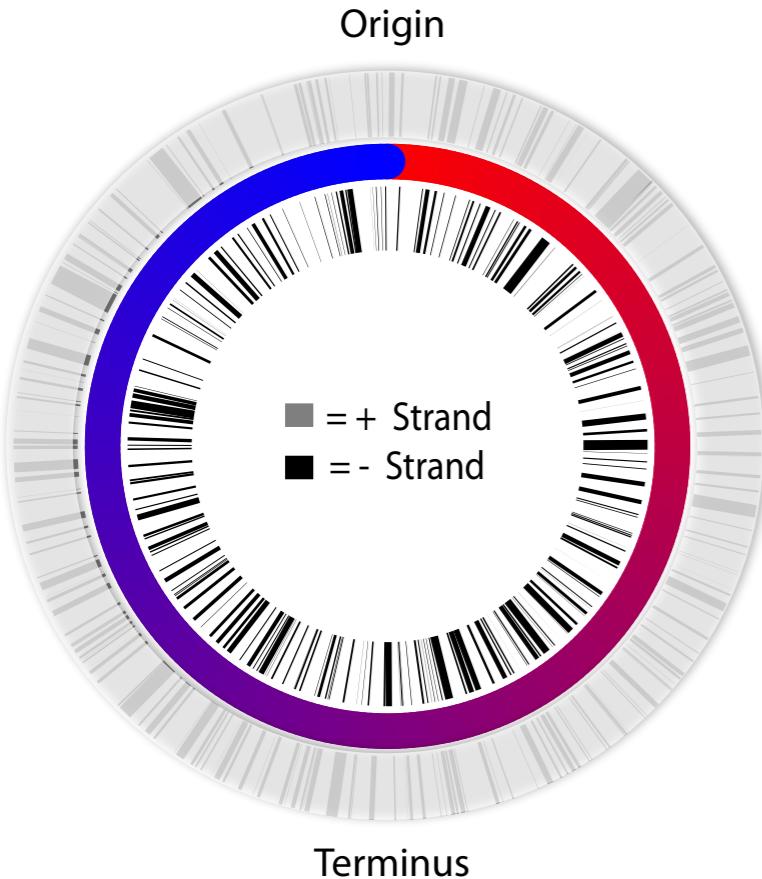
D. Baù et al. Nat Struct Mol Biol (2011) 18:107-14
A. Sanyal et al. Current Opinion in Cell Biology (2011) 23:325–33.

Caulobacter crescentus genome



The 3D architecture of *Caulobacter Crescentus*

4,016,942 bp & 3,767 genes

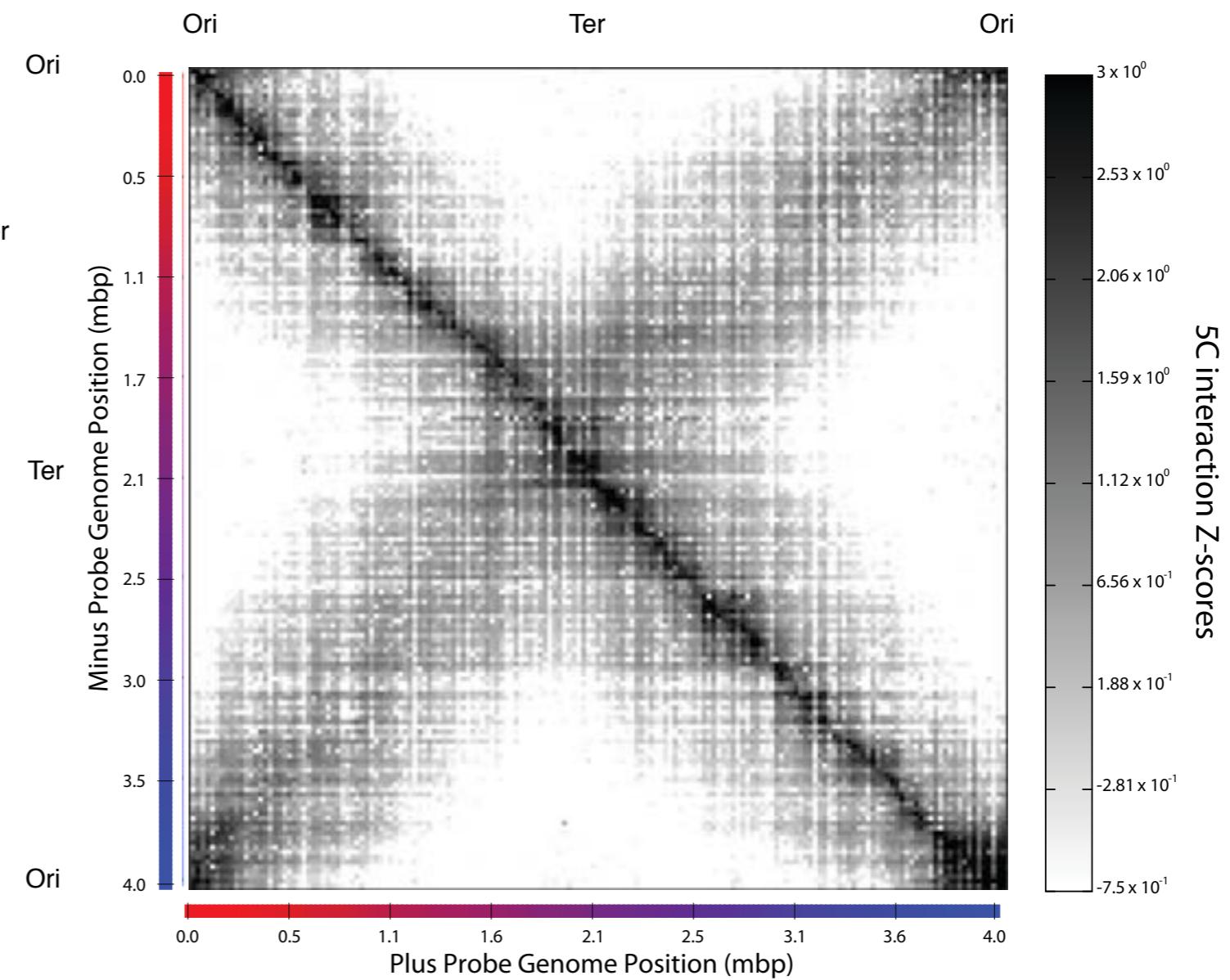
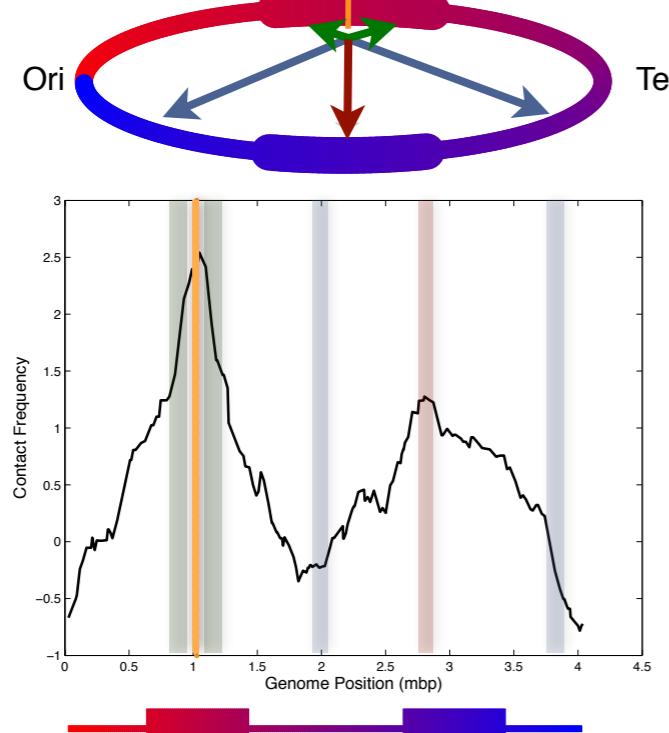
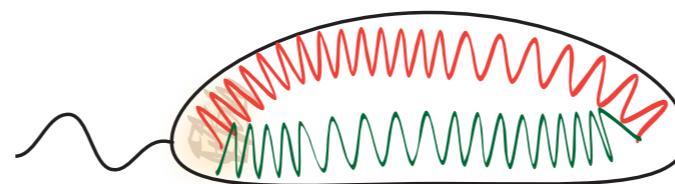


169 5C primers on + strand
170 5C primers on - strand
28,730 chromatin interactions

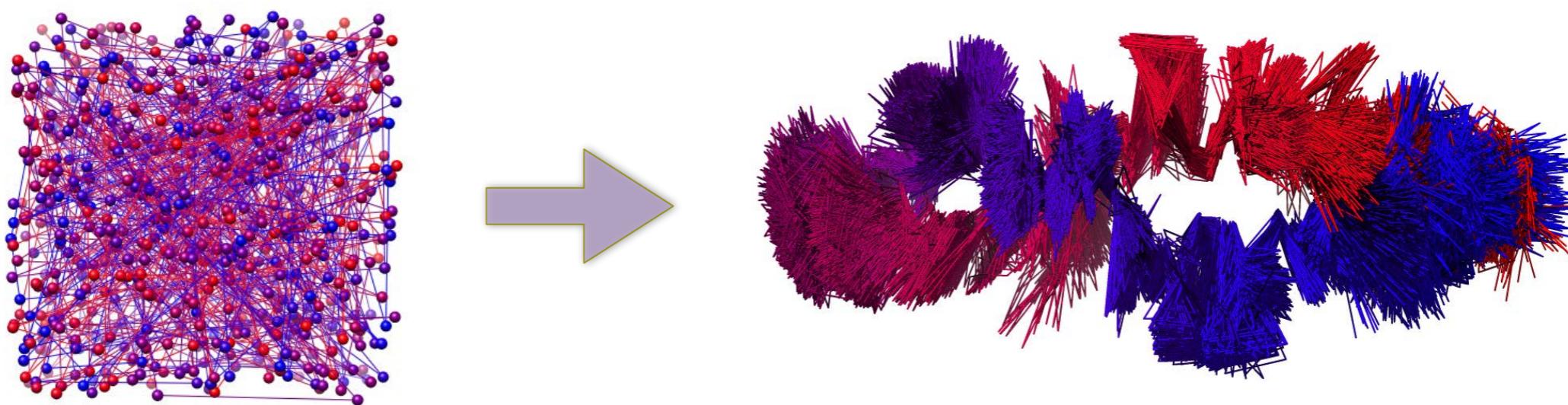
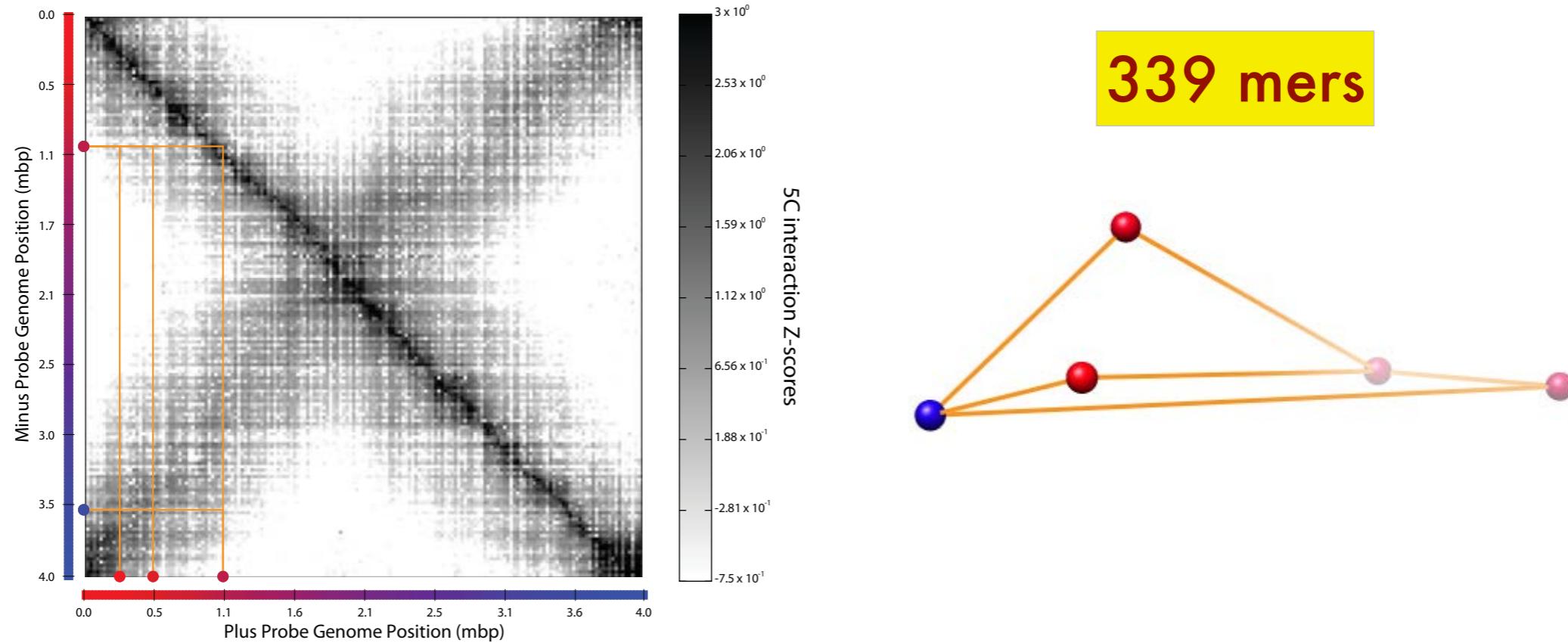
~13Kb

5C interaction matrix

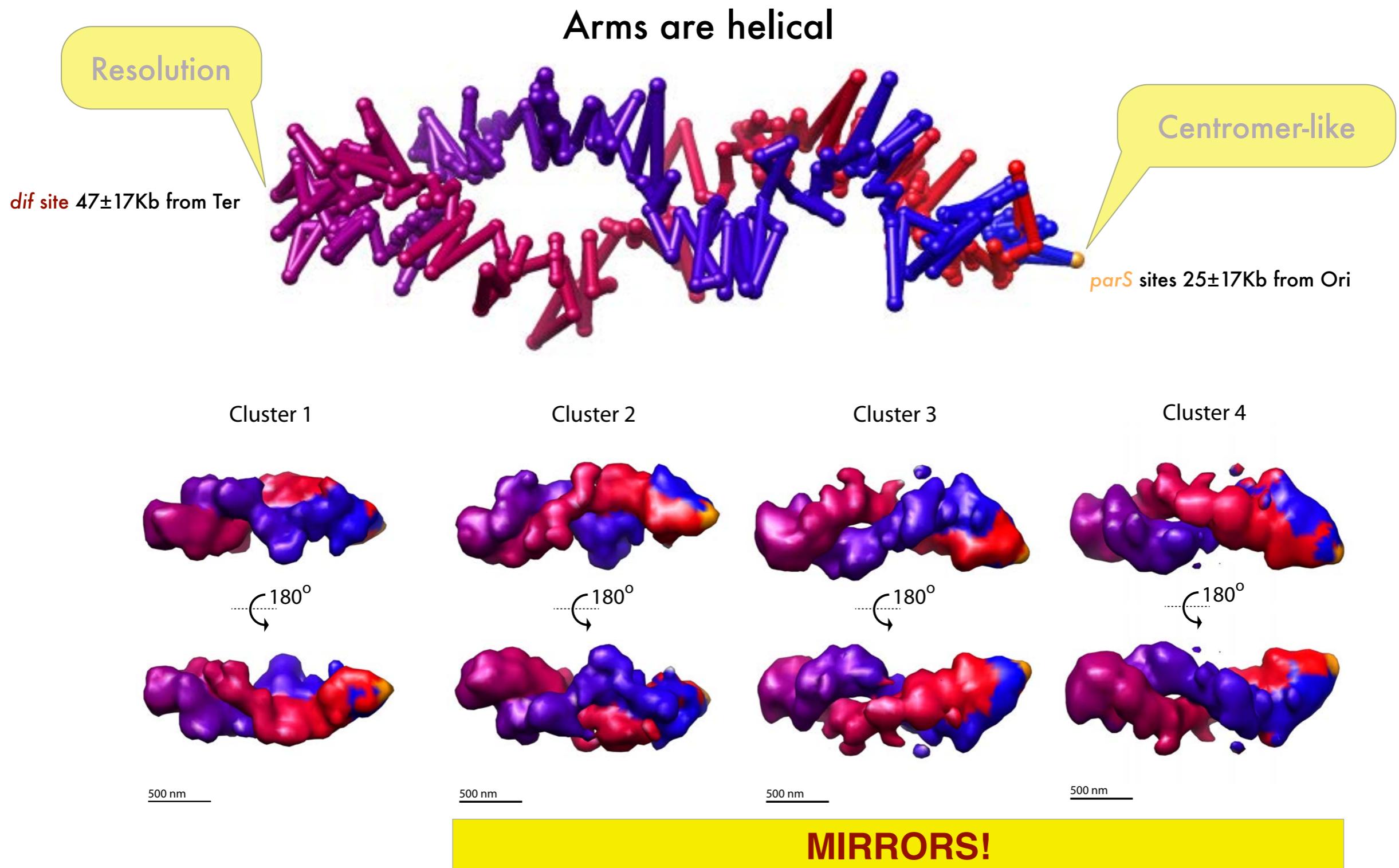
ELLIPSOID for *Caulobacter crescentus*



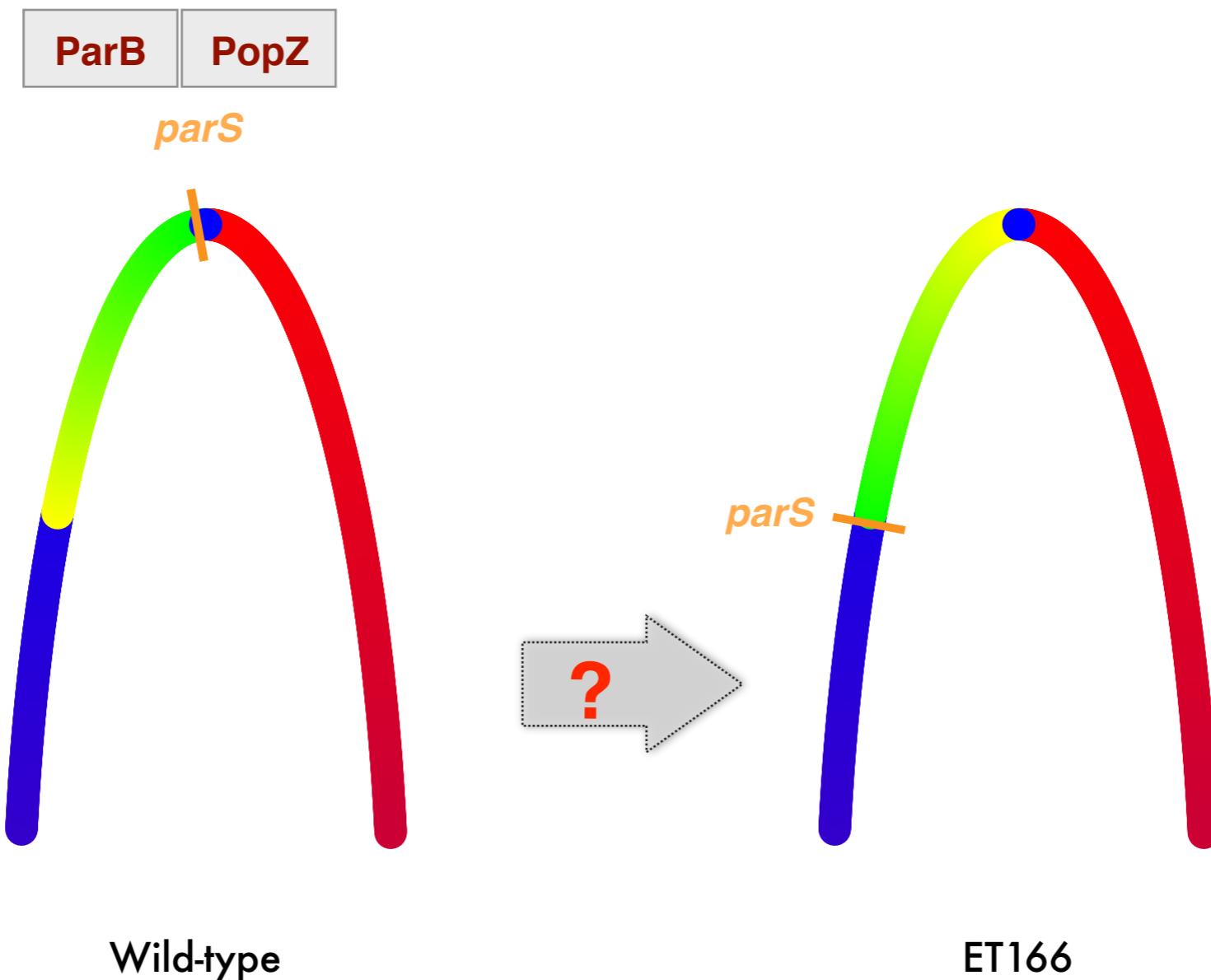
3D model building with the 5C + IMP approach



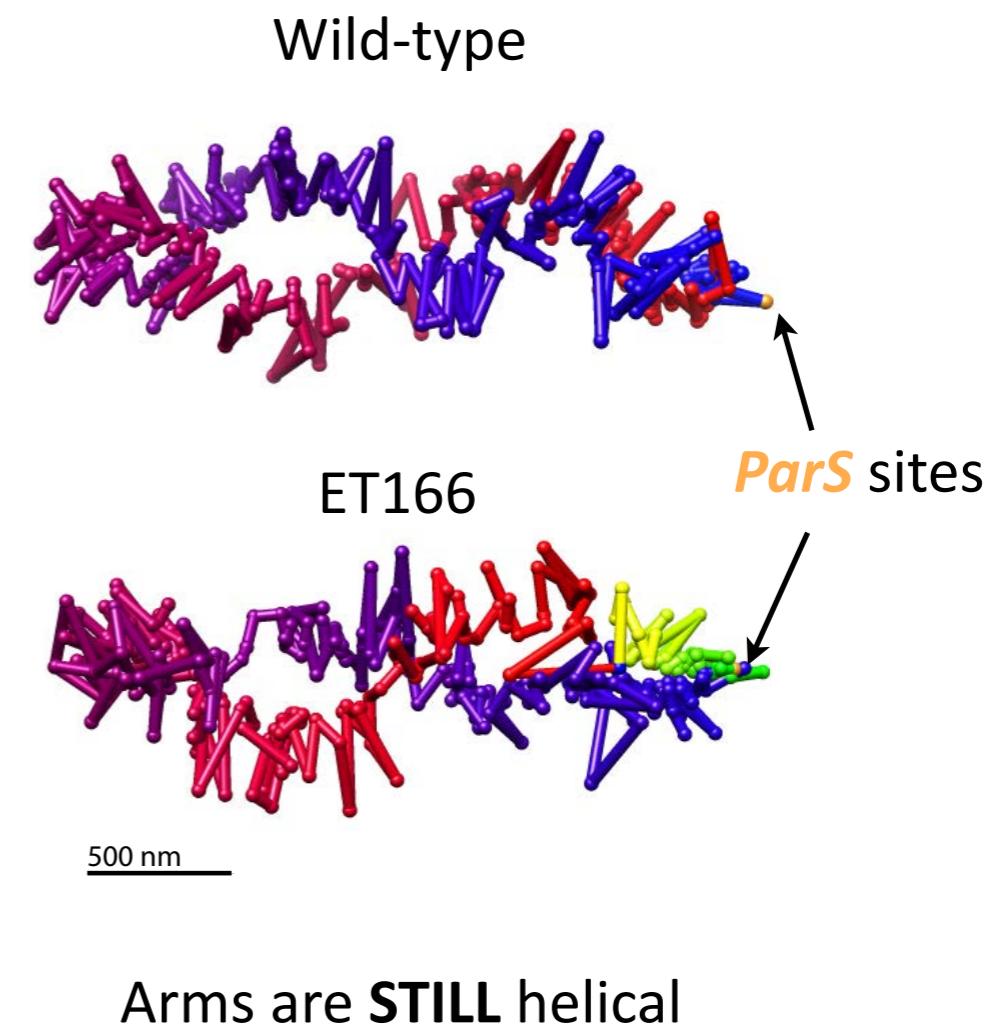
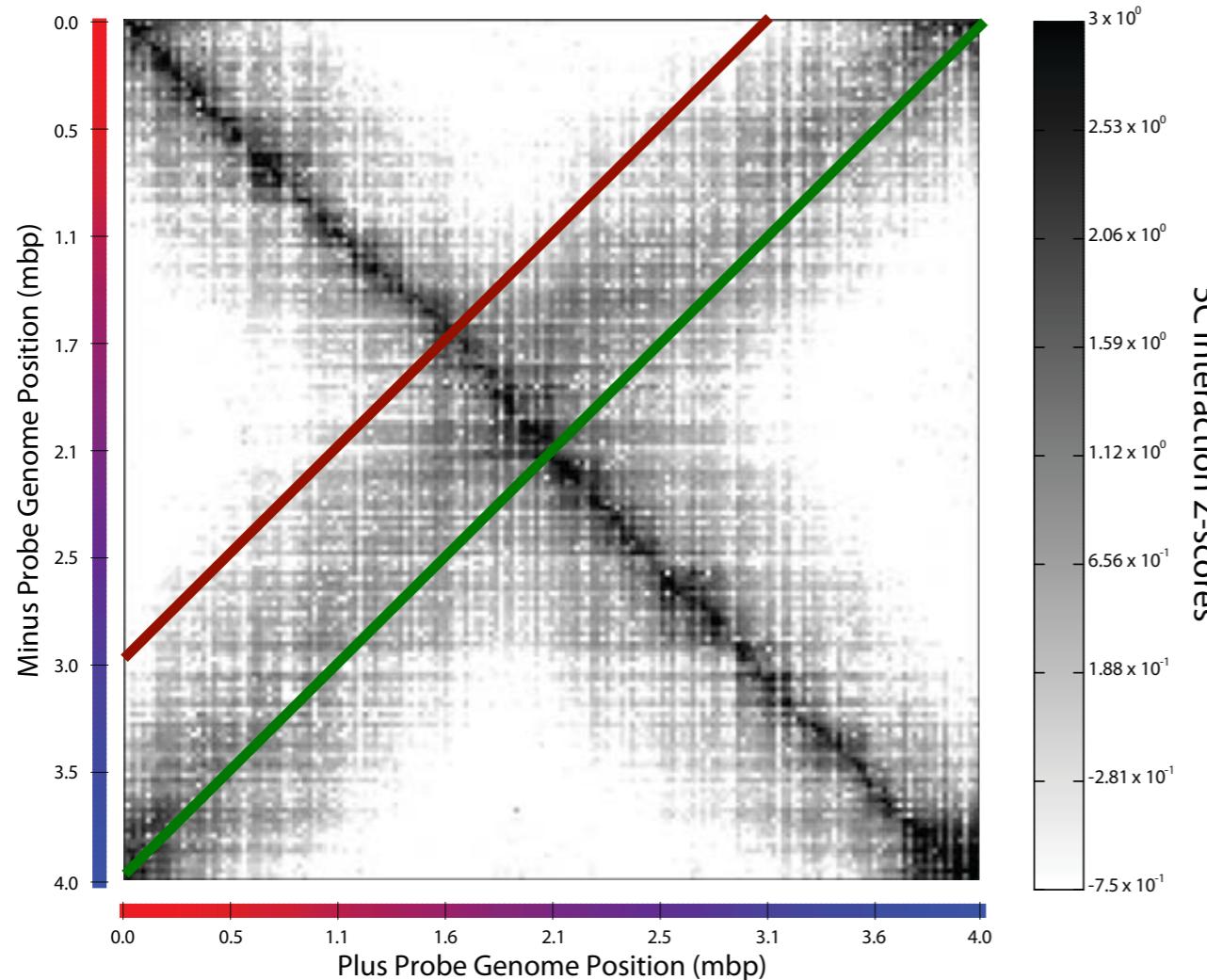
Genome organization in *Caulobacter crescentus*



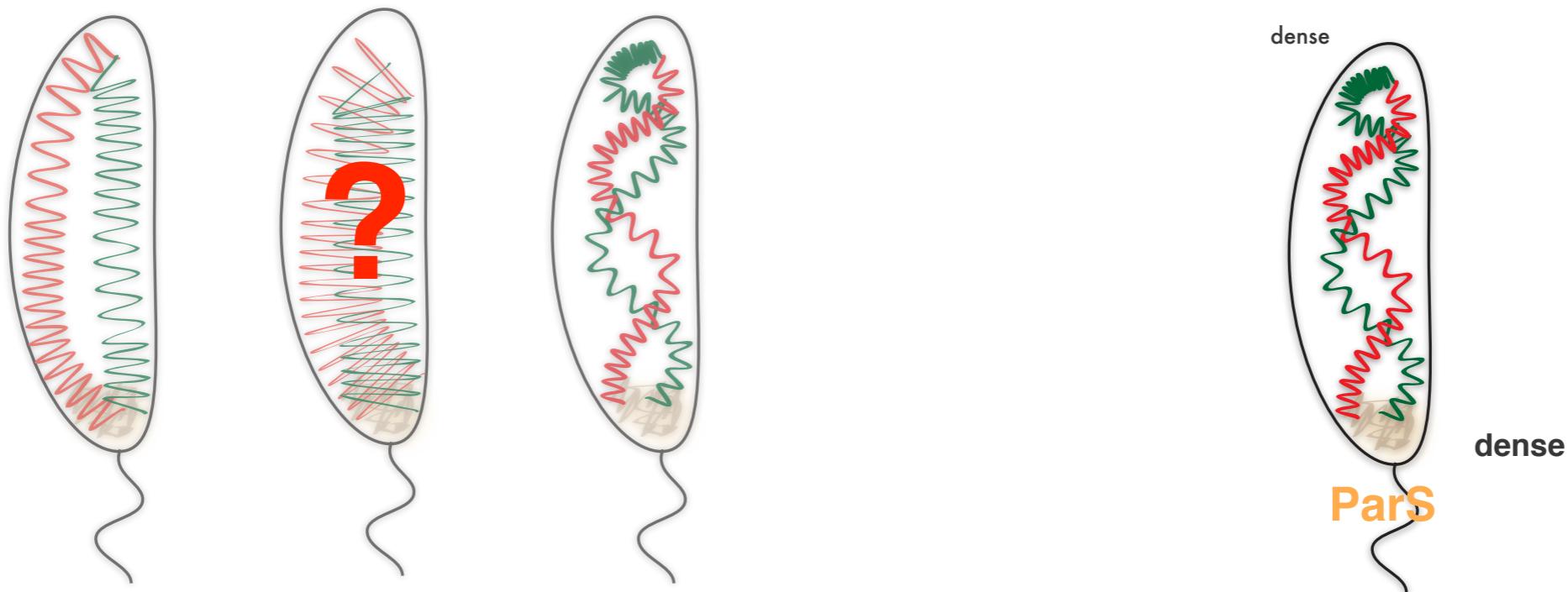
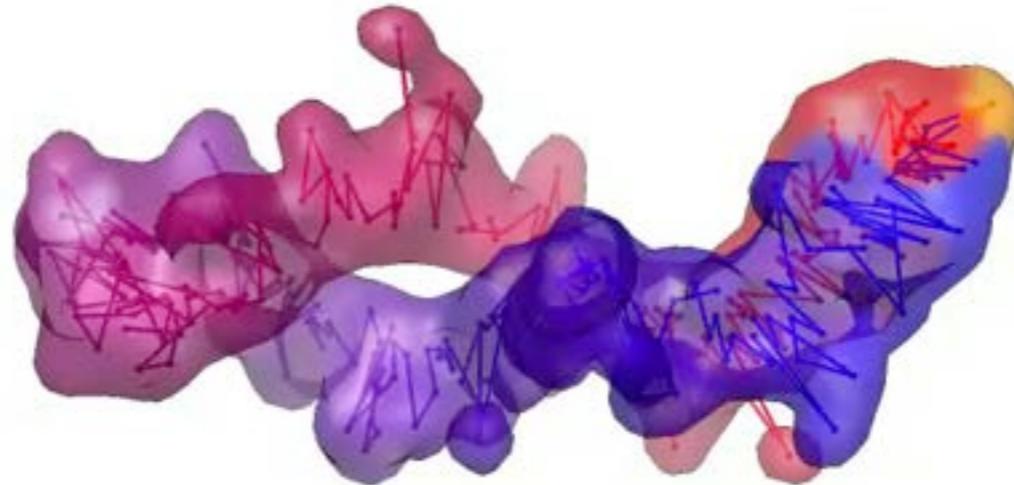
Moving the *parS* sites 400 Kb away from Ori



Moving the *parS* sites results in whole genome rotation!

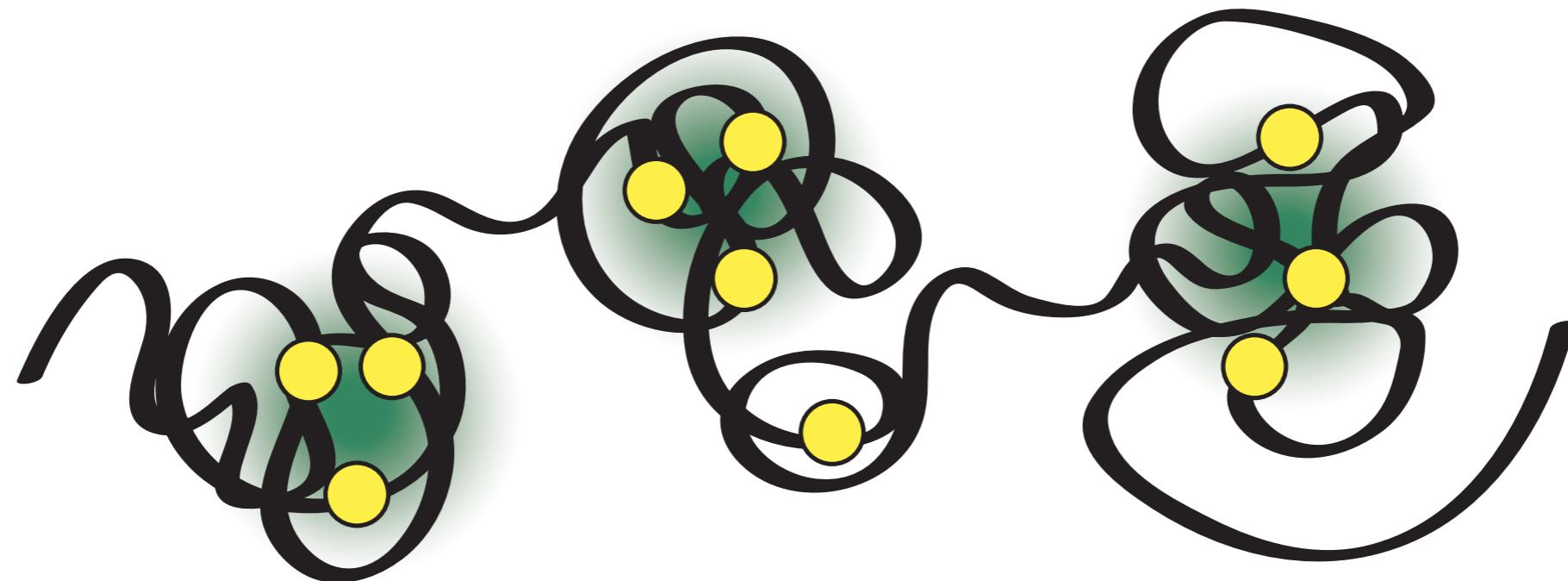


Genome architecture in Caulobacter

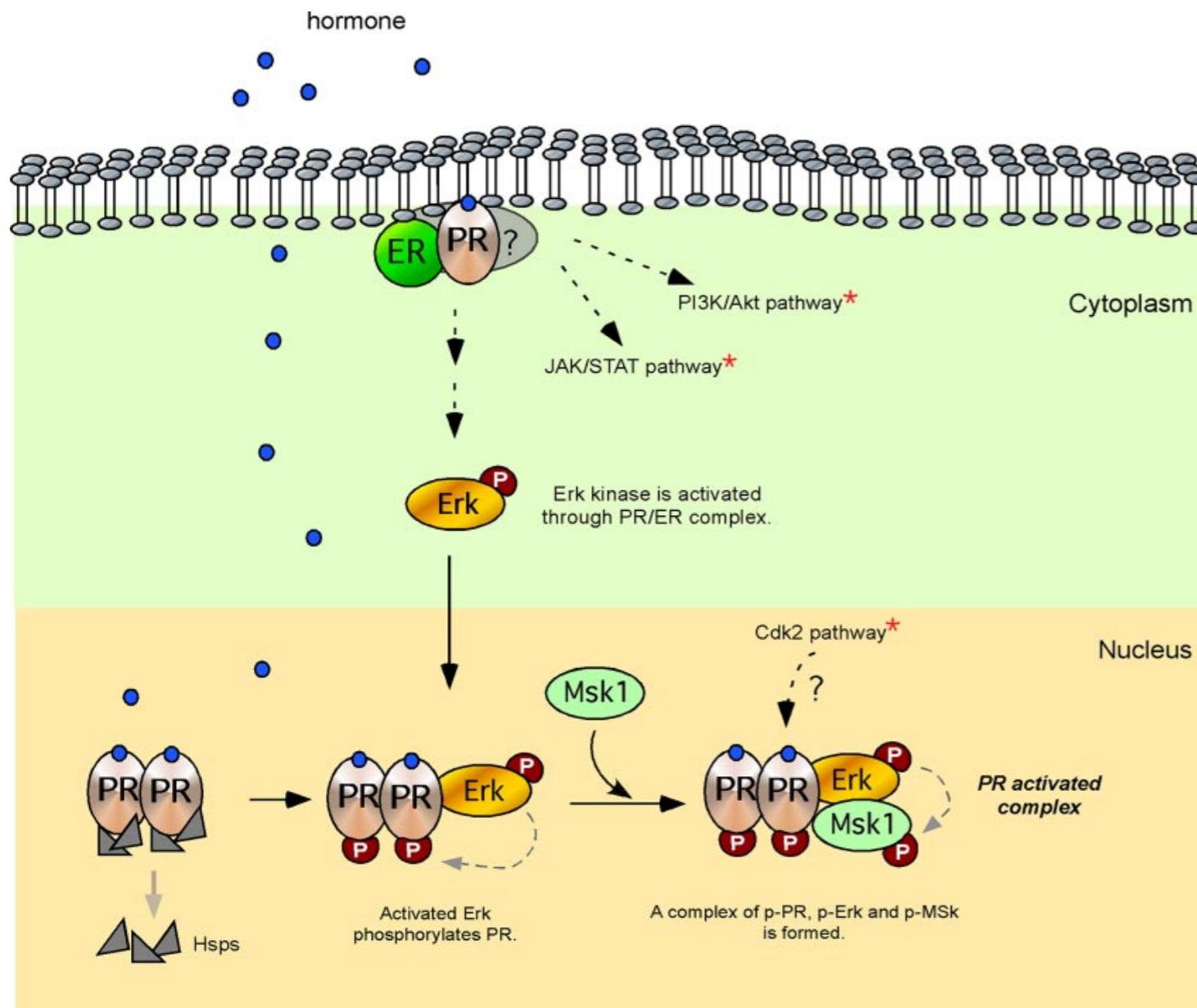


M.A. Umbarger, et al. Molecular Cell (2011) 44:252–264

On TADs and hormones



Progesterone-regulated transcription in breast cancer

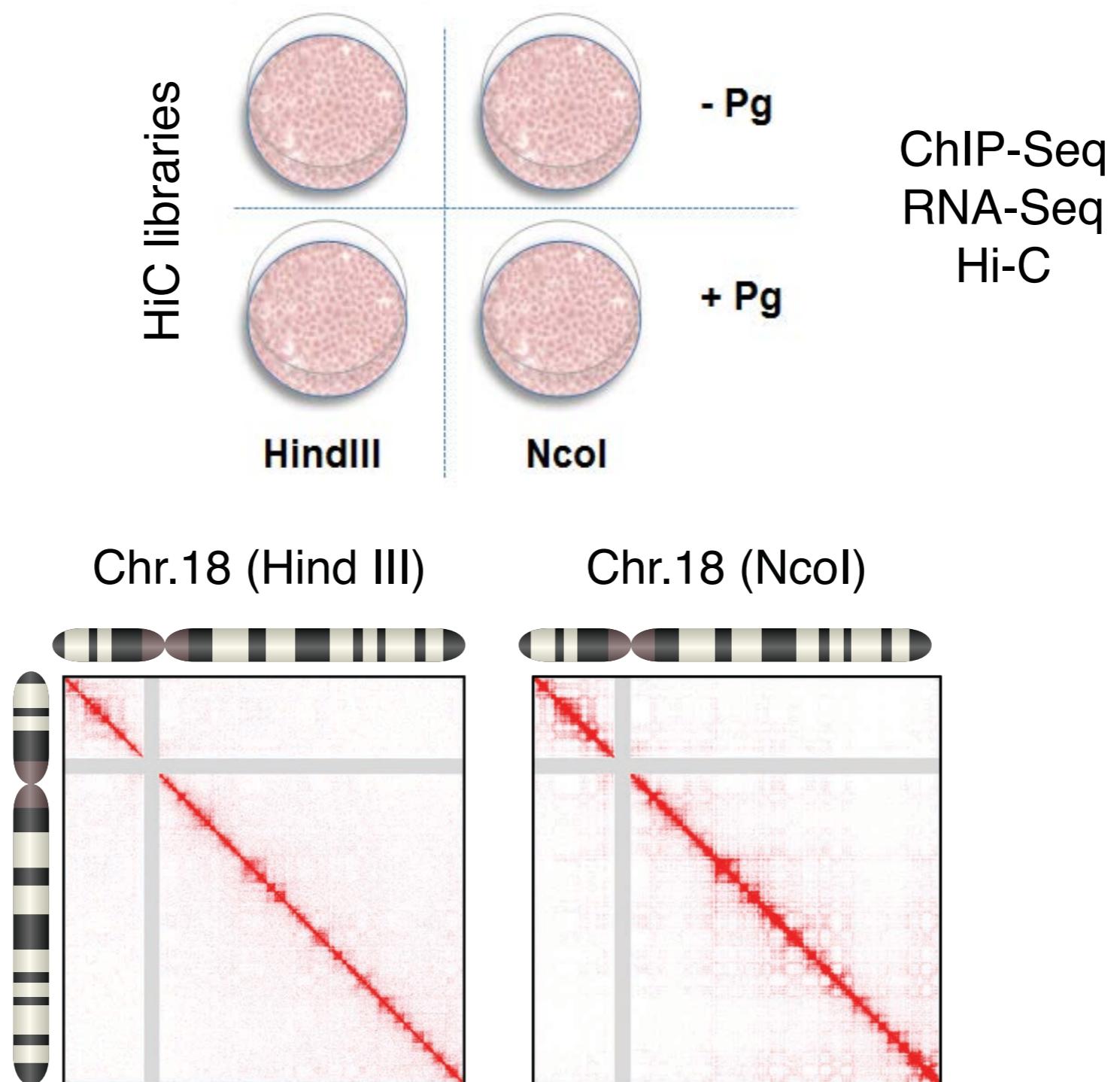


> 2,000 genes Up-regulated
> 2,000 genes Down-regulated

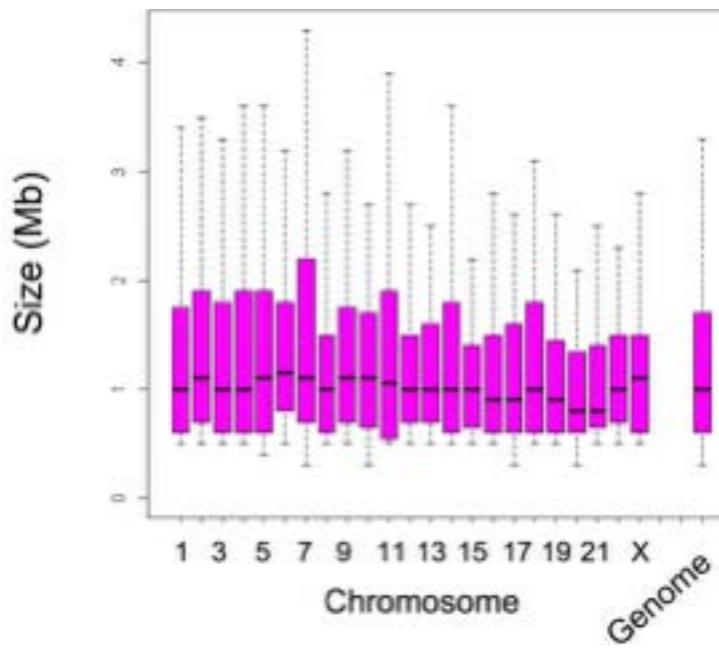
Regulation in 3D?

Vicent *et al* 2011, Wright *et al* 2012, Ballare *et al* 2012

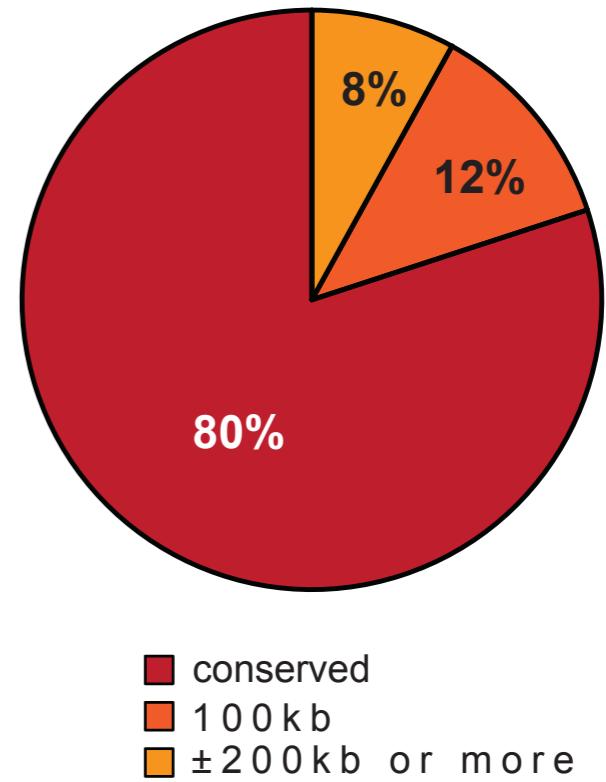
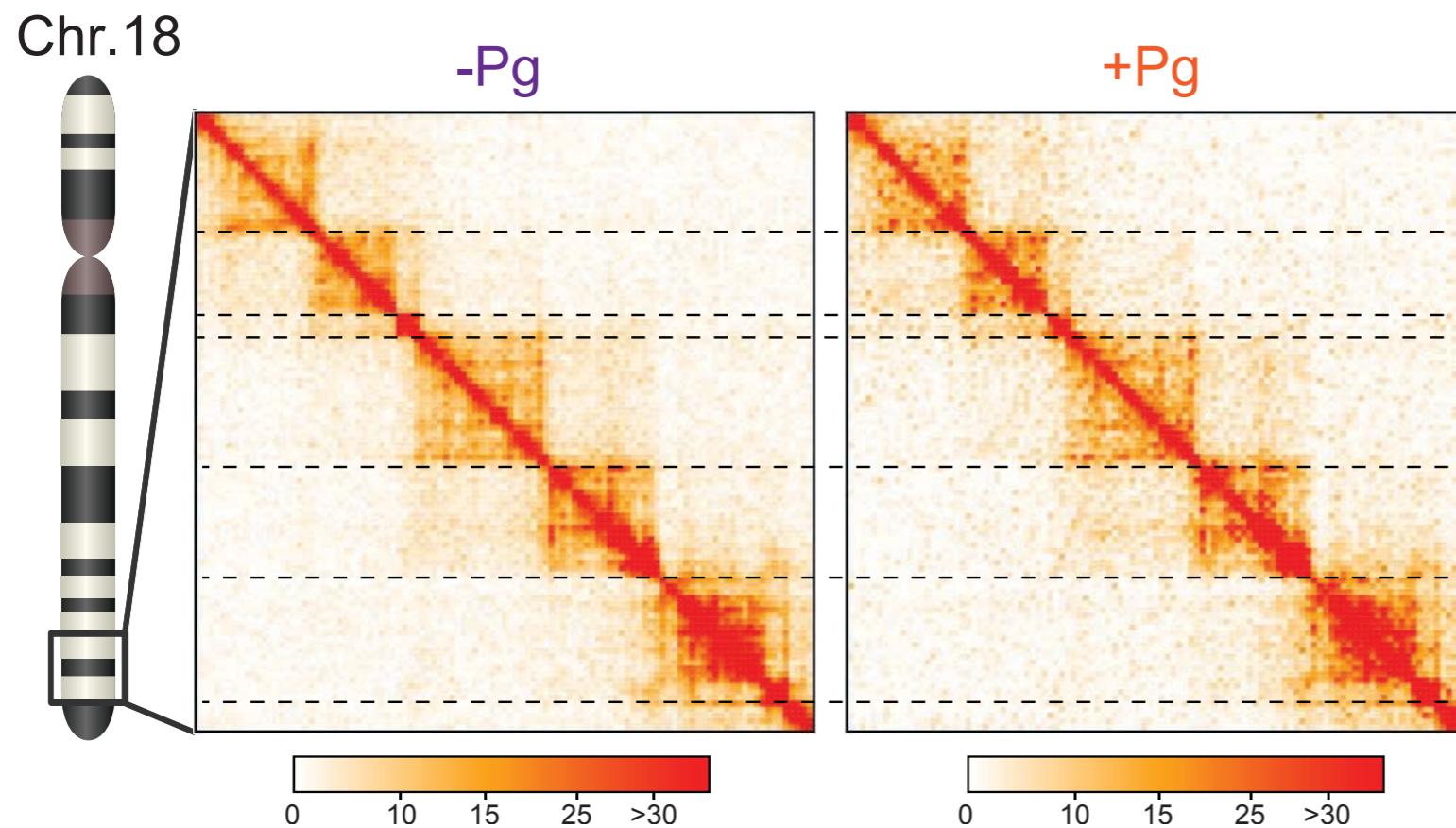
Experimental design



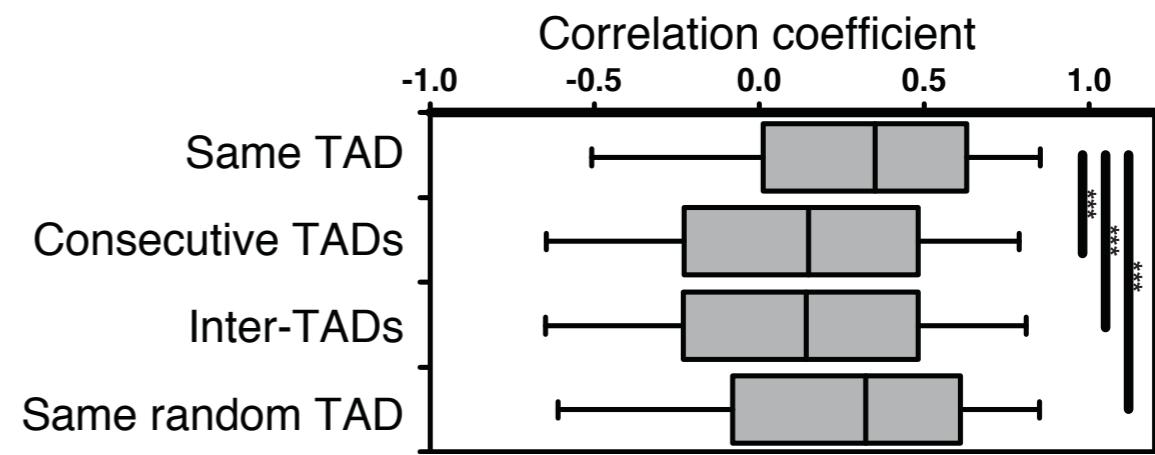
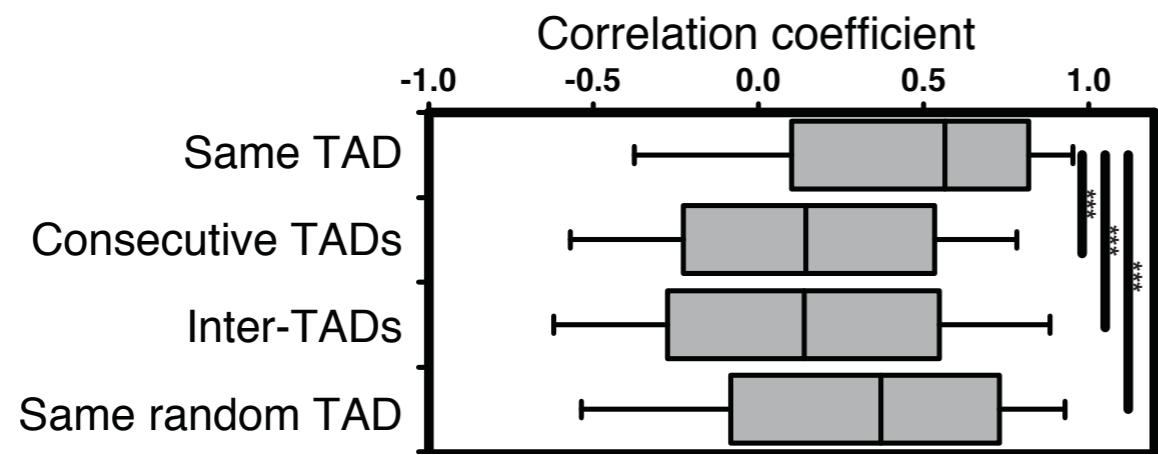
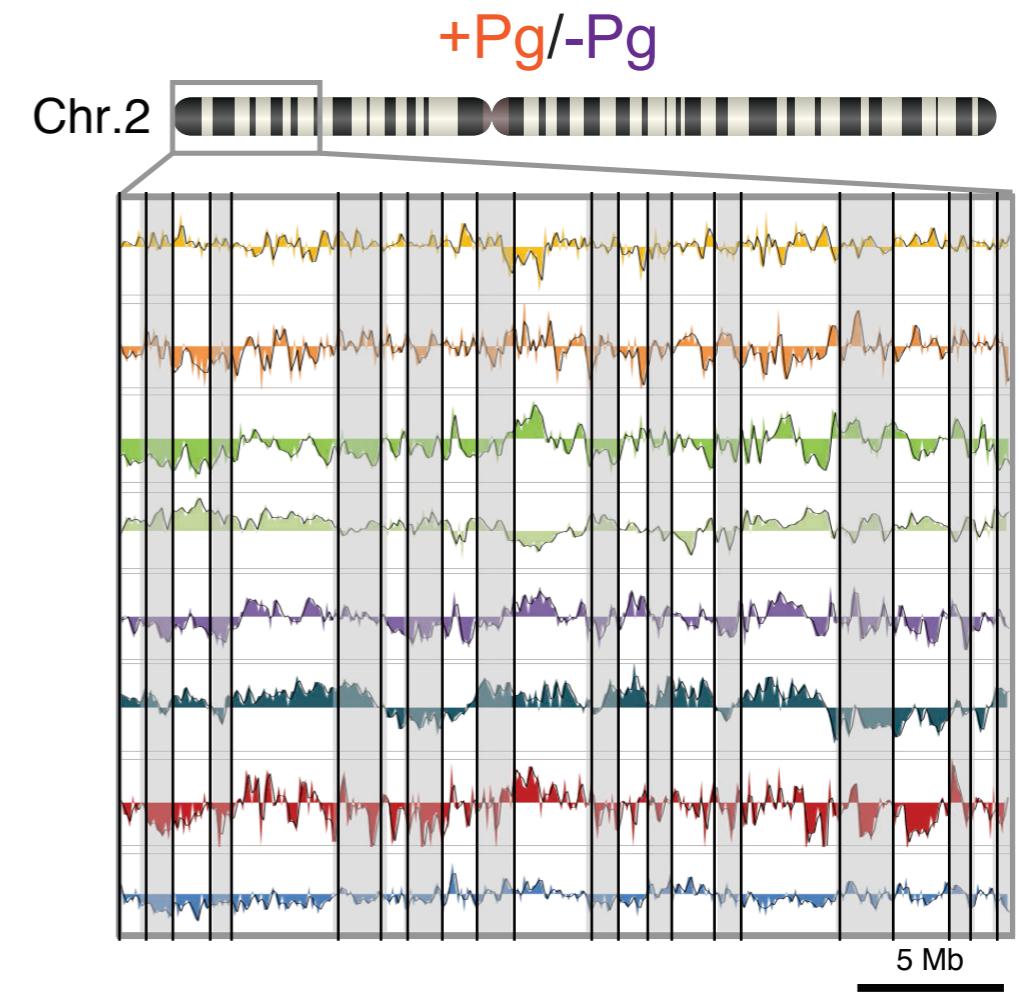
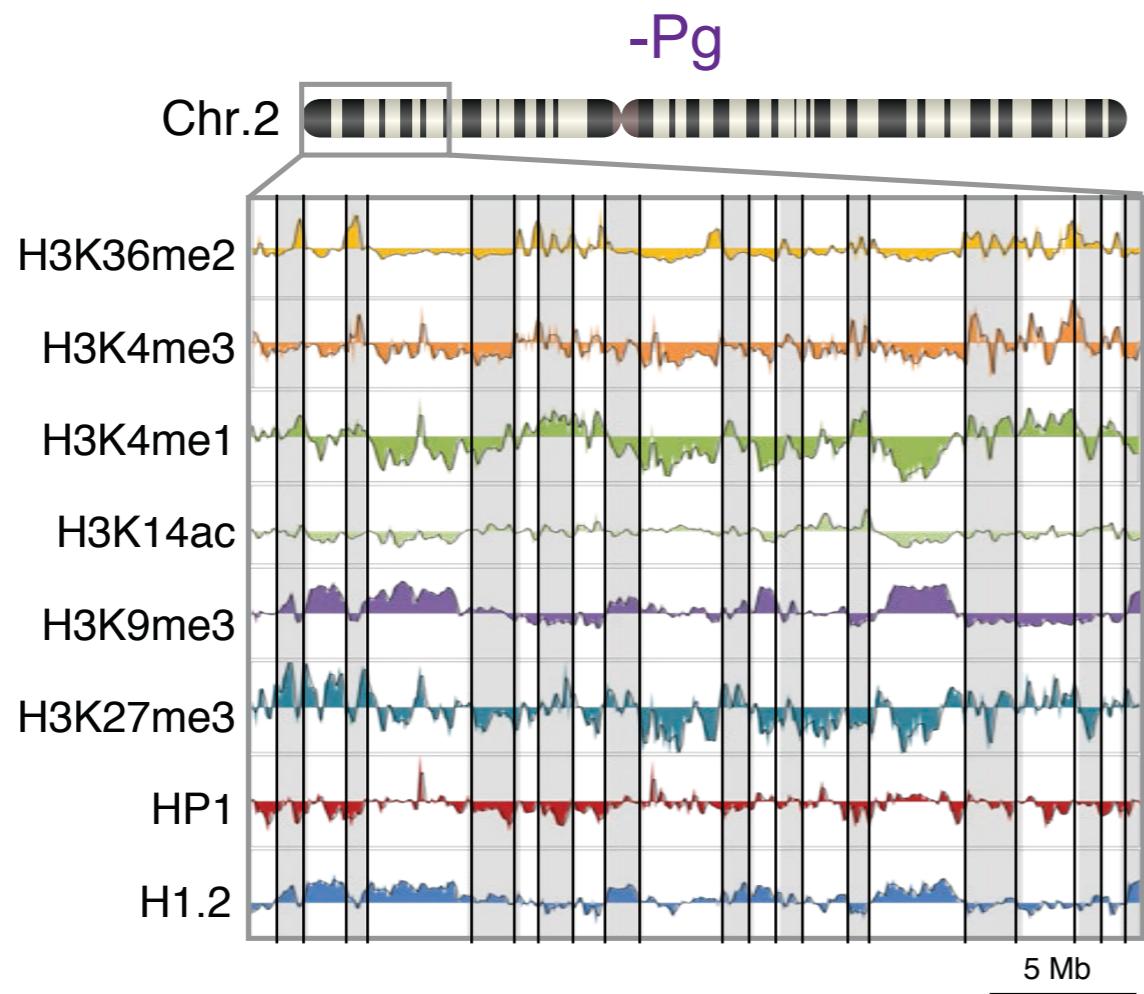
Are there TADs? how robust?



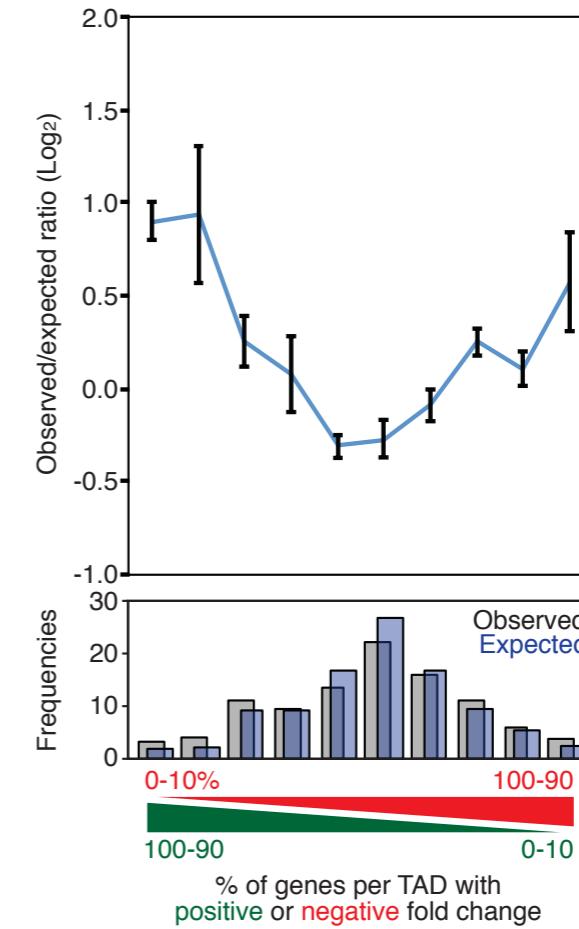
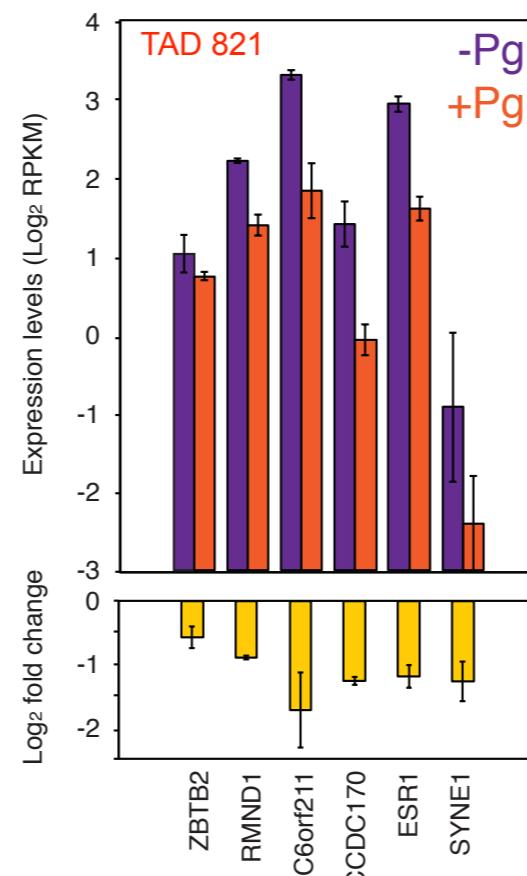
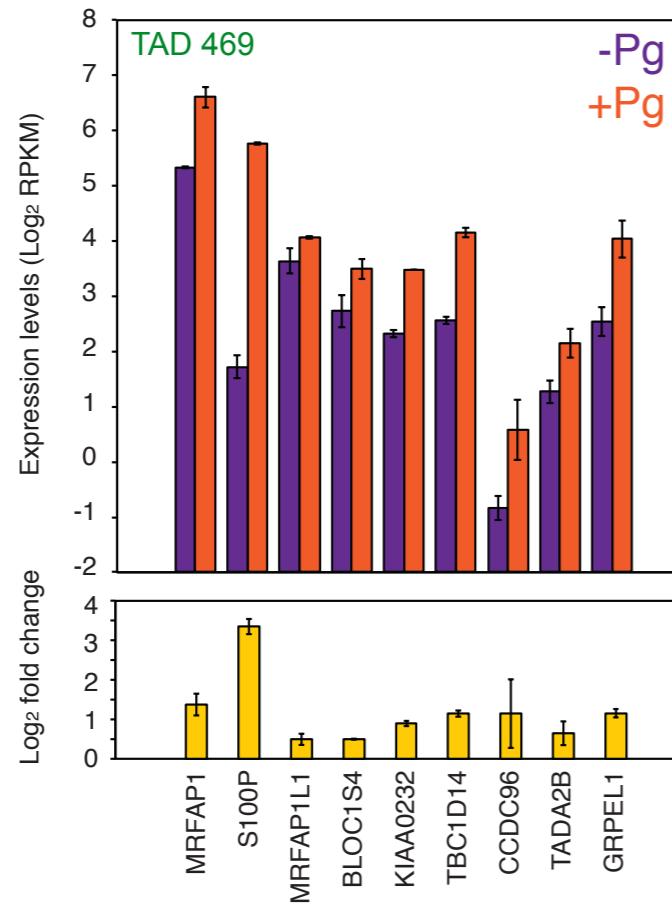
>2,000 detected TADs



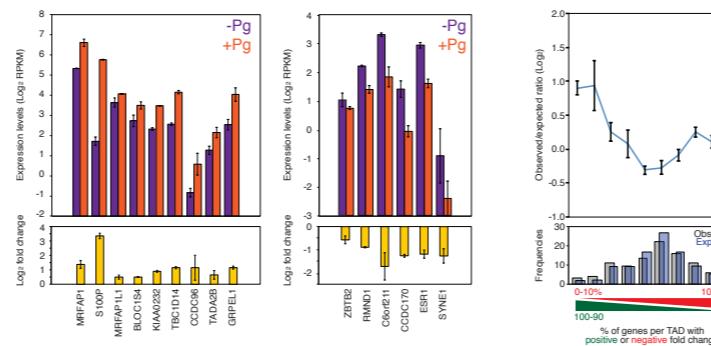
Are TADs homogeneous?



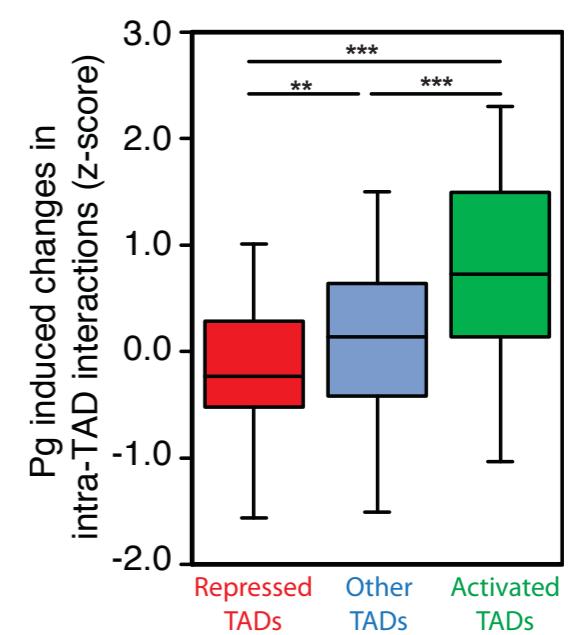
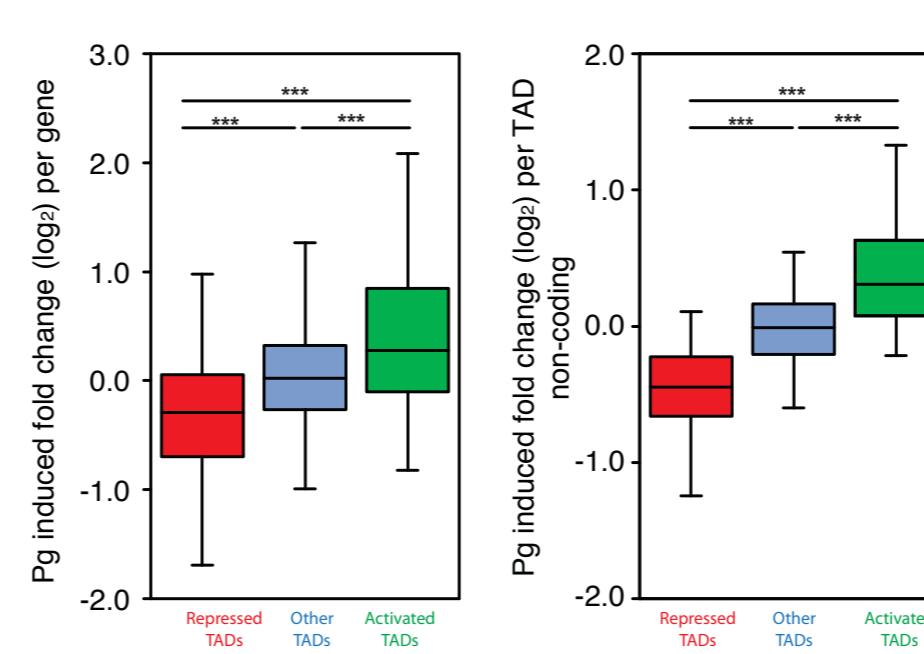
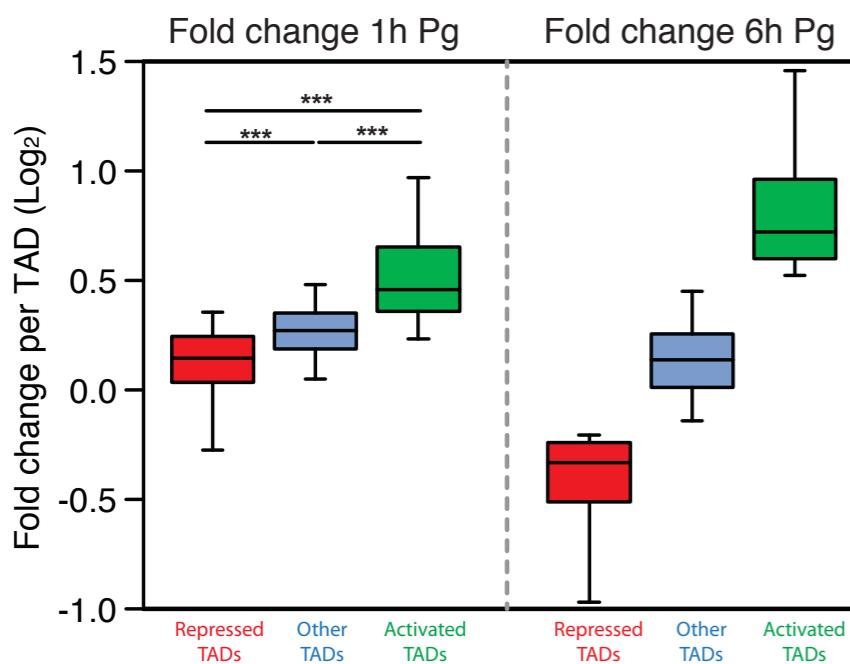
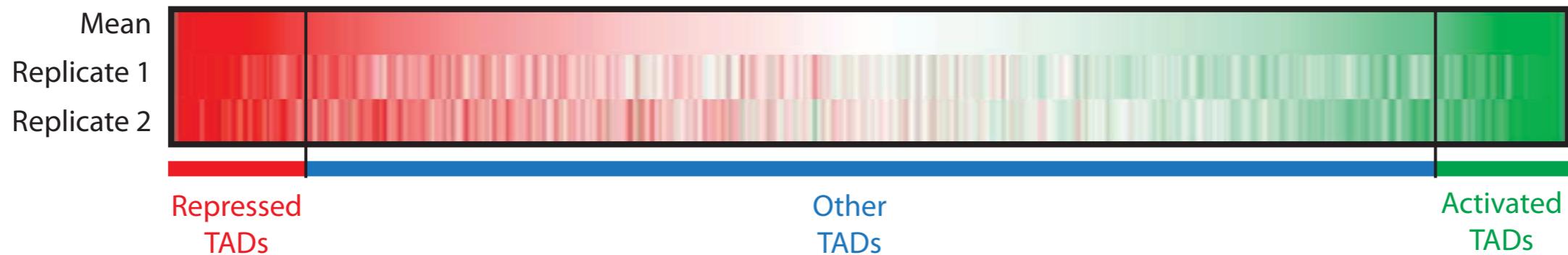
Do TADs respond differently to Pg treatment?



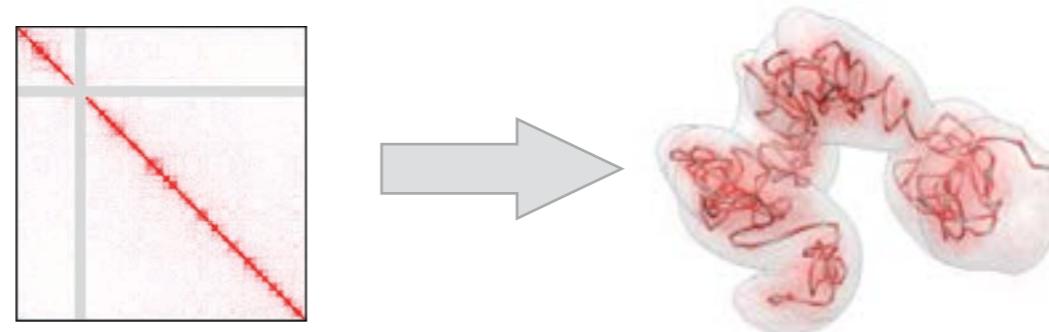
Do TADs respond differently to Pg treatment?



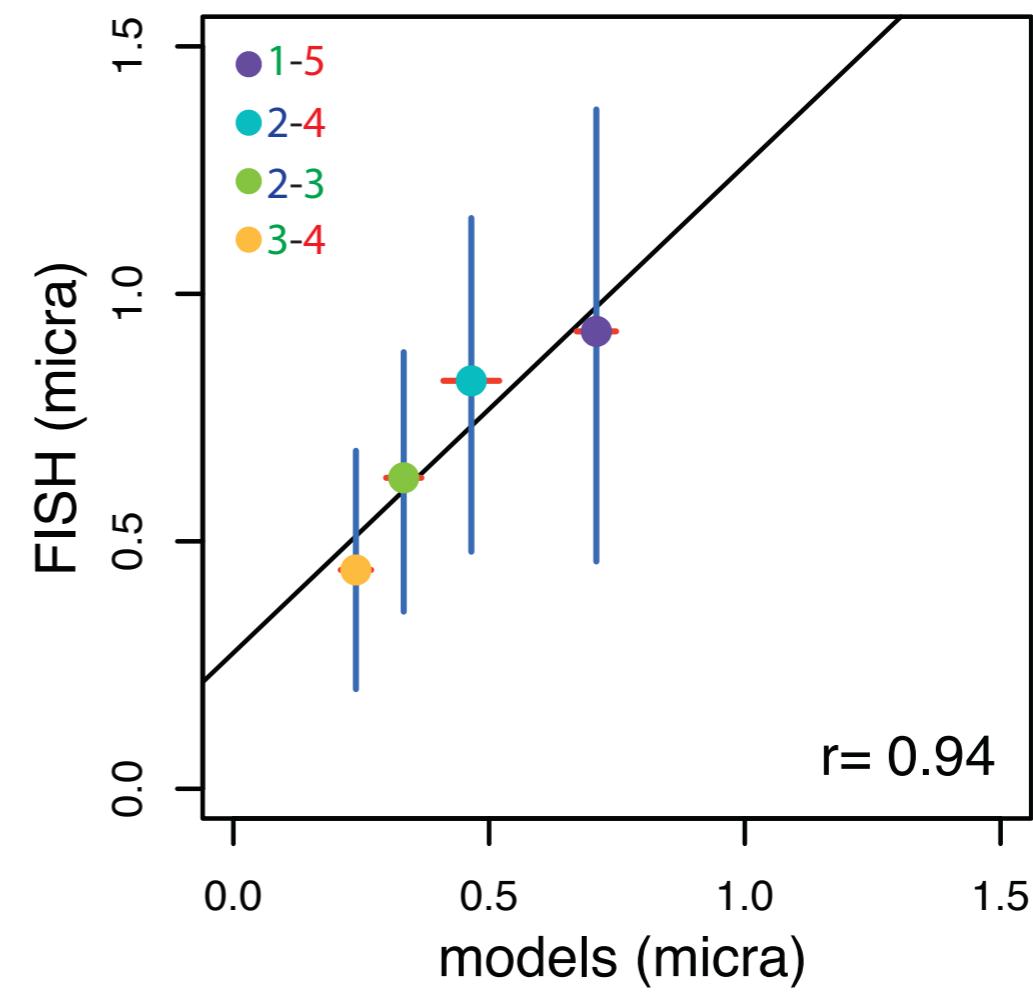
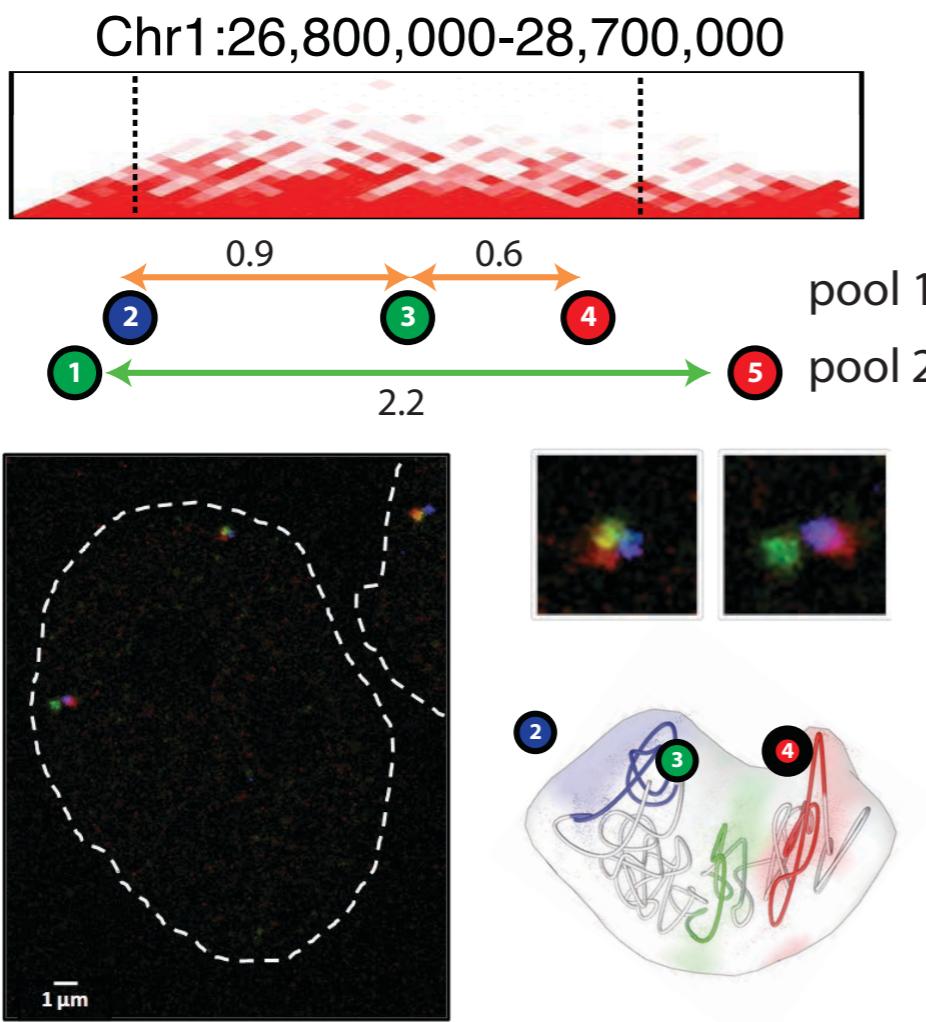
Pg induced fold change per TAD (6h)



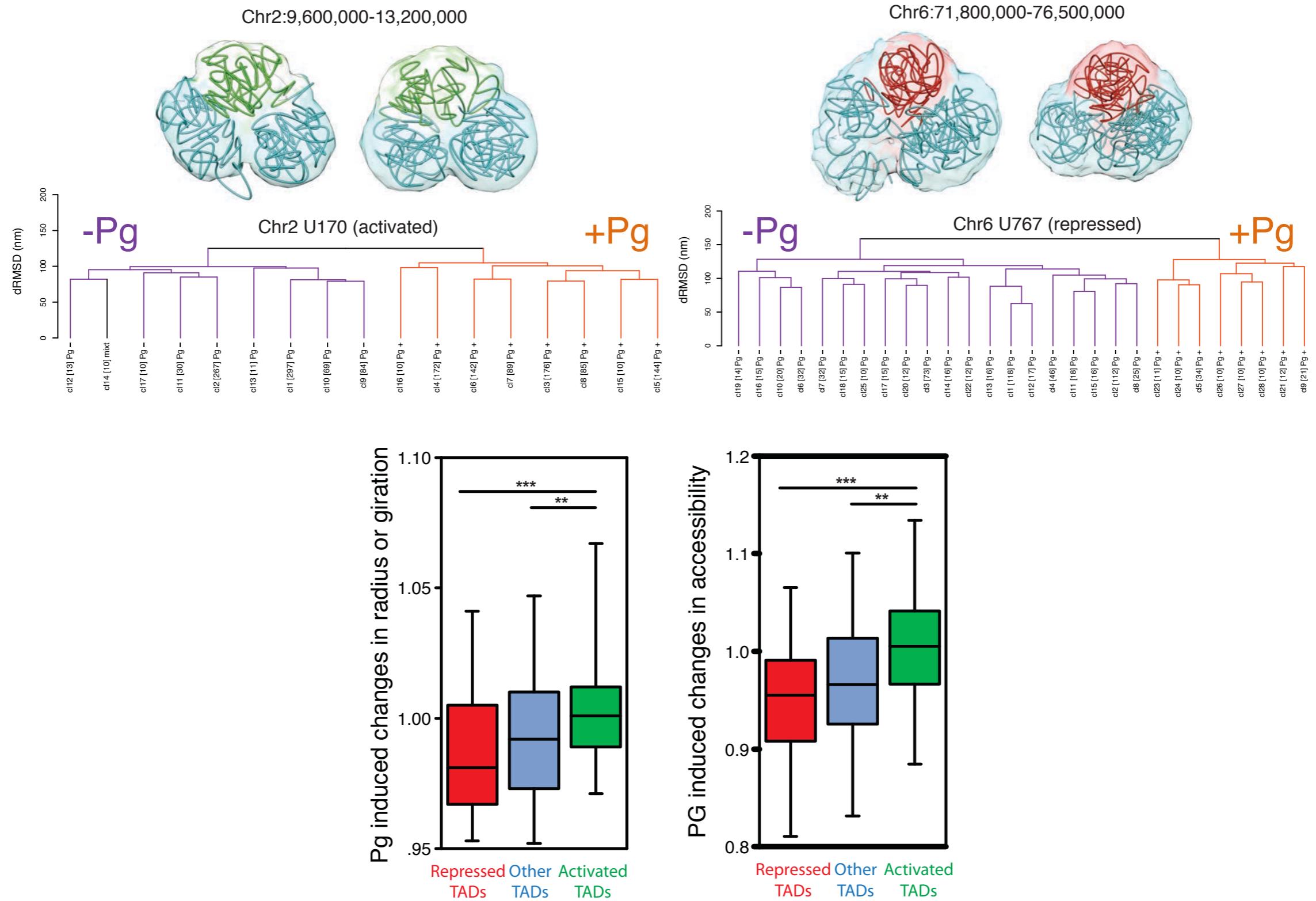
Modeling 3D TADs



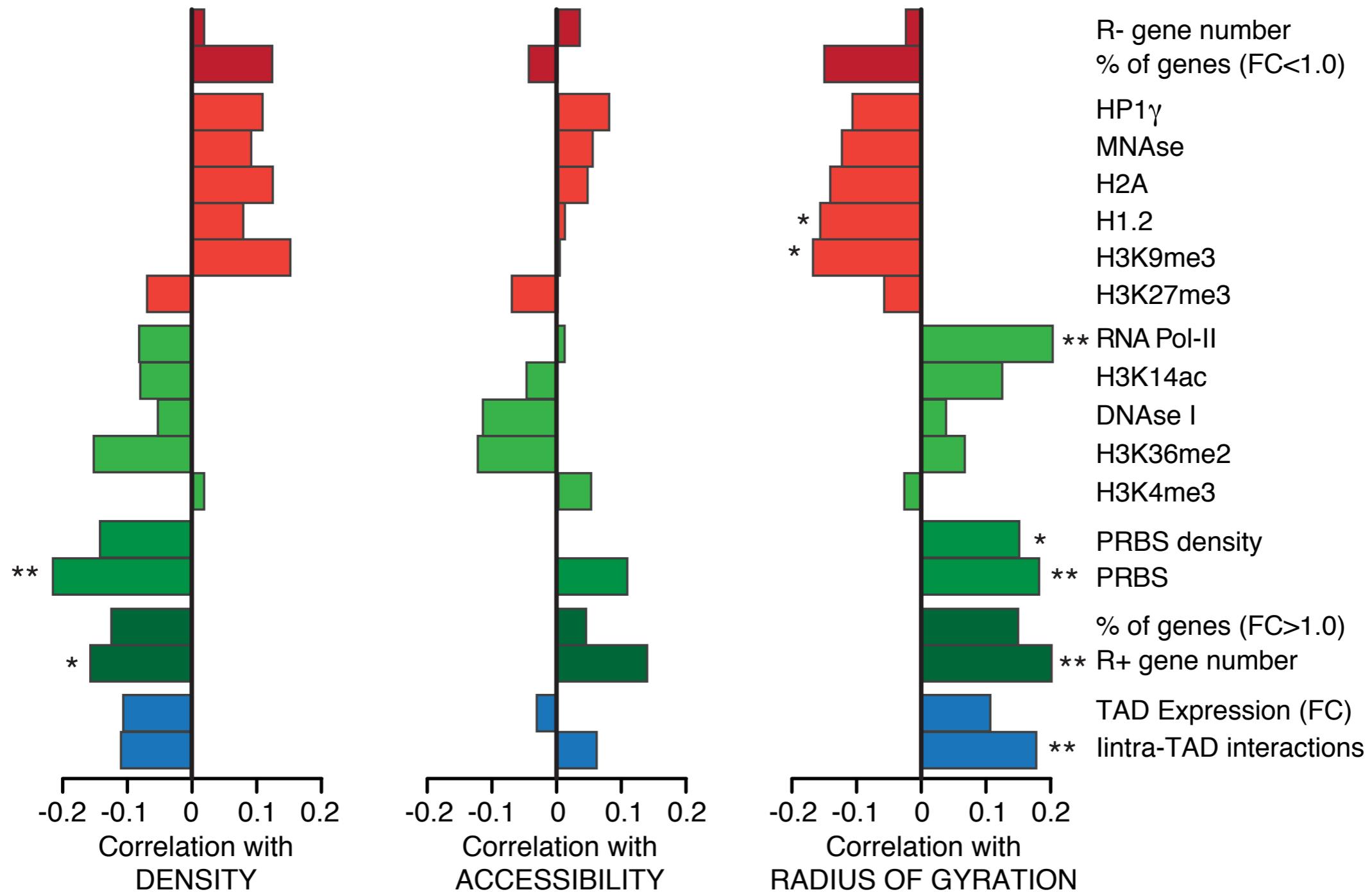
61 genomic regions containing 209 TADs covering 267Mb

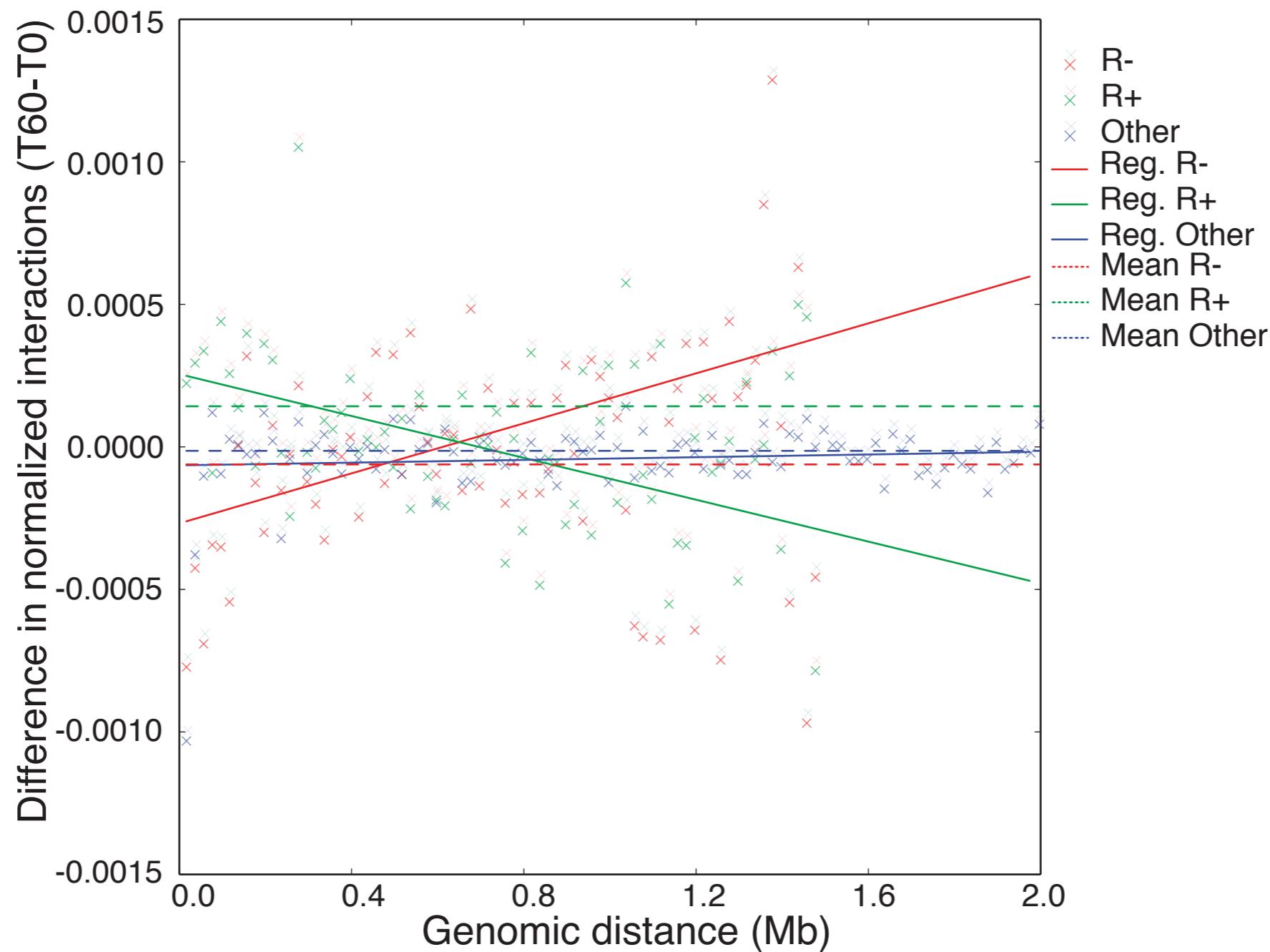
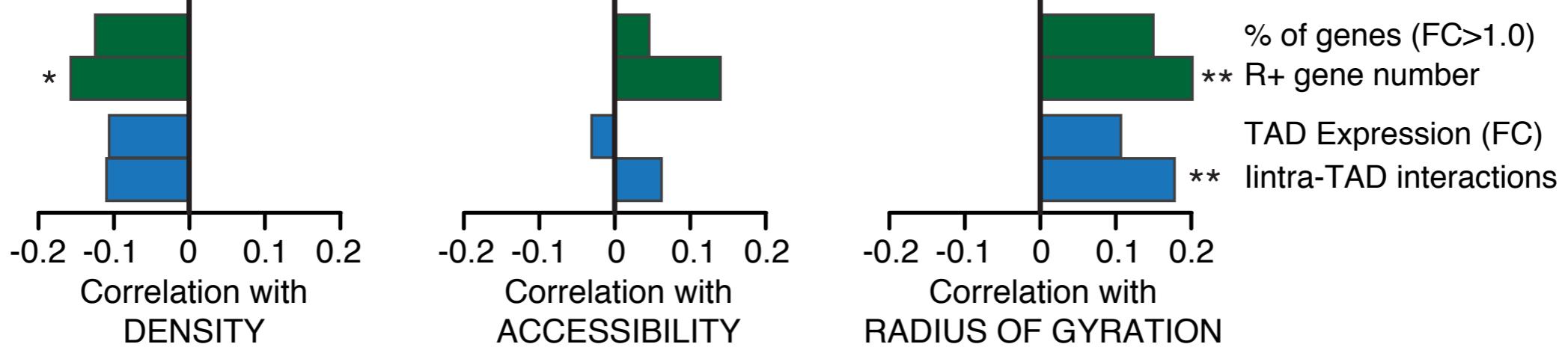


How TADs respond structurally to Pg?



How TADs respond structurally to Pg?

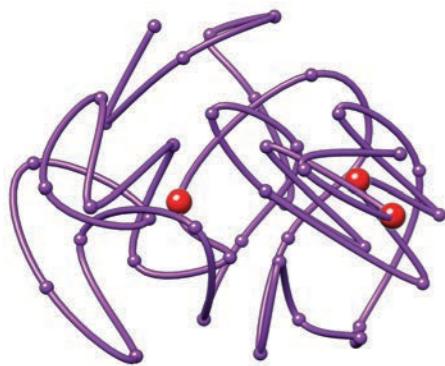




Model for TAD regulation

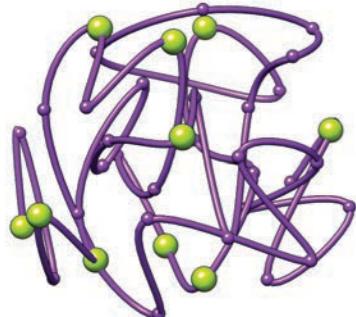
Repressed TAD

chr1 U41



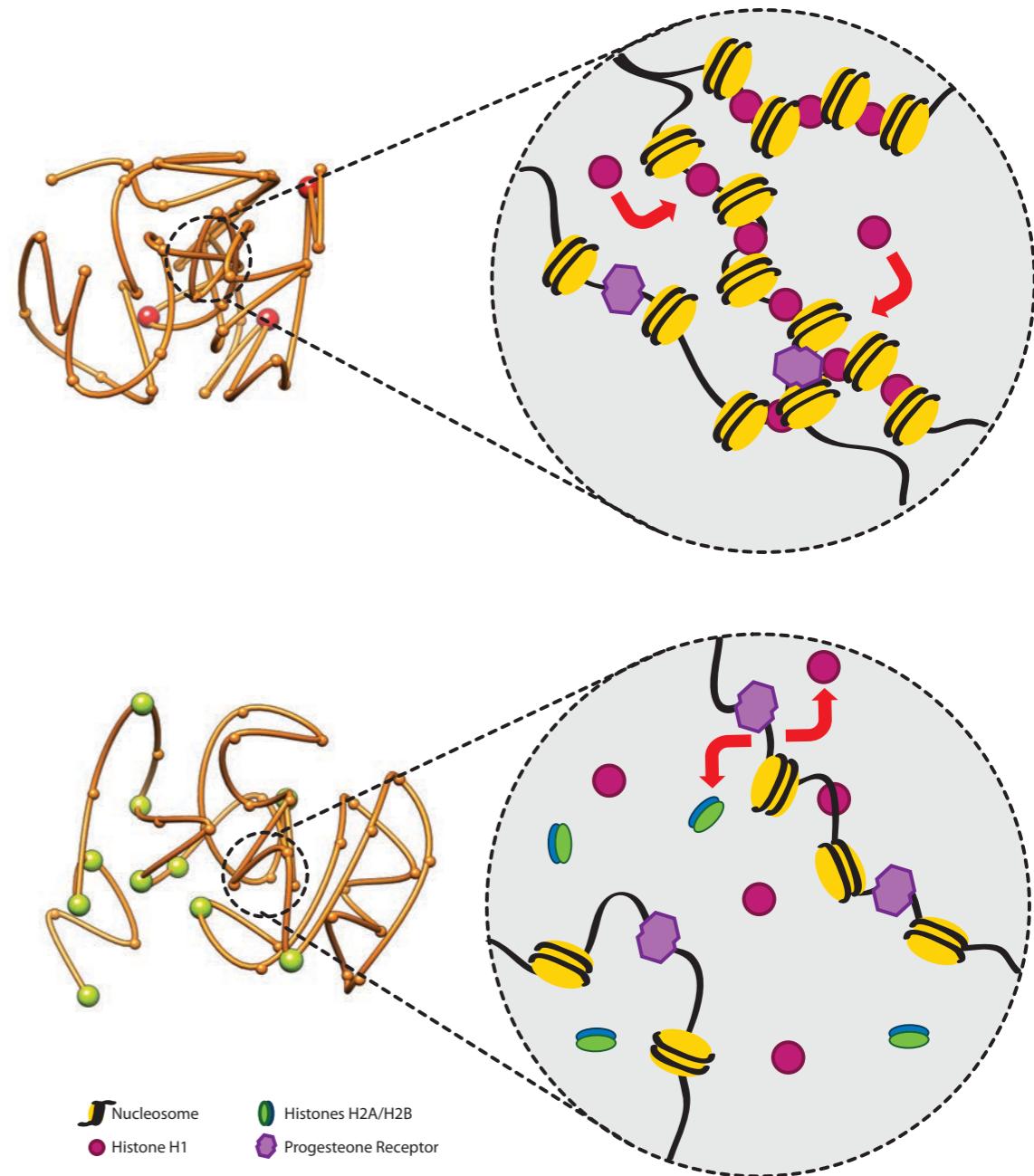
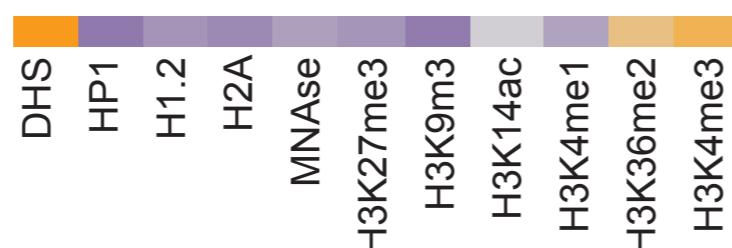
Activated TAD

chr2 U207

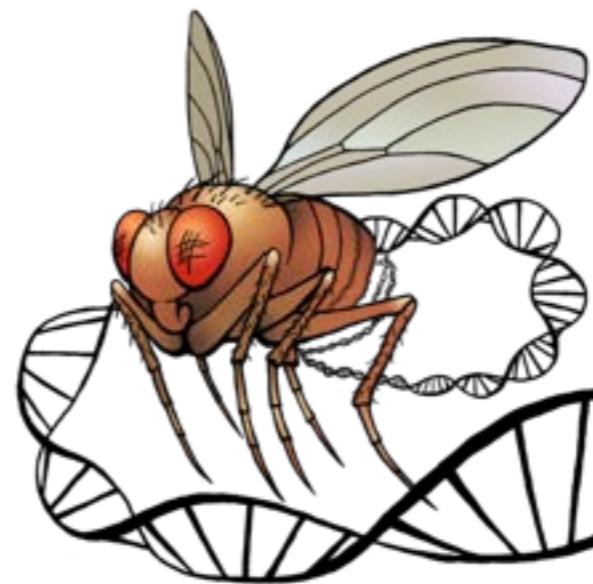


Structural transition

+Pg

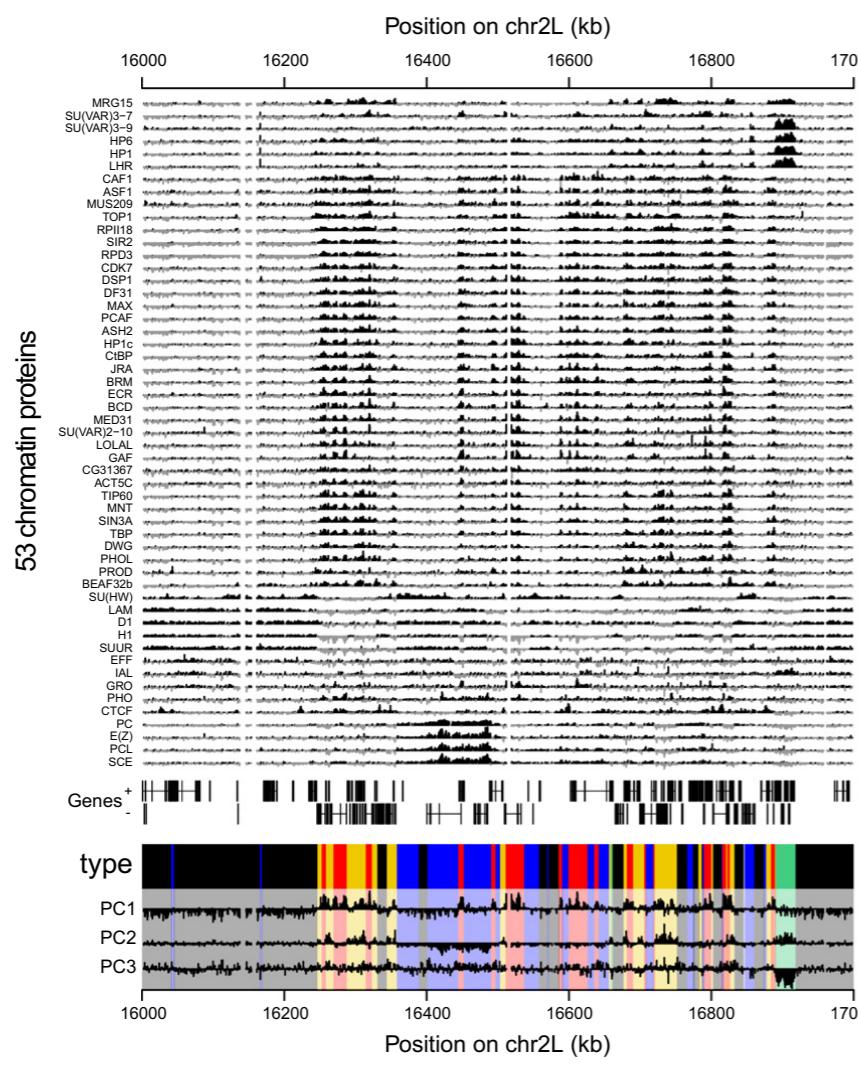
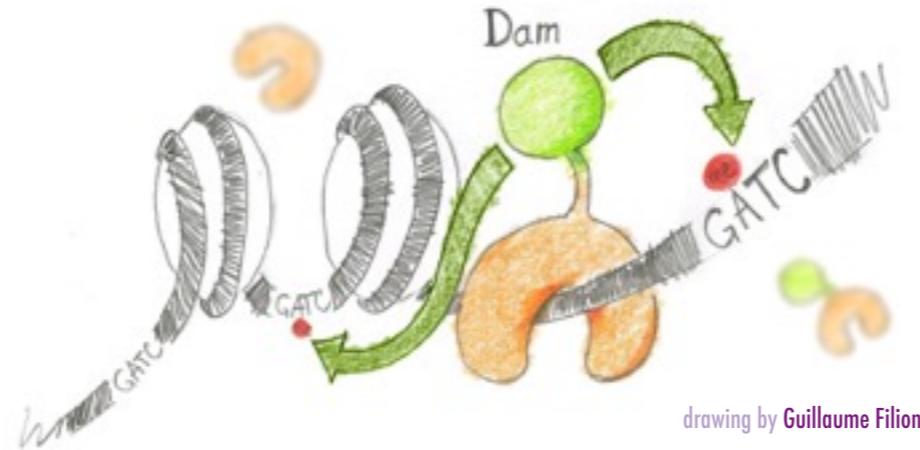


Structuring the **COLORs** of chromatin

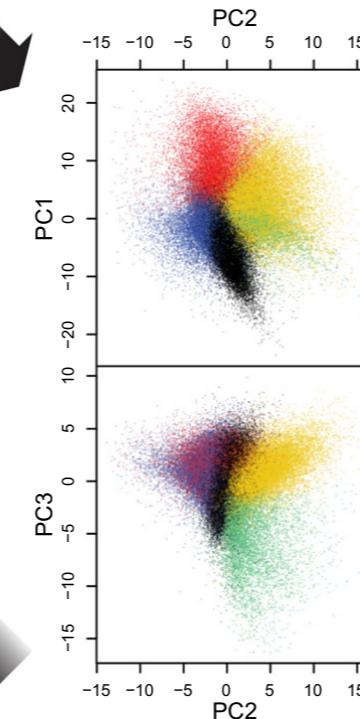


Fly Chromatin COLORs

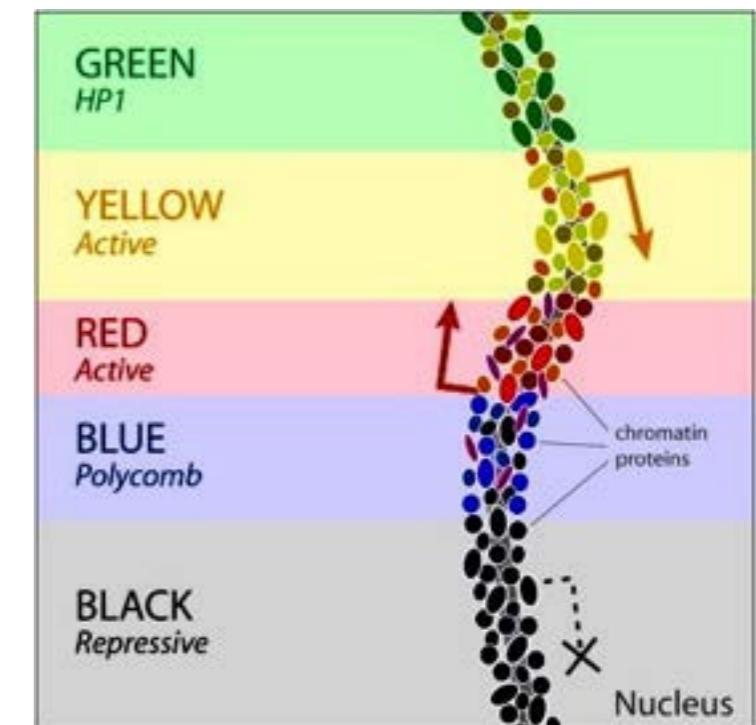
Filion et al. (2010). Cell, 143(2), 212–224.



Principal component analysis

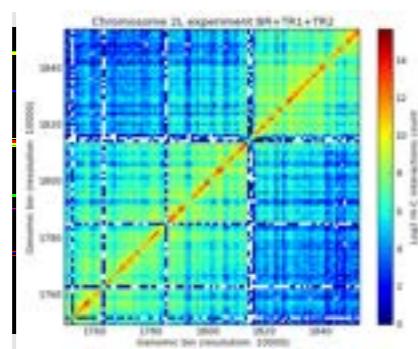
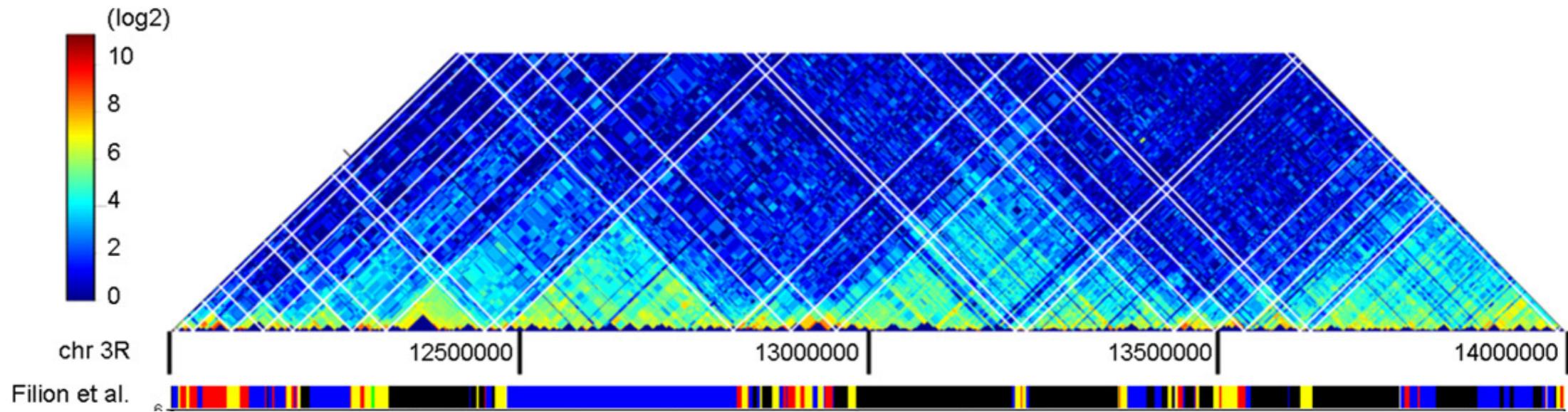


Hidden Markov model



Fly Chromatin COLORs

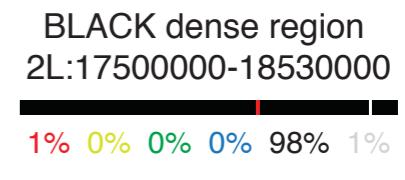
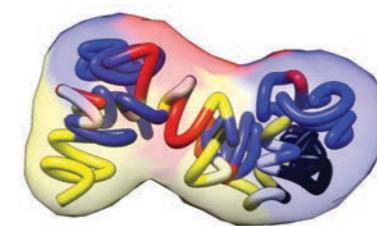
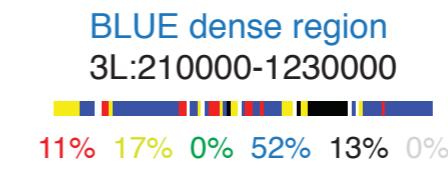
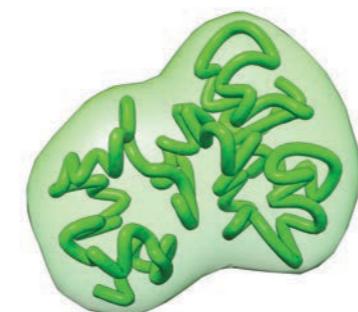
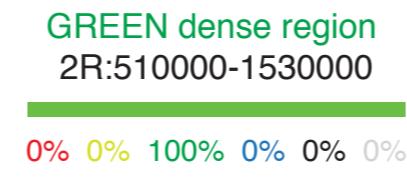
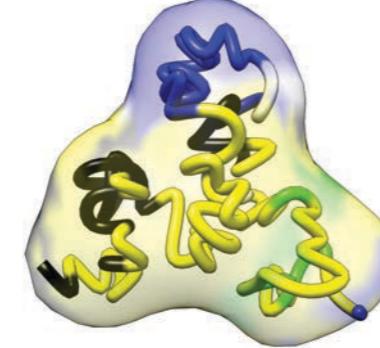
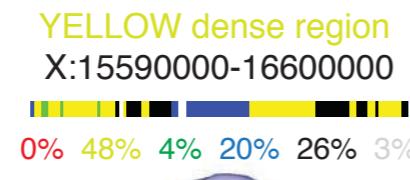
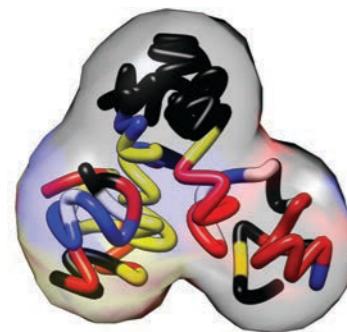
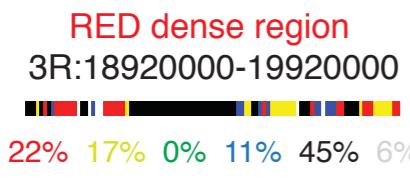
Hou et al. (2012). Molecular Cell, 48(3), 471–484.



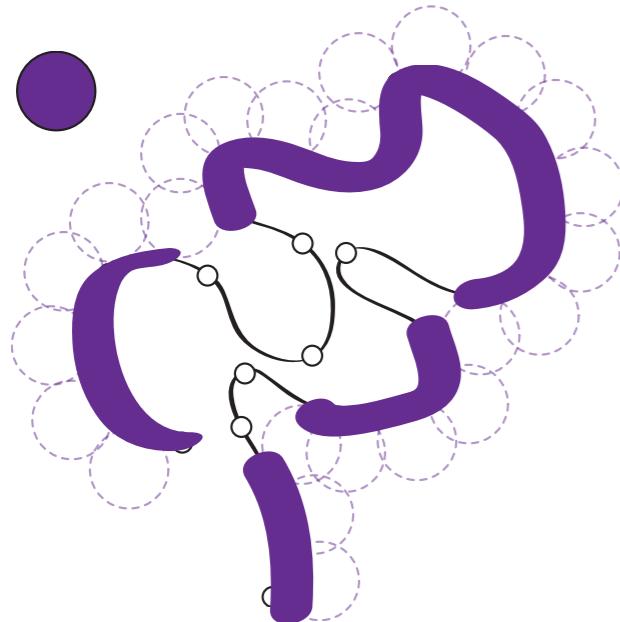
50 ~1Mb regions
10 for each color

Structural properties

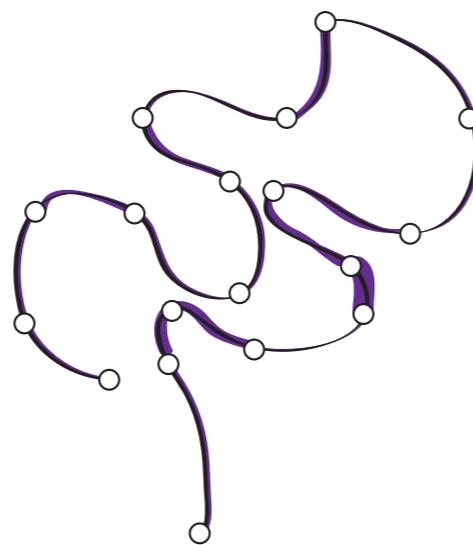
50 1Mb regions. 10 enriched for each color.



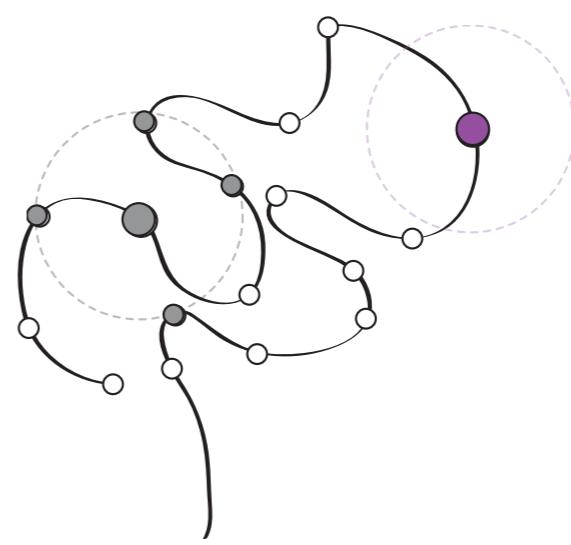
Accessibility (%)



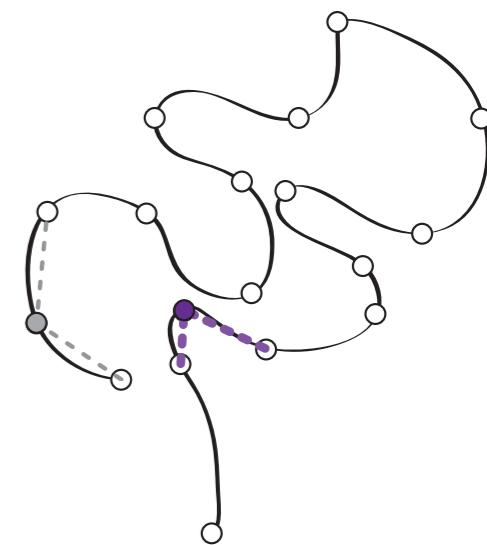
Density (bp/nm)



Interactions



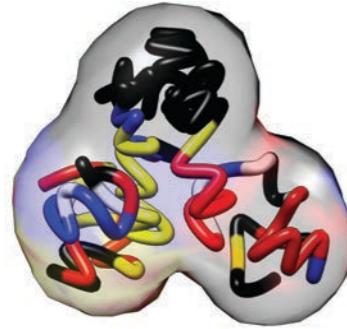
Angle



Structural COLORs

RED dense region
3R:18920000-19920000

 22% 17% 0% 11% 45% 6%



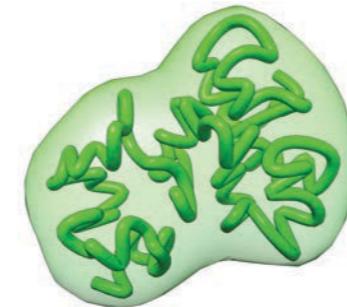
YELLOW dense region
X:15590000-16600000

 0% 48% 4% 20% 26% 3%



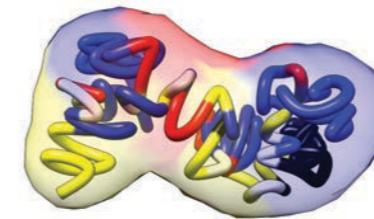
GREEN dense region
2R:510000-1530000

 0% 0% 100% 0% 0% 0%



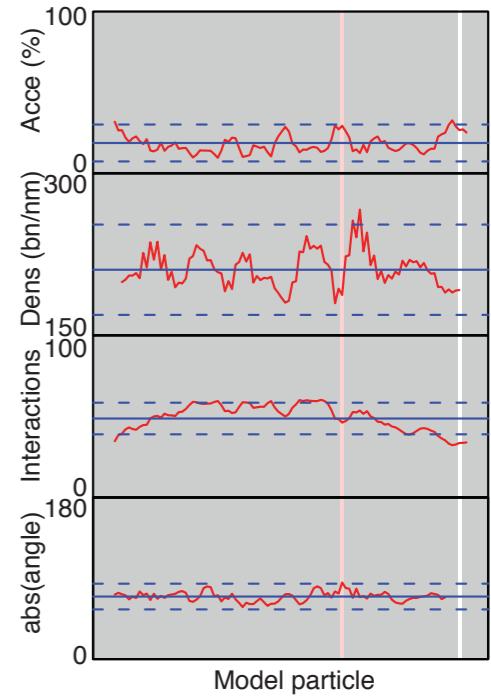
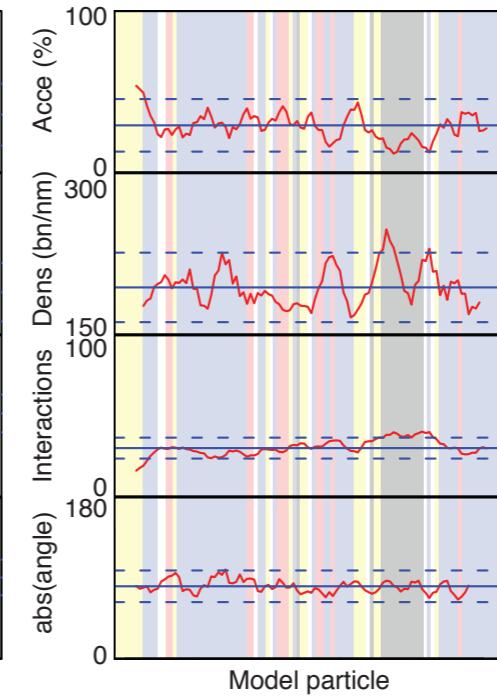
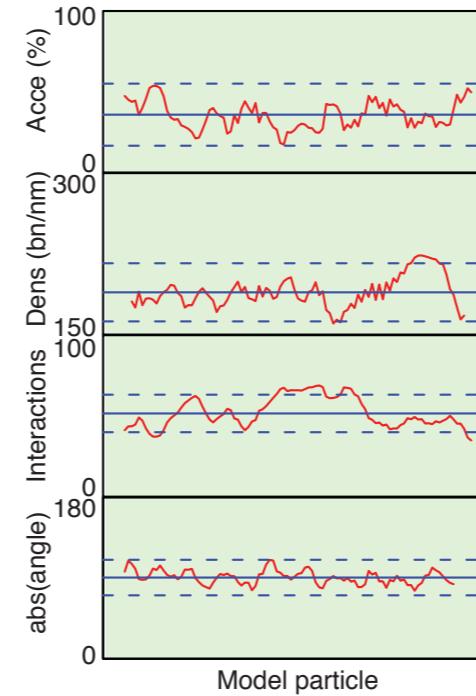
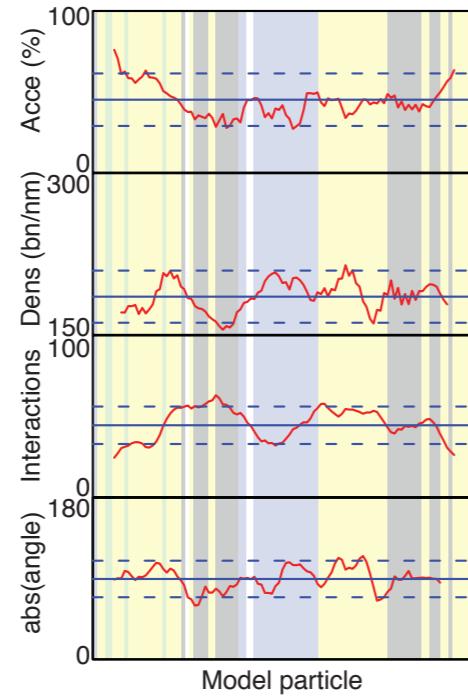
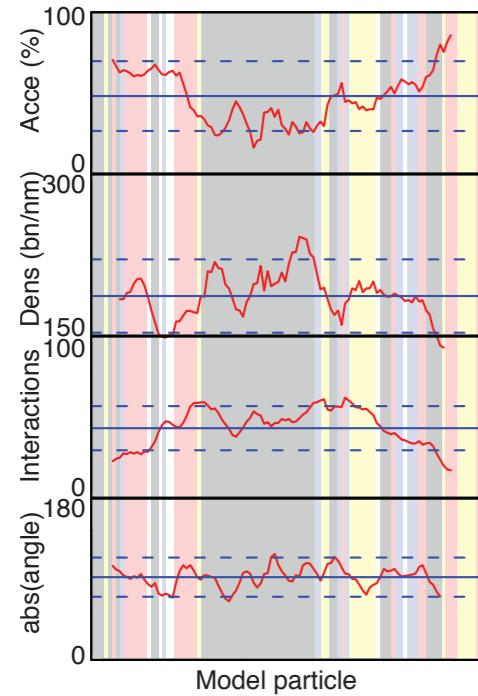
BLUE dense region
3L:210000-1230000

 11% 17% 0% 52% 13% 0%

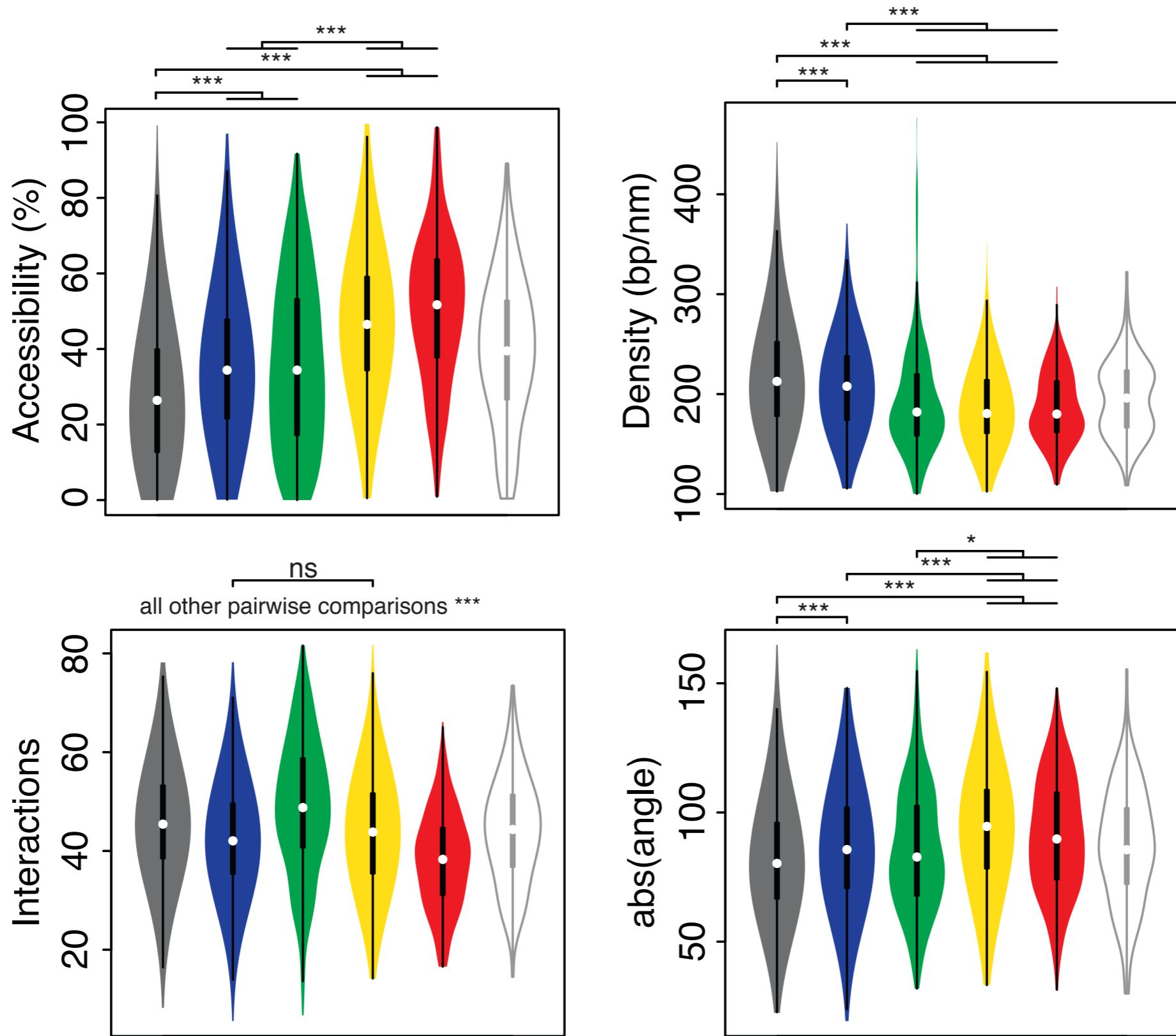


BLACK dense region
2L:17500000-18530000

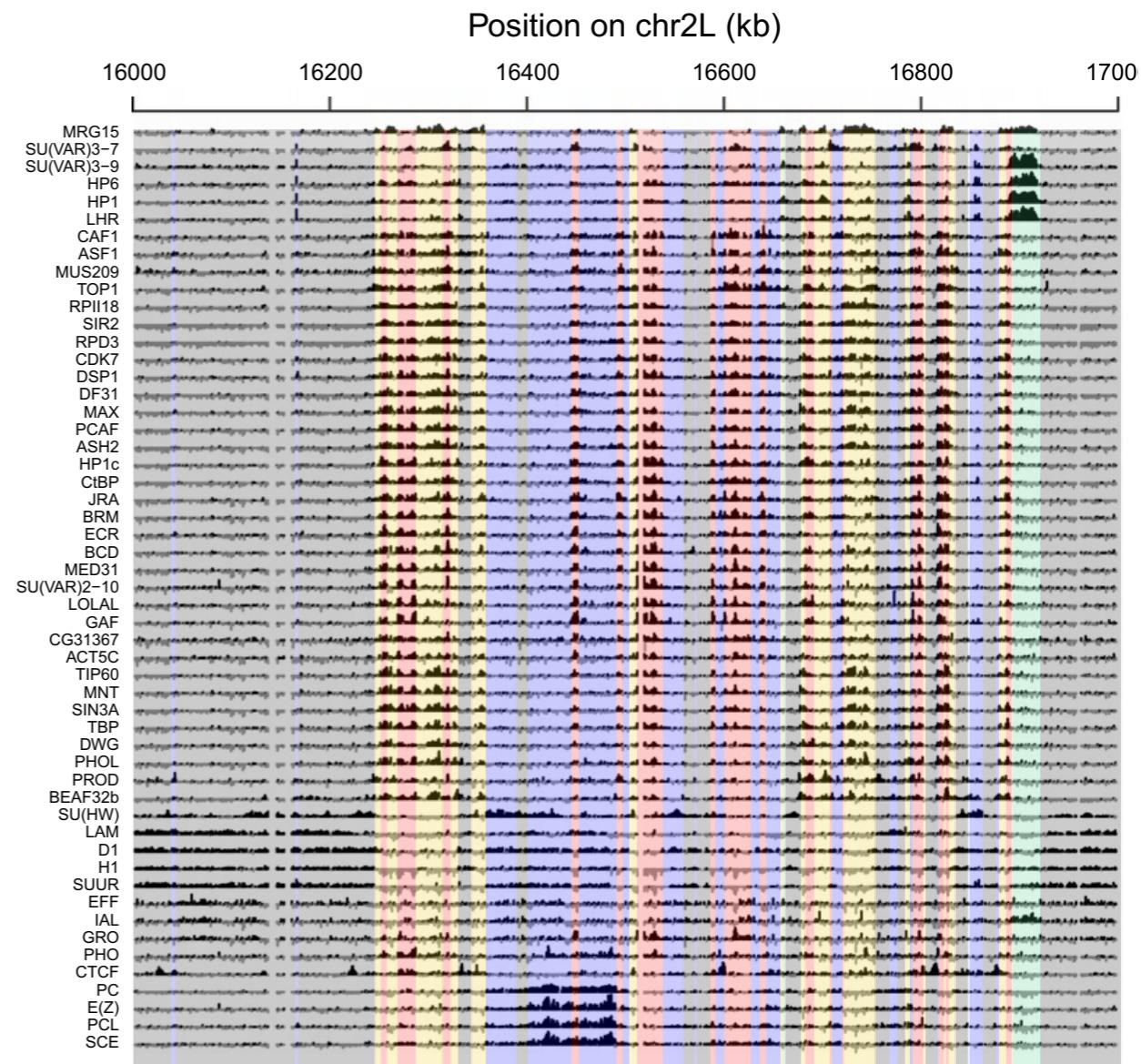
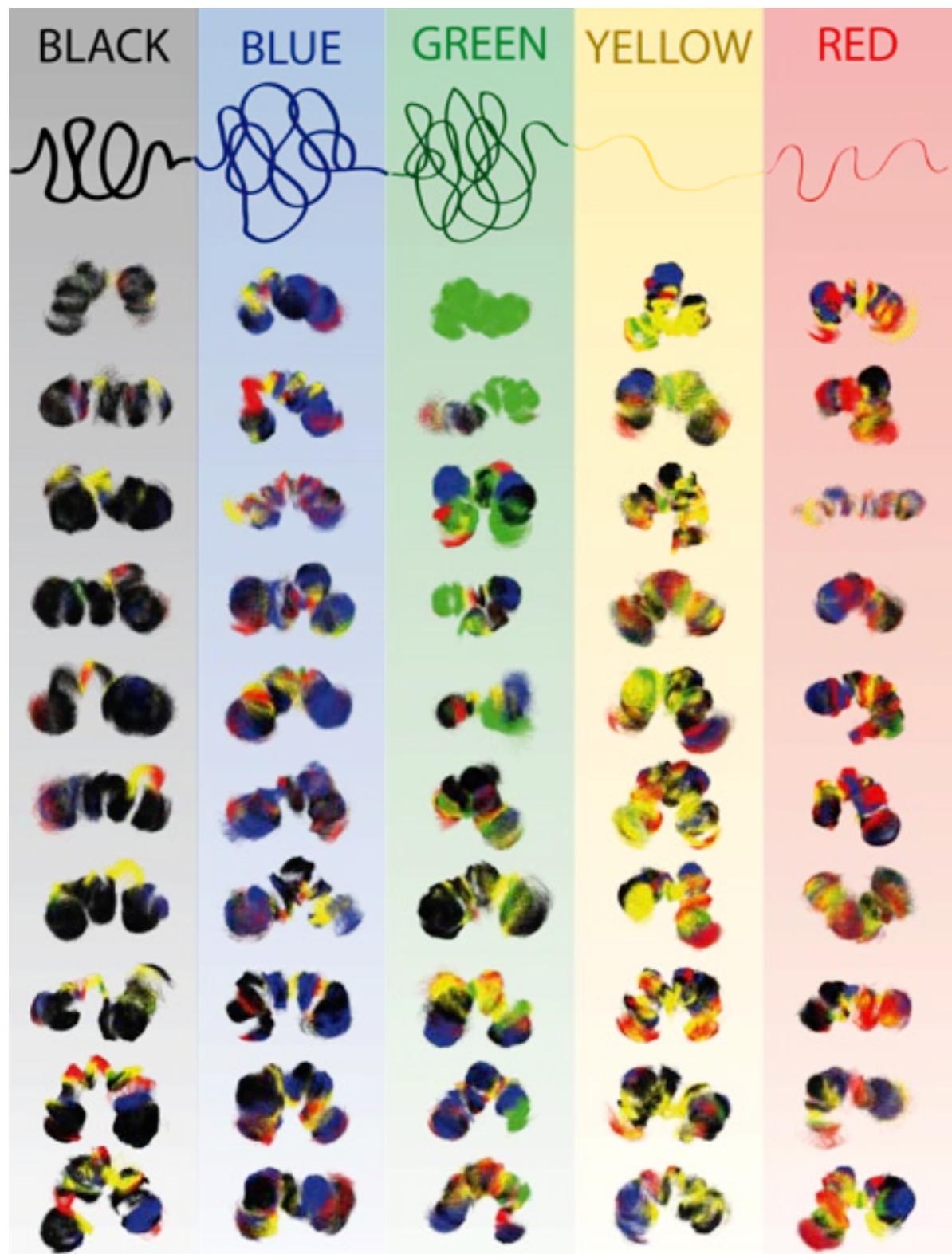
 1% 0% 0% 0% 98% 1%



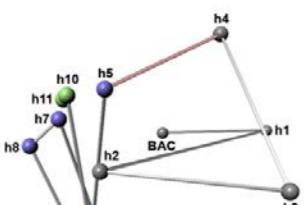
Structural COLORs



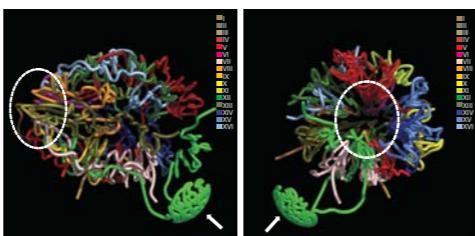
Structural COLORs



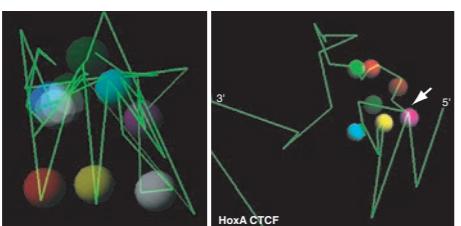
Are the models correct?



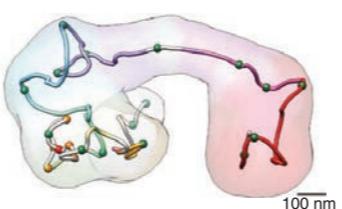
Jhunjhunwala (2008) Cell



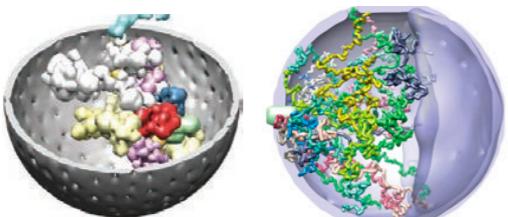
Duan (2010) Nature



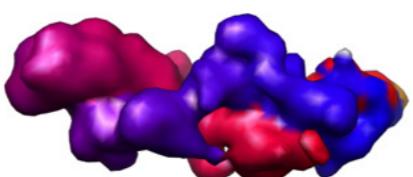
Fraser (2009) Genome Biology
Ferraiuolo (2010) Nucleic Acids Research



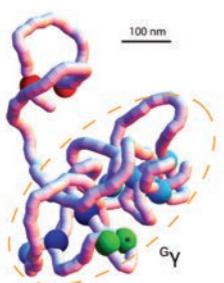
Bàu (2011) Nature Structural & Molecular Biology



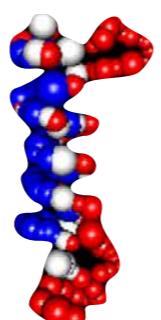
Kalhor (2011) Nature Biotechnology
Tjong (2012) Genome Research



Umbarger (2011) Molecular Cell



Junier (2012) Nucleic Acids Research



Hu (2013) PLoS Computational Biology

Nucleic Acids Research Advance Access published March 23, 2015

Nucleic Acids Research, 2015, 1
doi: 10.1093/nar/gkv221

Assessing the limits of restraint-based 3D modeling of genomes and genomic domains

Marie Trussart^{1,2}, François Serra^{3,4}, Davide Bàu^{3,4}, Ivan Junier^{2,3}, Luís Serrano^{1,2,5} and Marc A. Martí-Renom^{3,4,5,*}

¹EMBL/CRG Systems Biology Research Unit, Centre for Genomic Regulation (CRG), Barcelona, Spain, ²Universitat Pompeu Fabra (UPF), Barcelona, Spain, ³Gene Regulation, Stem Cells and Cancer Program, Centre for Genomic Regulation (CRG), Barcelona, Spain, ⁴Genome Biology Group, Centre Nacional d'Anàlisi Genòmica (CNAG), Barcelona, Spain and ⁵Institució Catalana de Recerca i Estudis Avançats (ICREA), Barcelona, Spain

Received January 16, 2015; Revised February 16, 2015; Accepted February 22, 2015

ABSTRACT

Restraint-based modeling of genomes has been recently explored with the advent of Chromosome Conformation Capture (3C-based) experiments. We previously developed a reconstruction method to resolve the 3D architecture of both prokaryotic and eukaryotic genomes using 3C-based data. These models were congruent with fluorescent imaging validation. However, the limits of such methods have not systematically been assessed. Here we propose the first evaluation of a mean-field restraint-based reconstruction of genomes by considering diverse chromosome architectures and different levels of data noise and structural variability. The results show that: first, current scoring functions for 3D reconstruction correlate with the accuracy of the models; second, reconstructed models are robust to noise but sensitive to structural variability; third, the local structure organization of genomes, such as Topologically Associating Domains, results in more accurate models; fourth, to a certain extent, the models capture the intrinsic structural variability in the input matrices and fifth, the accuracy of the models can be *a priori* predicted by analyzing the properties of the interaction matrices. In summary, our work provides a systematic analysis of the limitations of a mean-field restraint-based method, which could be taken into consideration in further development of methods as well as their applications.

INTRODUCTION

Recent studies of the three-dimensional (3D) conformation of genomes are revealing insights into the organization and the regulation of biological processes, such as gene

expression regulation and replication (1–6). The advent of the so-called Chromosome Conformation Capture (3C) assays (7), which allowed identifying chromatin-looping interactions between pairs of loci, helped deciphering some of the key elements organizing the genomes. High-throughput derivations of genome-wide 3C-based assays were established with Hi-C technologies (8) for an unbiased identification of chromatin interactions. The resulting genome interaction matrices from Hi-C experiments have been extensively used for computationally analyzing the organization of genomes and genomic domains (5). In particular, a significant number of new approaches for modeling the 3D organization of genomes have recently flourished (9–14). The main goal of such approaches is to provide an accurate 3D representation of the bi-dimensional interaction matrices, which can then be more easily explored to extract biological insights. One type of methods for building 3D models from interaction matrices relies on the existence of a limited number of conformational states in the cell. Such methods are regarded as mean-field approaches and are able to capture, to a certain degree, the structural variability around these mean structures (15).

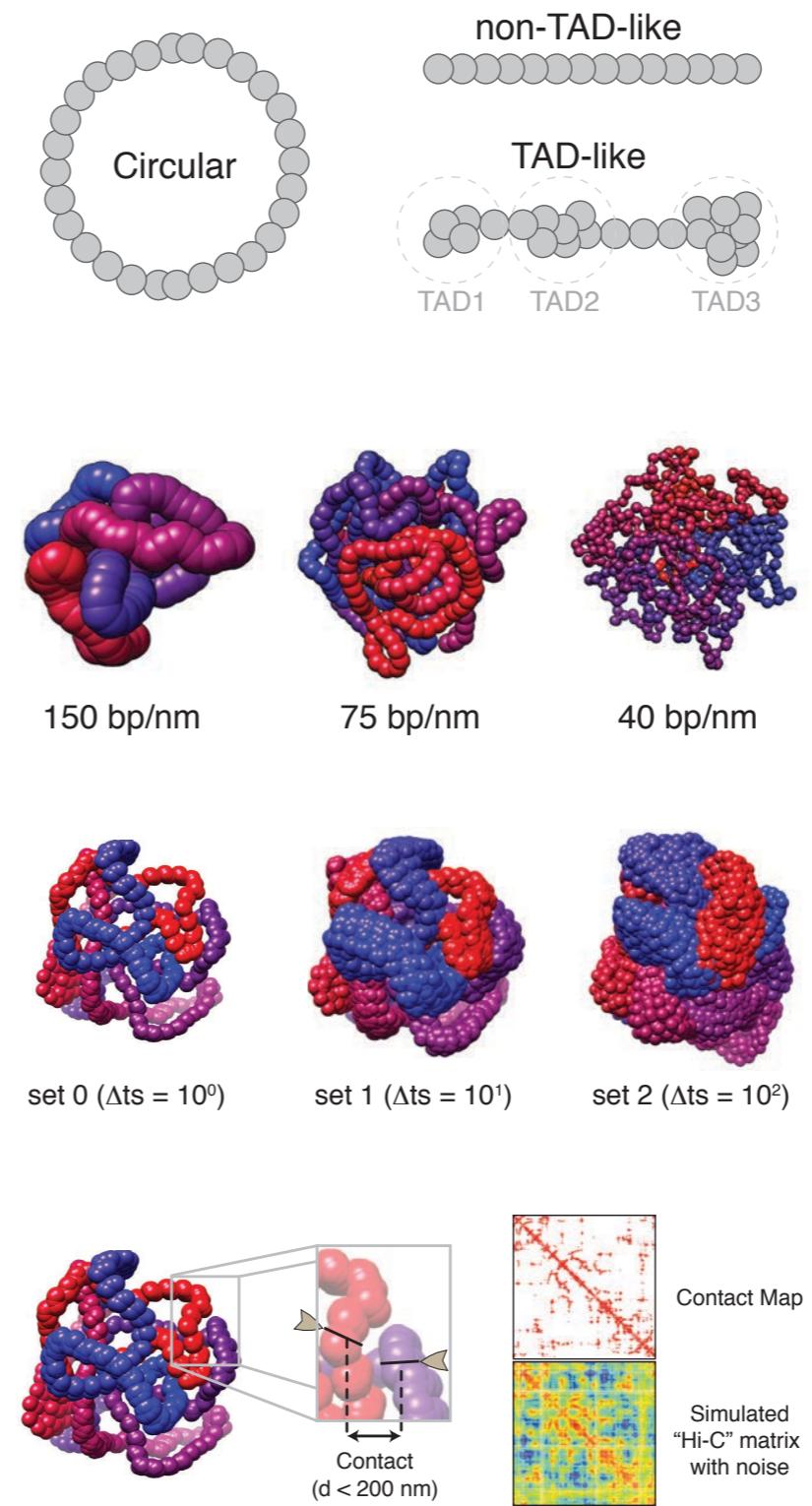
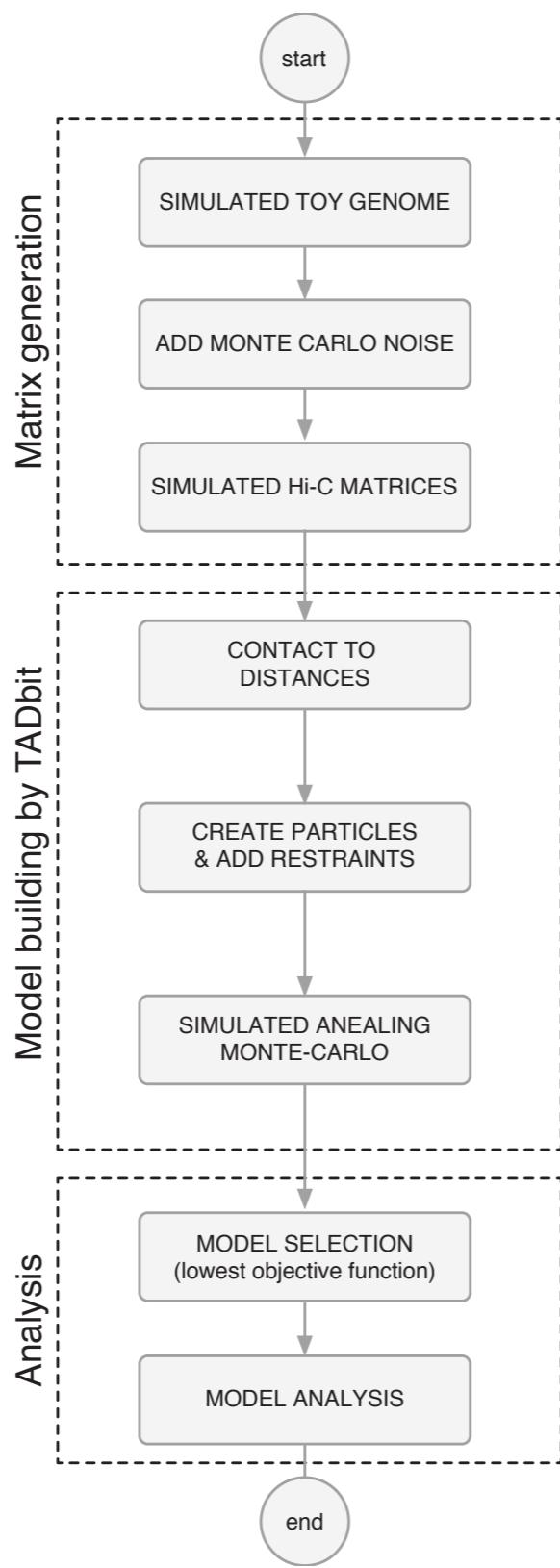
We recently developed a mean-field method for modeling 3D structures of genomes and genomic domains based on 3C interaction data (9). Our approach, called TADbit, was developed around the Integrative Modeling Platform (IMP, <http://integrativemodeling.org>), a general framework for restraint-based modeling of 3D bio-molecular structures (16). Briefly, our method uses chromatin interaction frequencies derived from experiments as a proxy of spatial proximity between the ligation products of the 3C libraries. Two fragments of DNA that interact with high frequency are dynamically placed close in space in our models while two fragments that do not interact as often will be kept apart. Our method has been successfully applied to model the structures of genomes and genomic domains in eukaryote and prokaryote organisms (17–19). In all of our studies, the final models were partially validated by assessing their

*To whom correspondence should be addressed. Tel: +34 934 020 542; Fax: +34 934 037 279; Email: mmarti@pcb.ub.cat

© The Author(s) 2015. Published by Oxford University Press on behalf of Nucleic Acids Research. This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

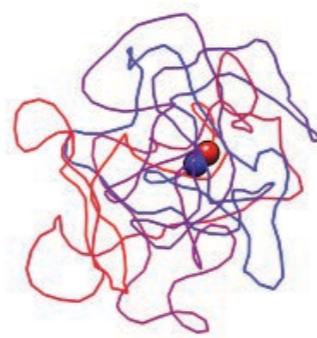
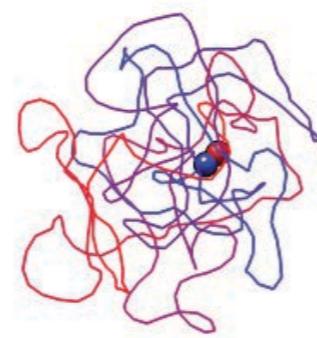
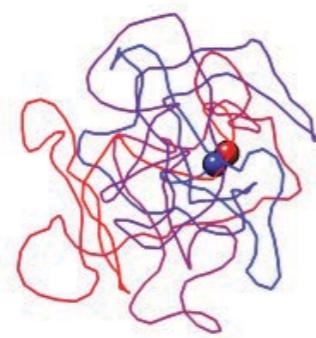
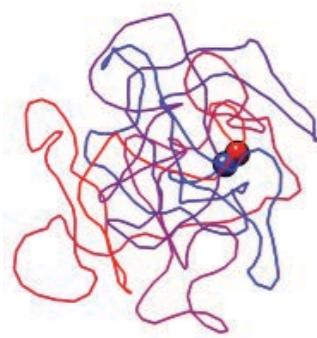
Trussart, et al. (2015). Nucleic Acids Research.

Toy models

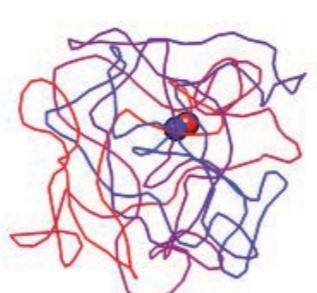
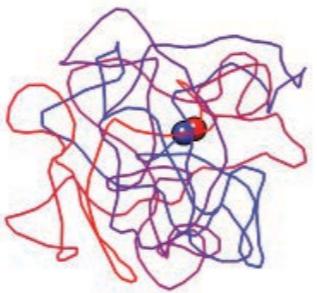
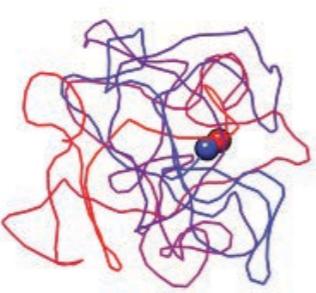
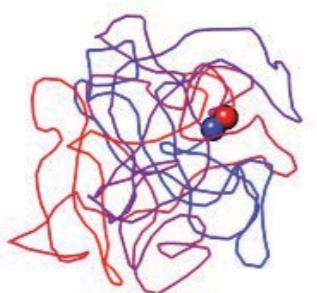
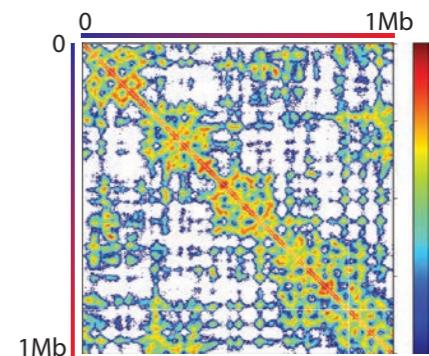


by Ivan Junier

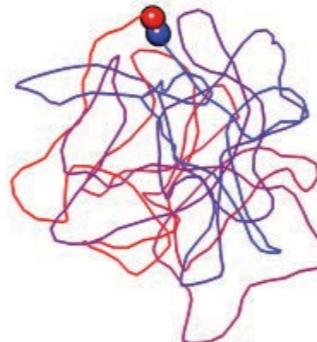
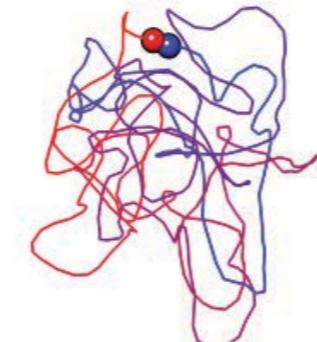
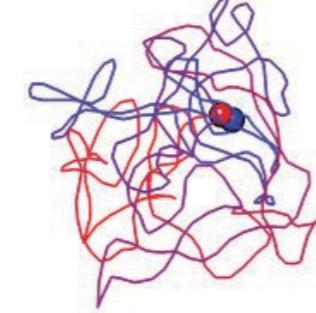
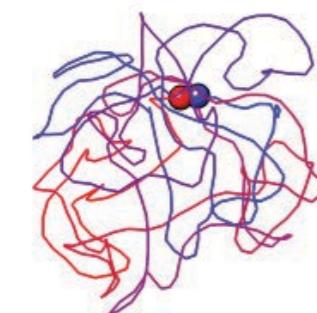
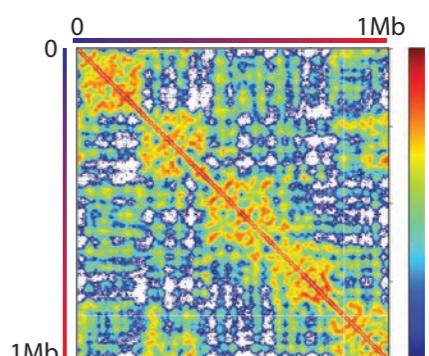
Toy interaction matrices



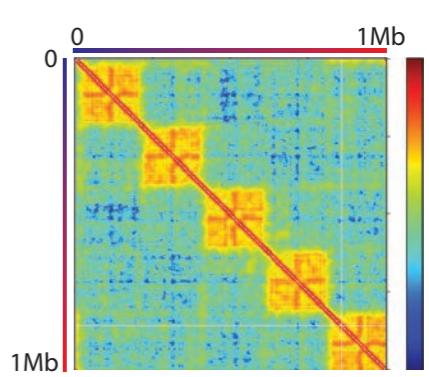
set 0 ($\Delta ts=10^0$)



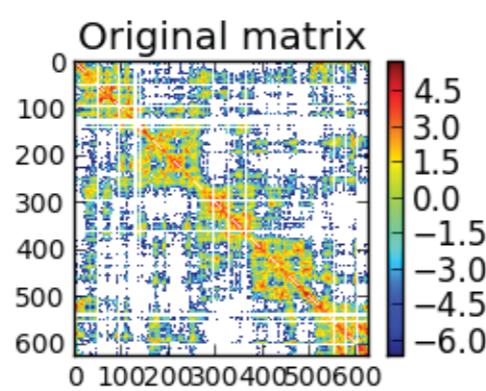
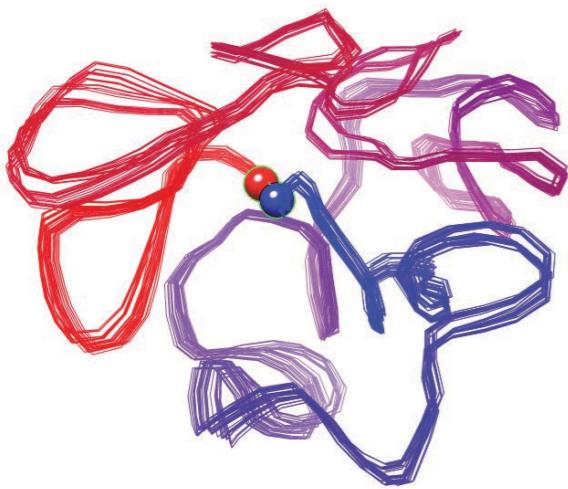
set 4 ($\Delta ts=10^4$)



set 6 ($\Delta ts=10^6$)



Reconstructing toy models



chr40_TAD

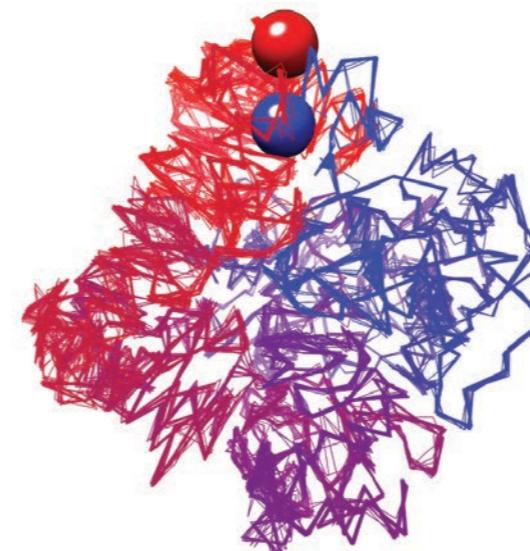
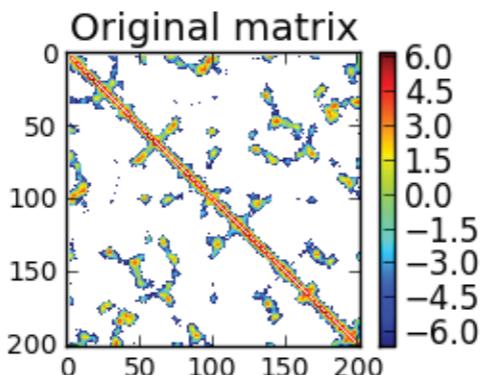
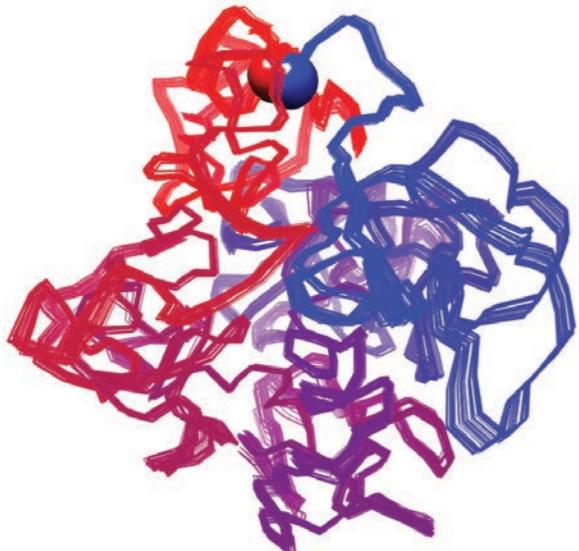
$\alpha=100$

$\Delta ts=10$

TADbit-SCC: 0.91

$\langle dRMSD \rangle$: 32.7 nm

$\langle dSCC \rangle$: 0.94



chr150_TAD

$\alpha=50$

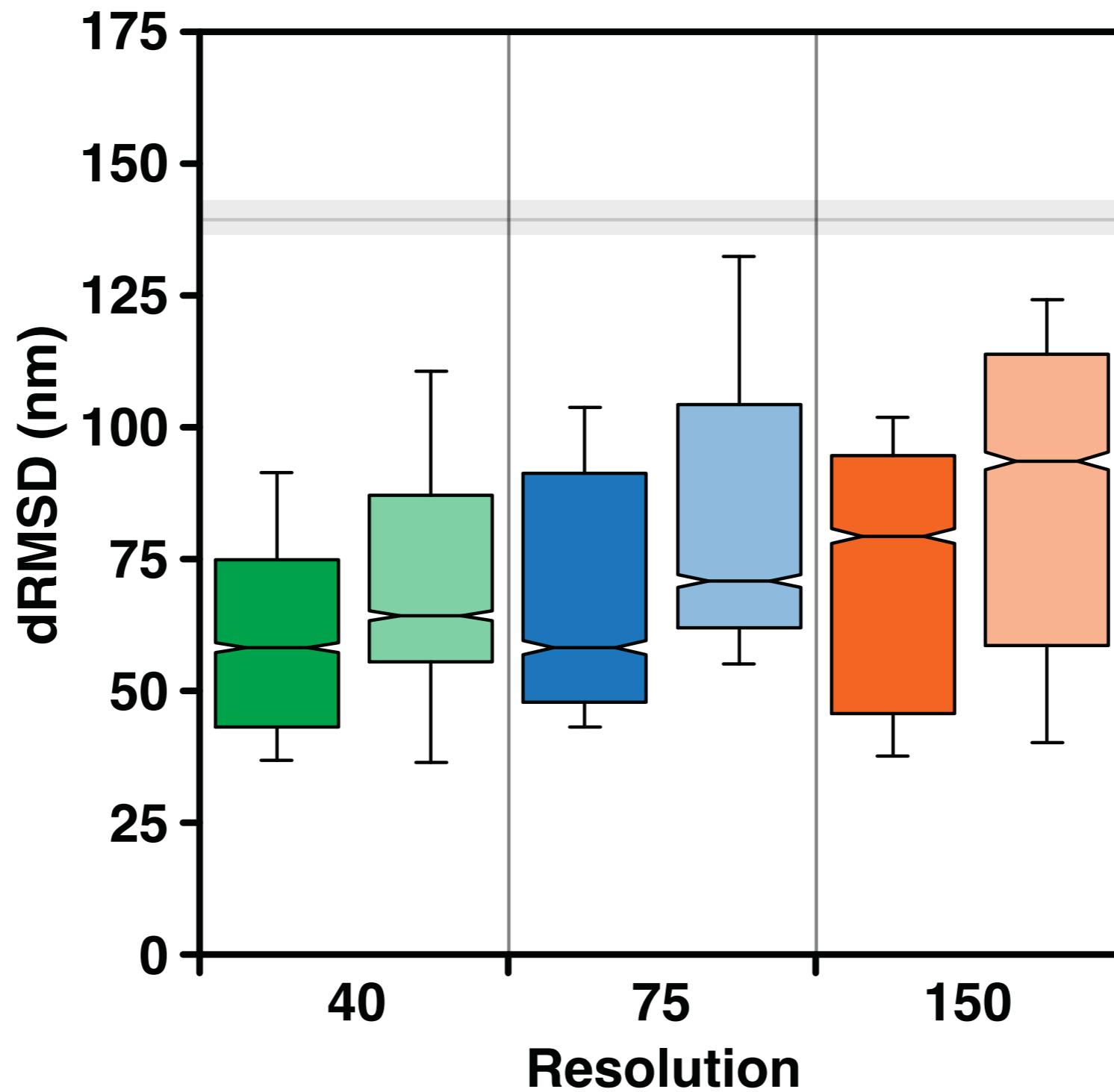
$\Delta ts=1$

TADbit-SCC: 0.82

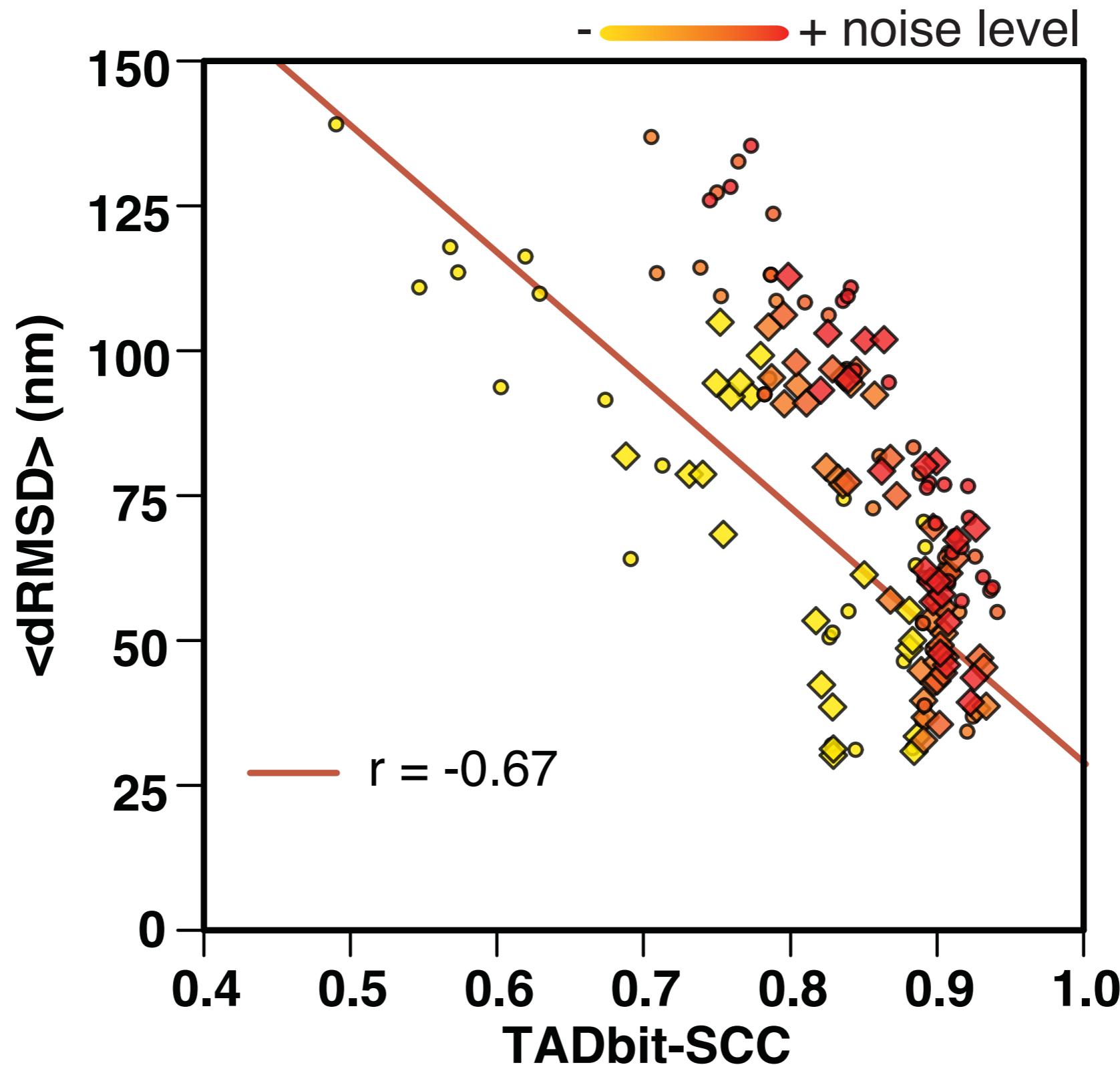
$\langle dRMSD \rangle$: 45.4 nm

$\langle dSCC \rangle$: 0.86

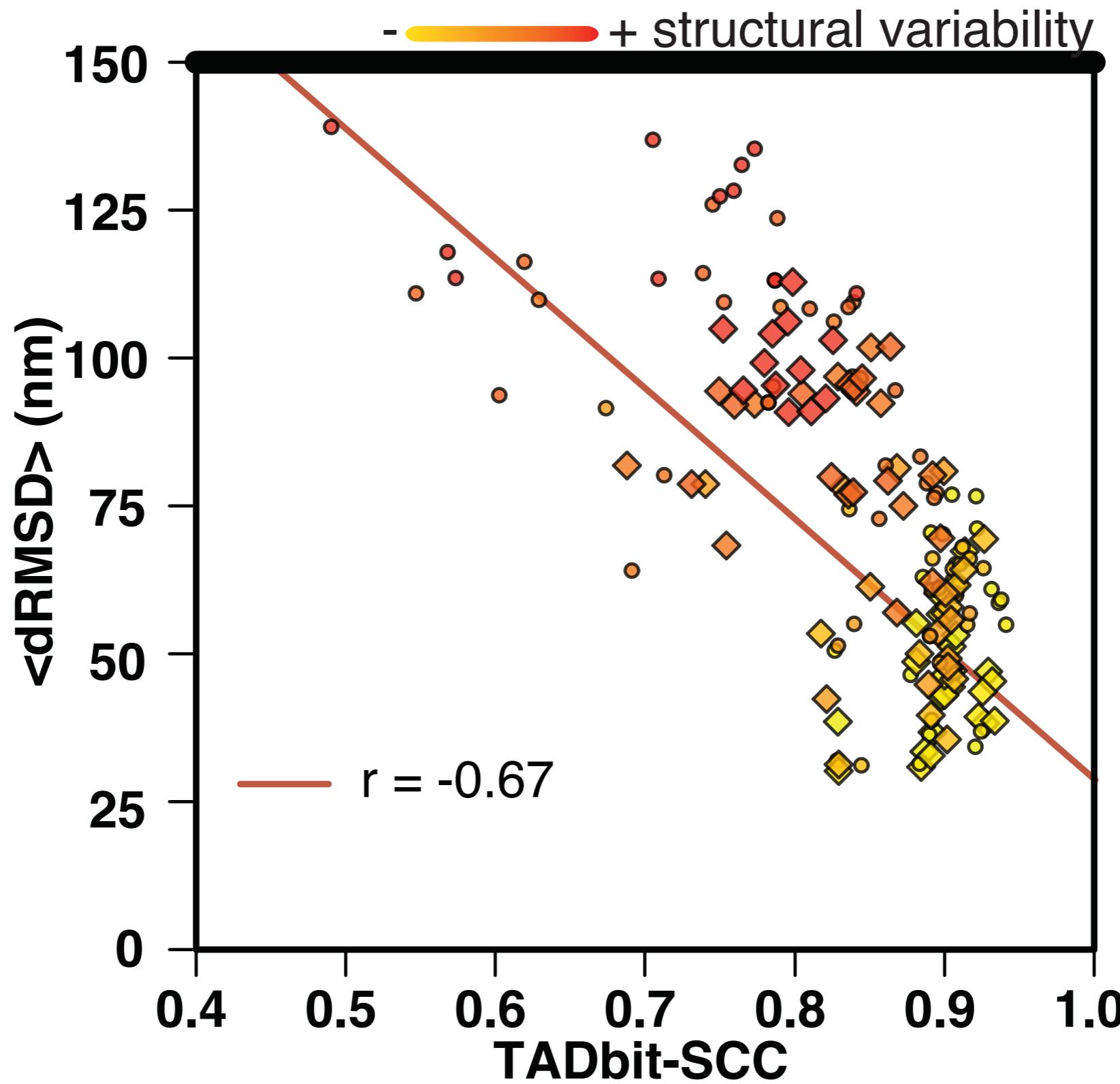
TADs & higher-res are "good"



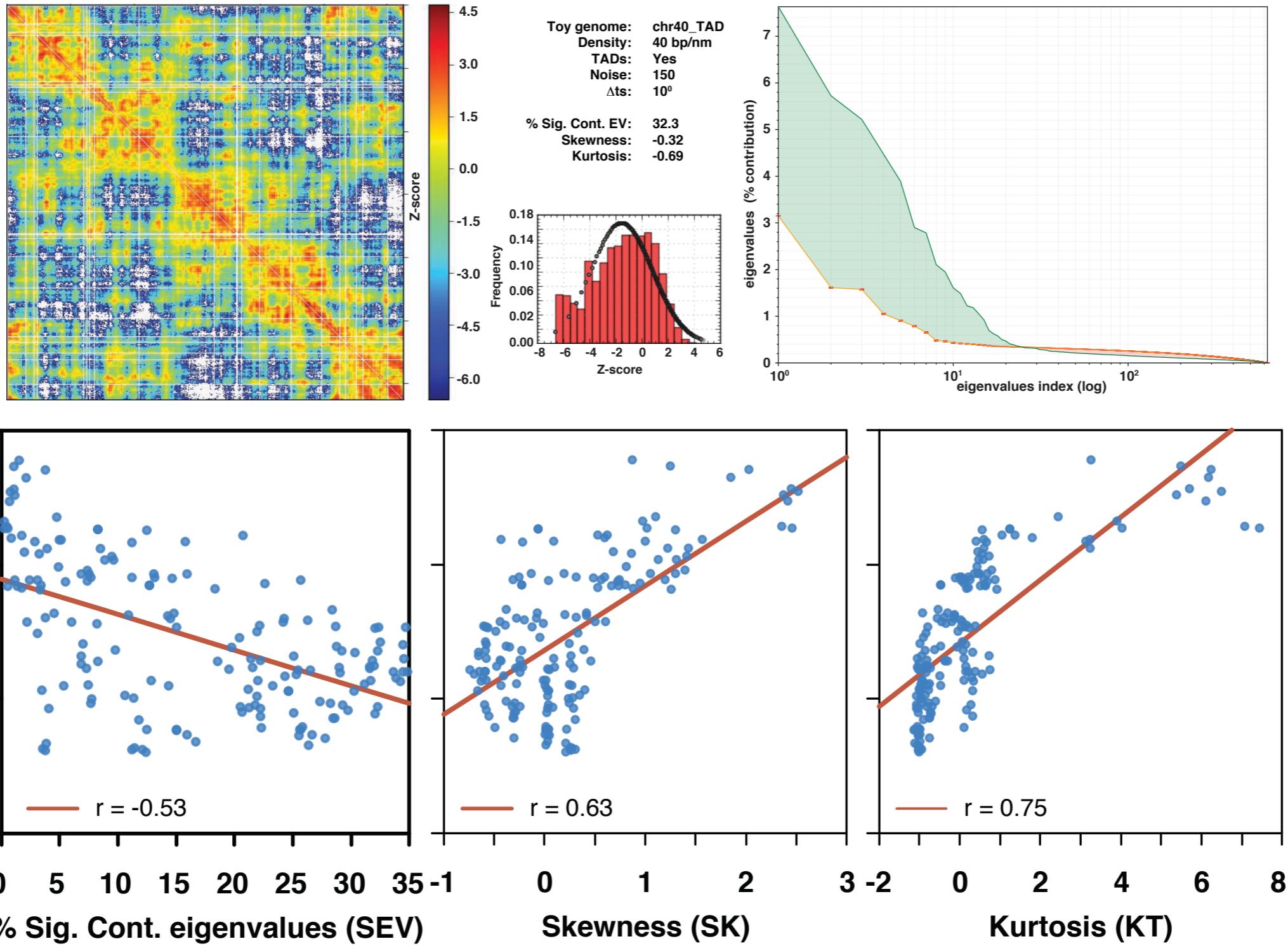
Noise is "OK"



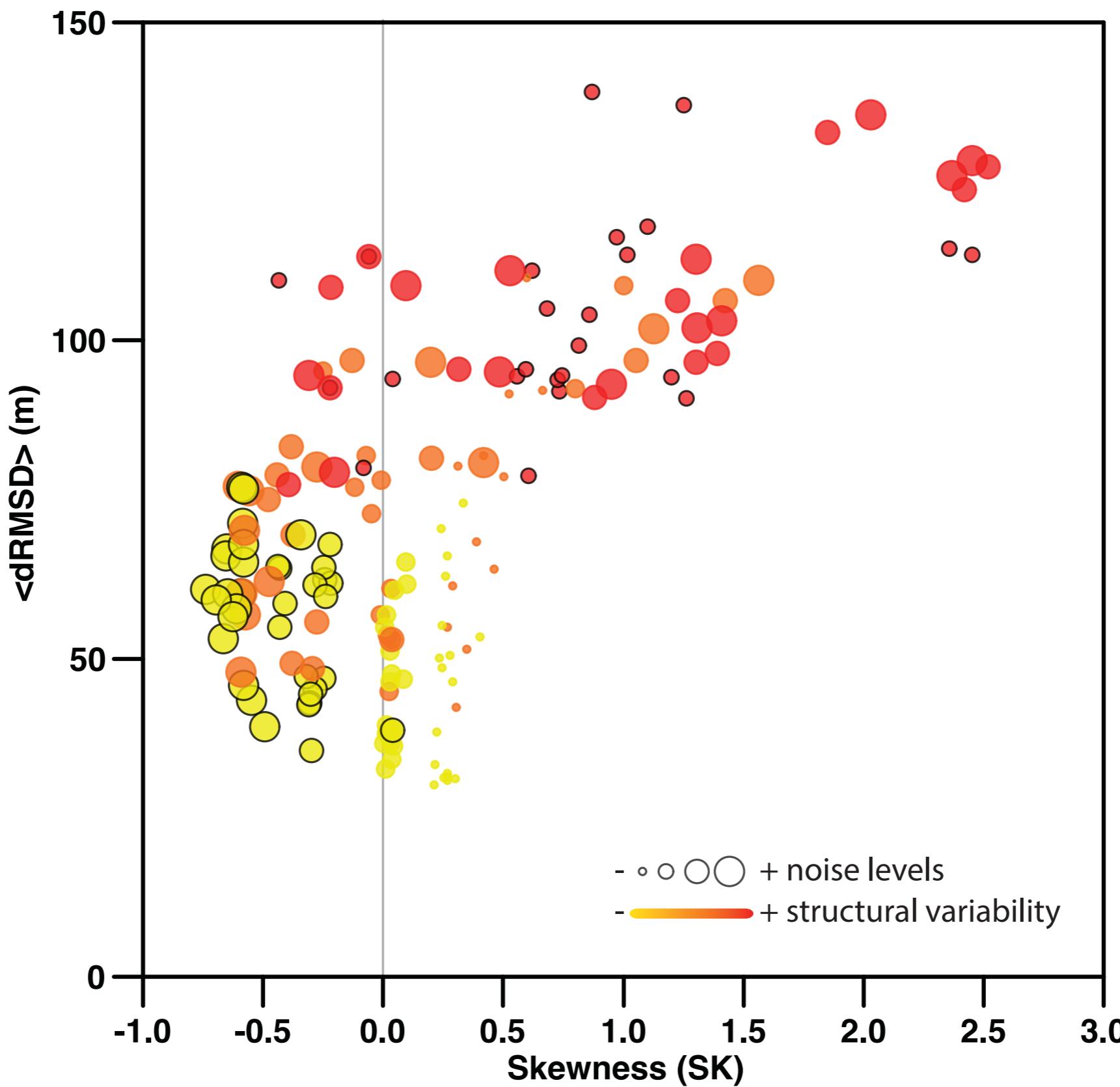
Structural variability is “NOT OK”



Can we predict the accuracy of the models?

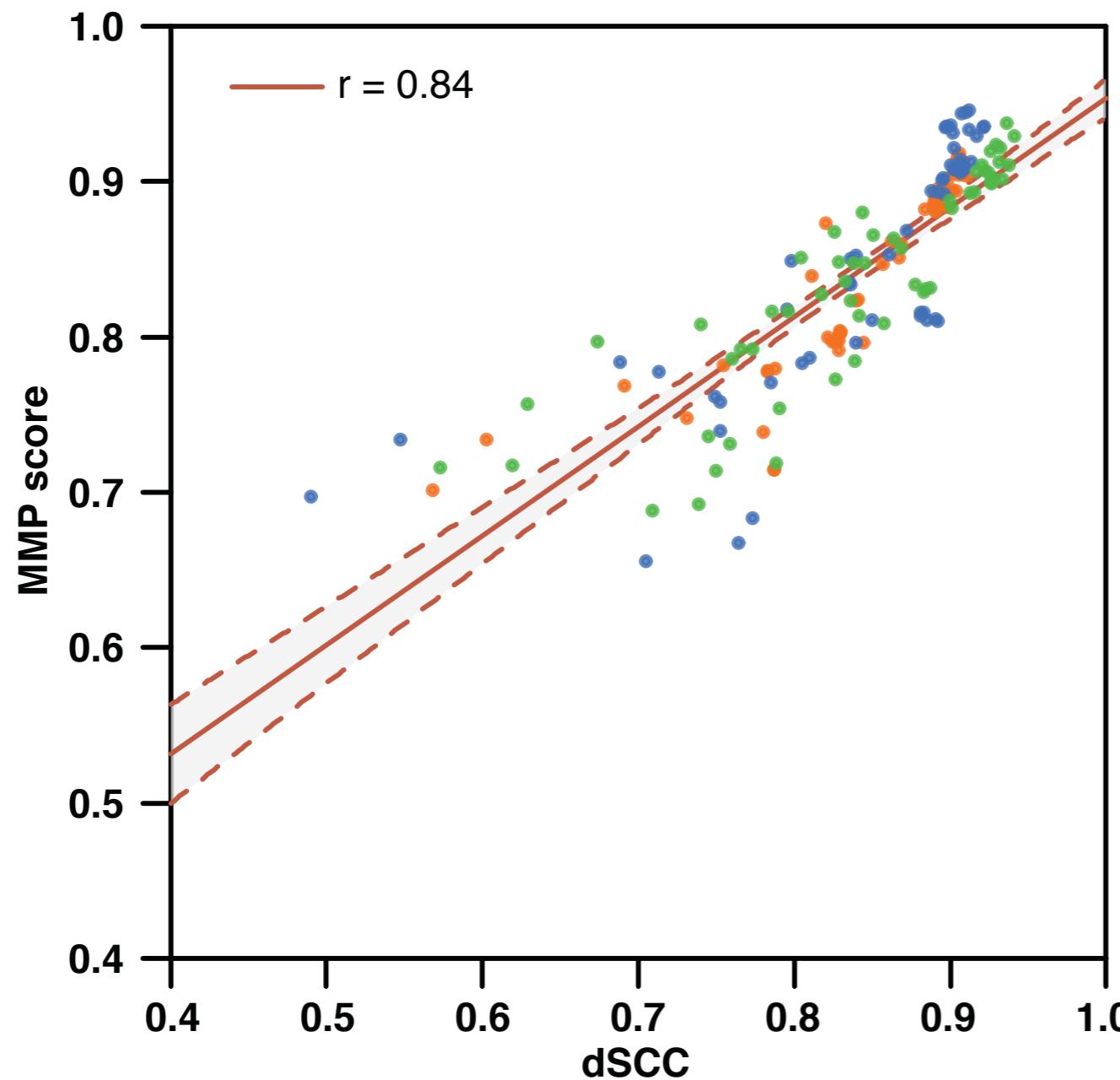


Skewness "side effect"

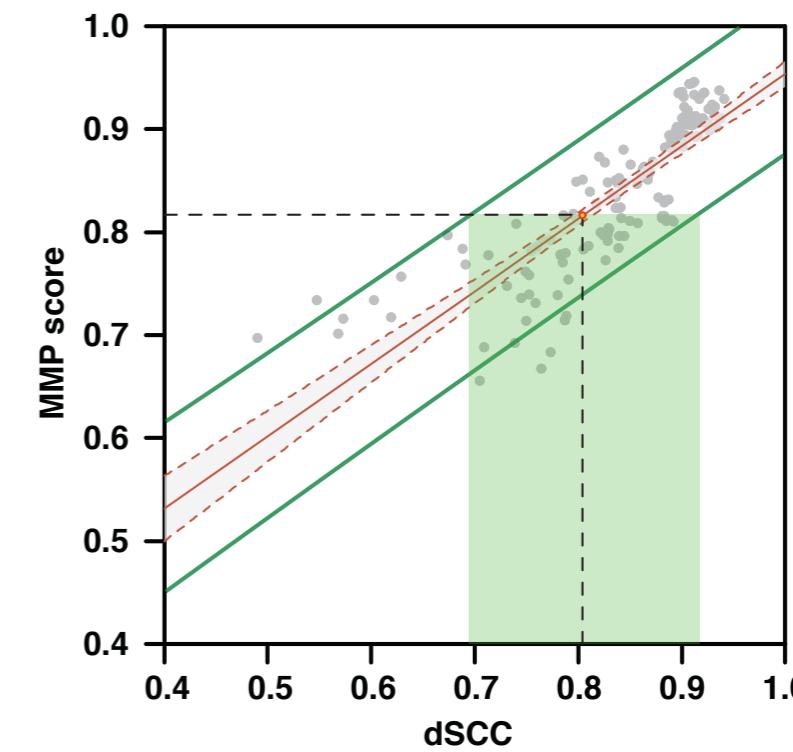
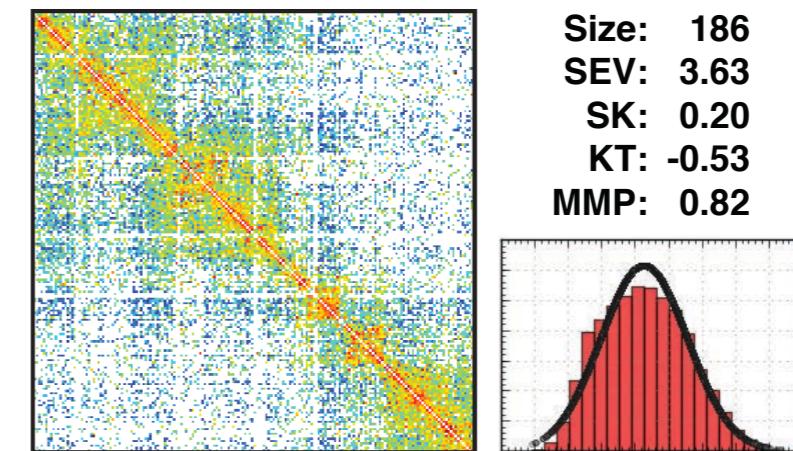


Can we predict the accuracy of the models?

$$\text{MMP} = -0.0002 * \text{Size} + 0.0335 * \text{SK} - 0.0229 * \text{KU} + 0.0069 * \text{SEV} + 0.8126$$



Human Chr1:120,640,000-128,040,000



Higher-res is "good"

put your \$\$ in sequencing

Noise is "OK"

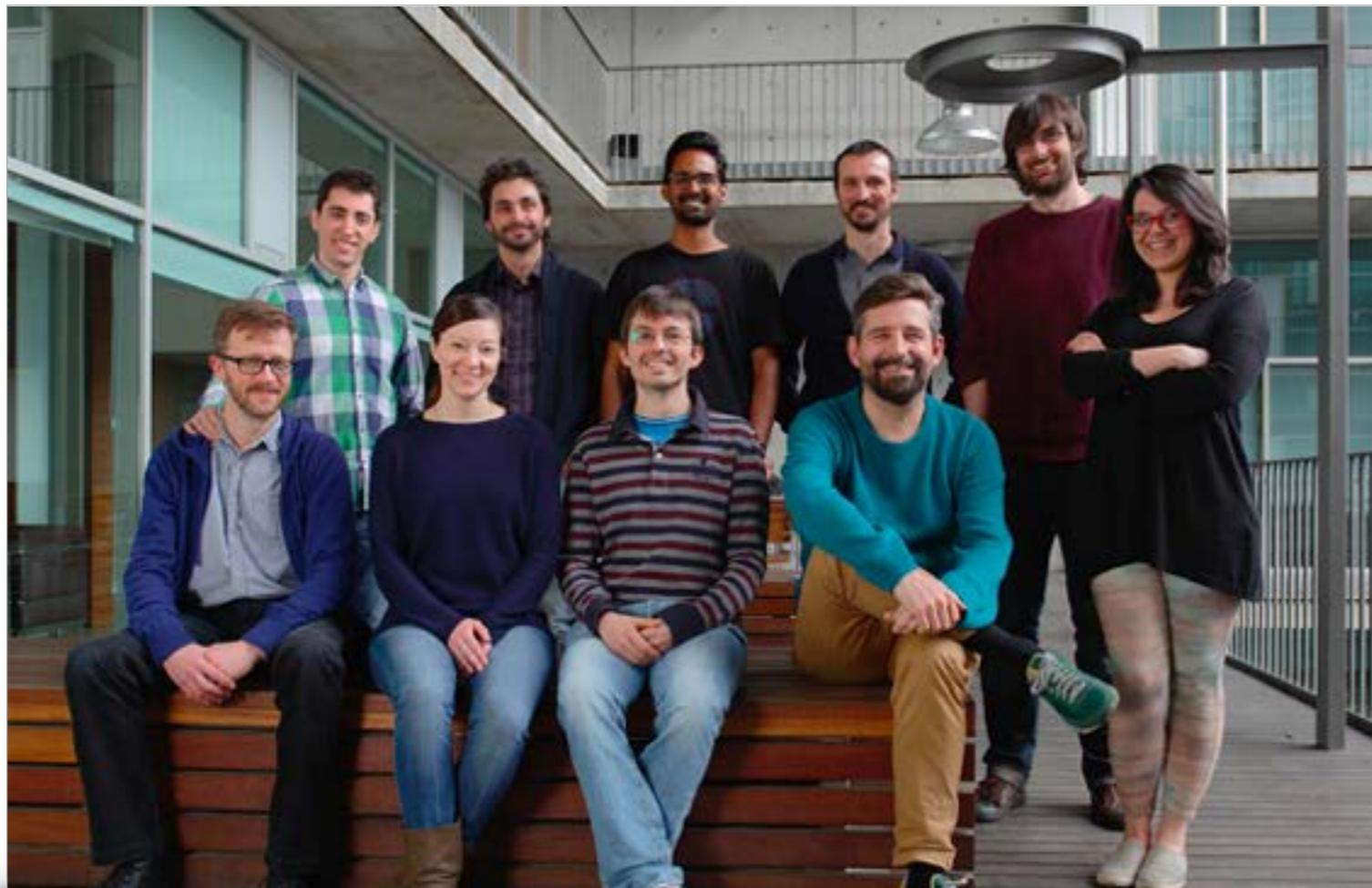
no need to worry much

Structural variability is "NOT OK"

homogenize your cell population!

...but we can differentiate between noise and structural variability

and we can a priori predict the accuracy of the models



Marie Trussart
François Serra
Davide Baù

Gireesh K. Bogu
Yasmina Cuartero
François le Dily
David Dufour
Irene Farabella
Silvia Galan
Mike Goodstadt
Francisco Martínez-Jiménez
Paula Soler
Yannick Spill
Marco di Stefano

in collaboration with Ivan Junier (Université Joseph Fourier) & Luís Serrano (CRG)

<http://marciuslab.org>
<http://3DGenomes.org>
<http://cnag.crg.eu>

