# Data integration for 3D structure determination.

**Marc A. Marti-Renom**
*Genome Biology Group (CNAG)*
*Structural Genomics Group (CRG)*
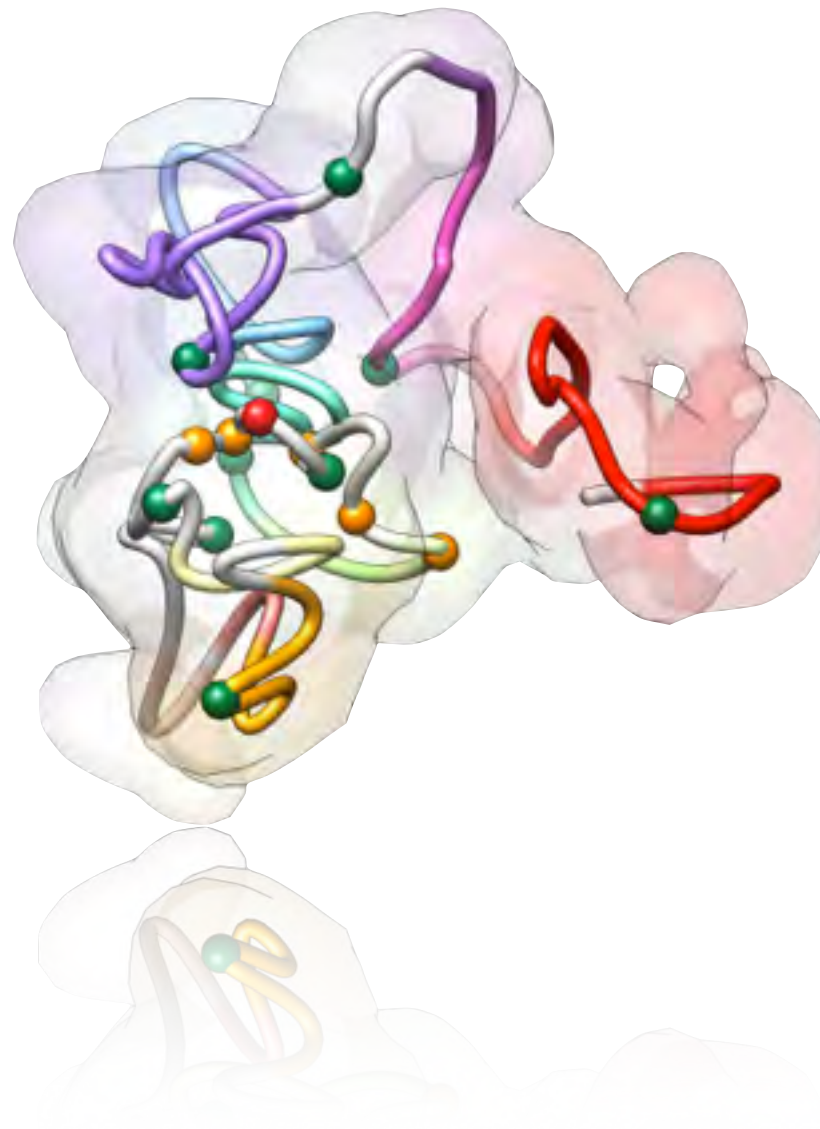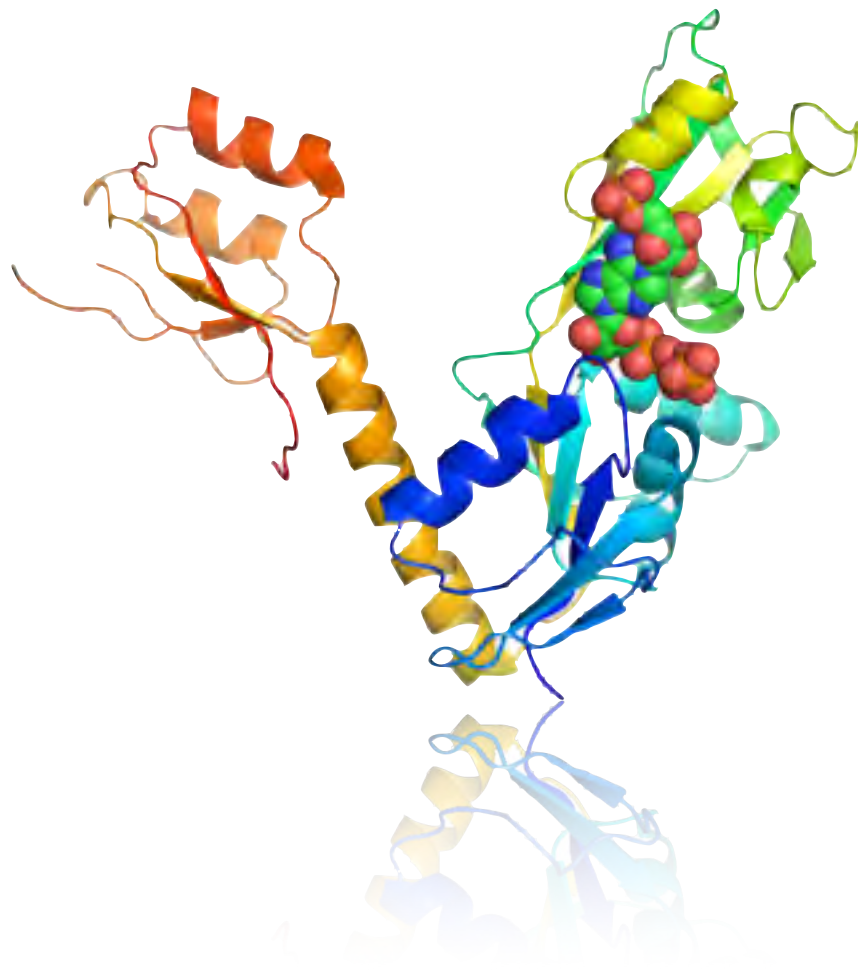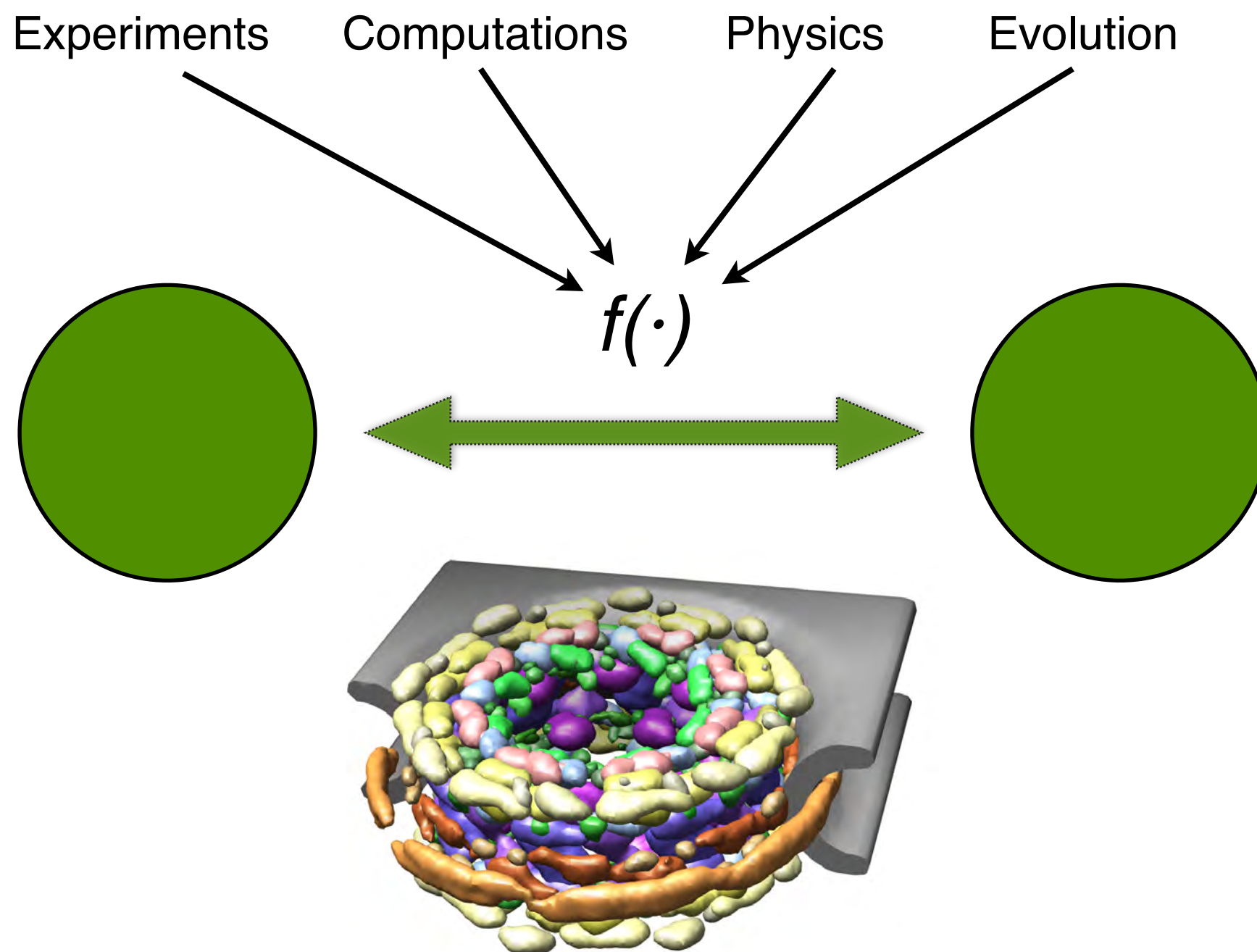
ICREA

cnag CRG
Centre
for Genomic
Regulation

# Structural Genomics Group

http://www.marciuslab.org

# Integrative Modeling Platform

http://www.integrativemodeling.org

Experiments    Computations    Physics    Evolution

$f(\cdot)$

# Stages

**Stage 1: Gathering Information.** Information is collected in the form of data from wet lab experiments, as well as statistical tendencies such as atomic statistical potentials, physical laws such as molecular mechanics force fields, and any other feature that can be converted into a score for use to assess features of a structural model.

**Stage 2: Choosing How To Represent And Evaluate Models.** The resolution of the representation depends on the quantity and resolution of the available information and should be commensurate with the resolution of the final models: different parts of a model may be represented at different resolutions, and one part of the model may be represented at several different resolutions simultaneously. The scoring function evaluates whether or not a given model is consistent with the input information, taking into account the uncertainty in the information.
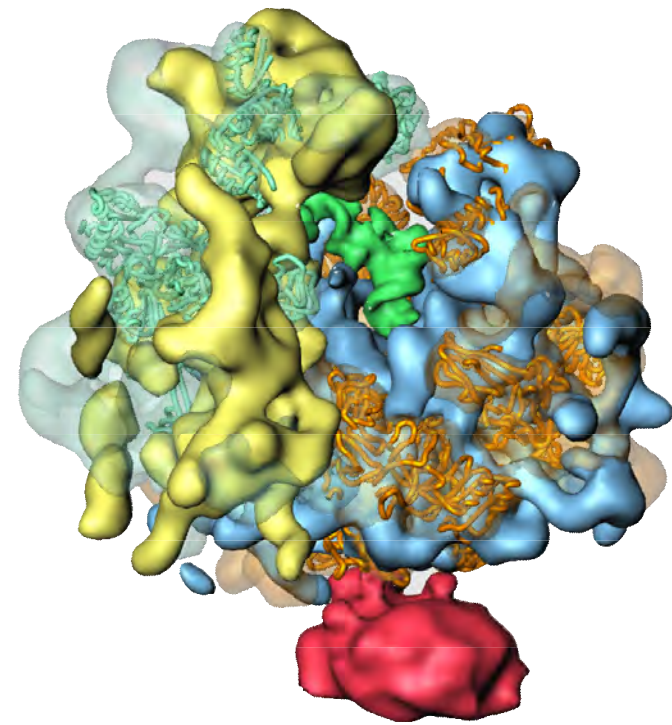
**Stage 3: Finding Models That Score Well.** The search for models that score well is performed using any of a variety of sampling and optimization schemes (such as the Monte Carlo method). There may be many models that score well if the data are incomplete or none if the data are inconsistent due to errors or unconsidered states of the assembly.

**Stage 4: Analyzing Resulting Models and Information.** The ensemble of good-scoring models needs to be clustered and analyzed to ascertain their precision and accuracy, and to check for inconsistent information. Analysis can also suggest what are likely to be the most informative experiments to perform in the next iteration.
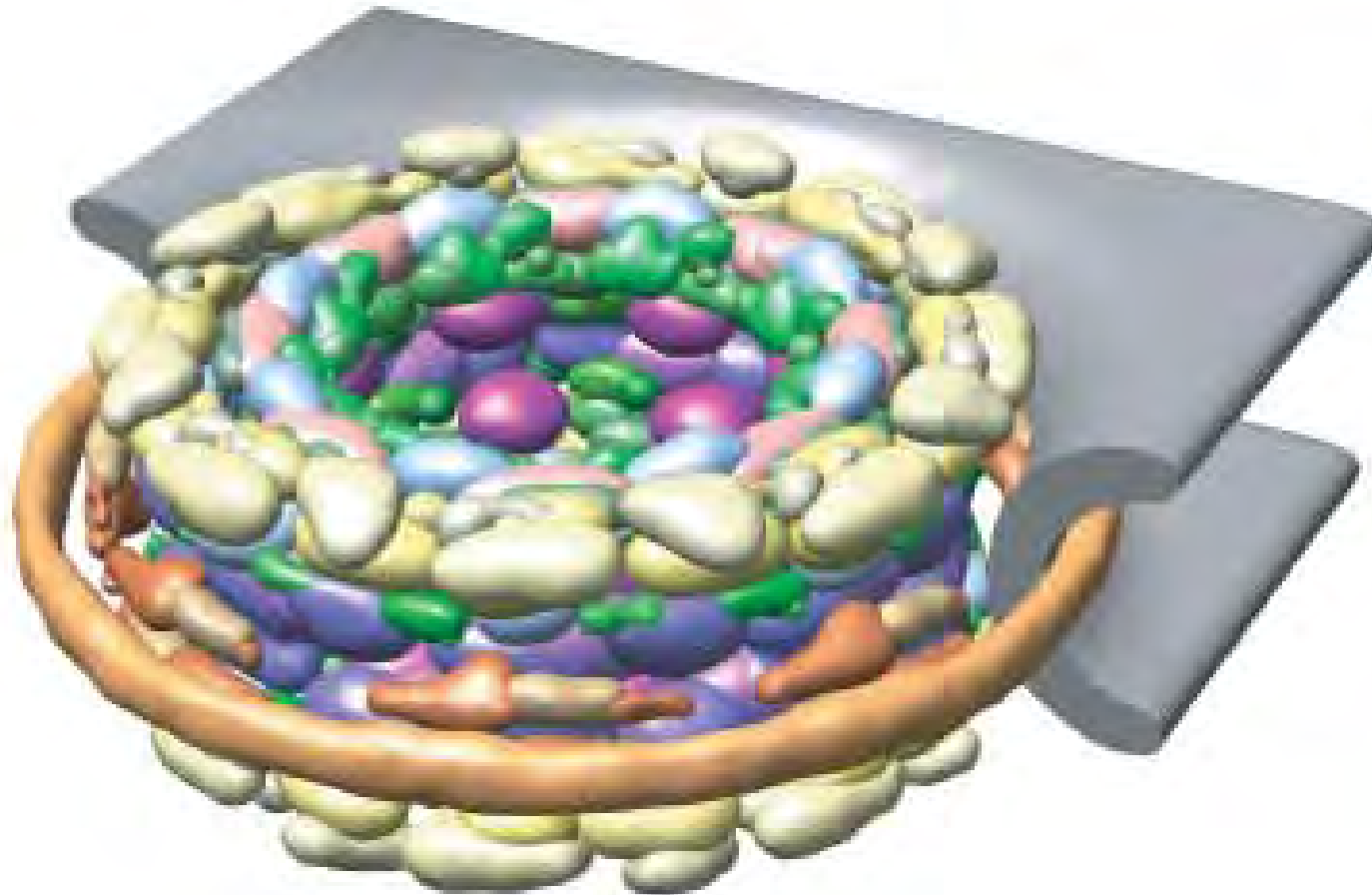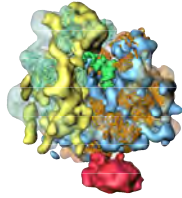
Integrative modeling iterates through these stages until a satisfactory model is built. Many iterations of the cycle may be required, given the need to gather more data as well as to resolve errors and inconsistent data.

Russel, D., Lasker, K., Webb, B., Velázquez-Muriel, J., Tjioe, E., Schneidman-Duhovny, D., Peterson, B., et al. (2012). *PLoS Biology*, *10*(1), e1001244
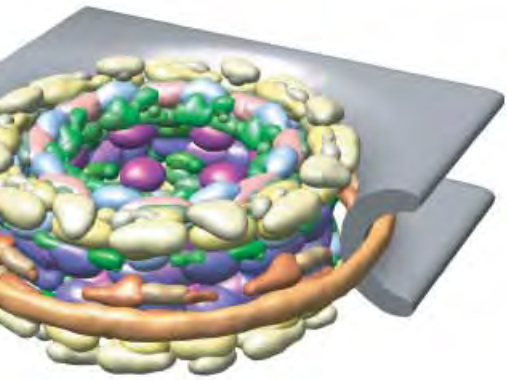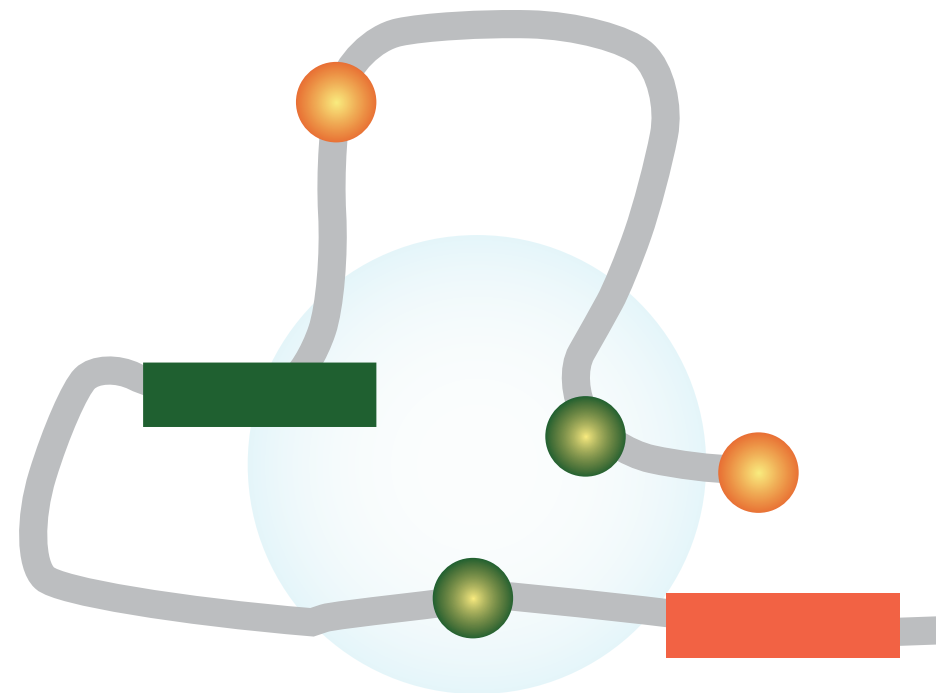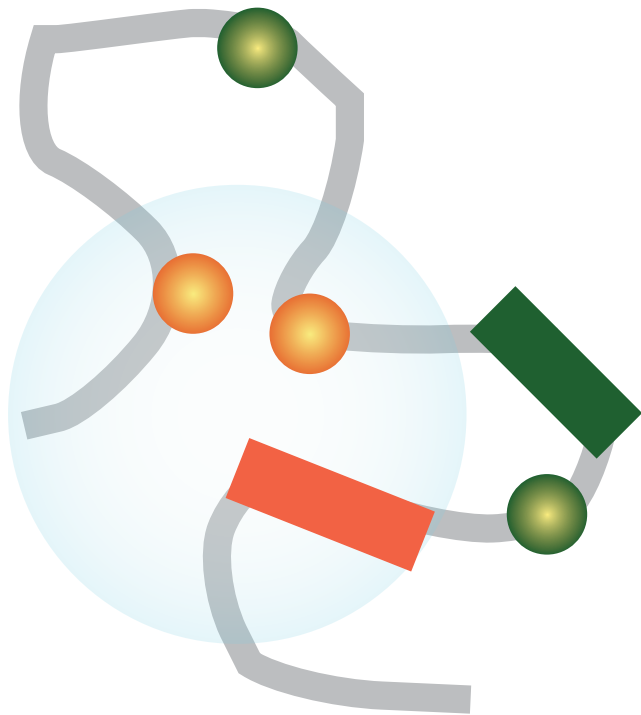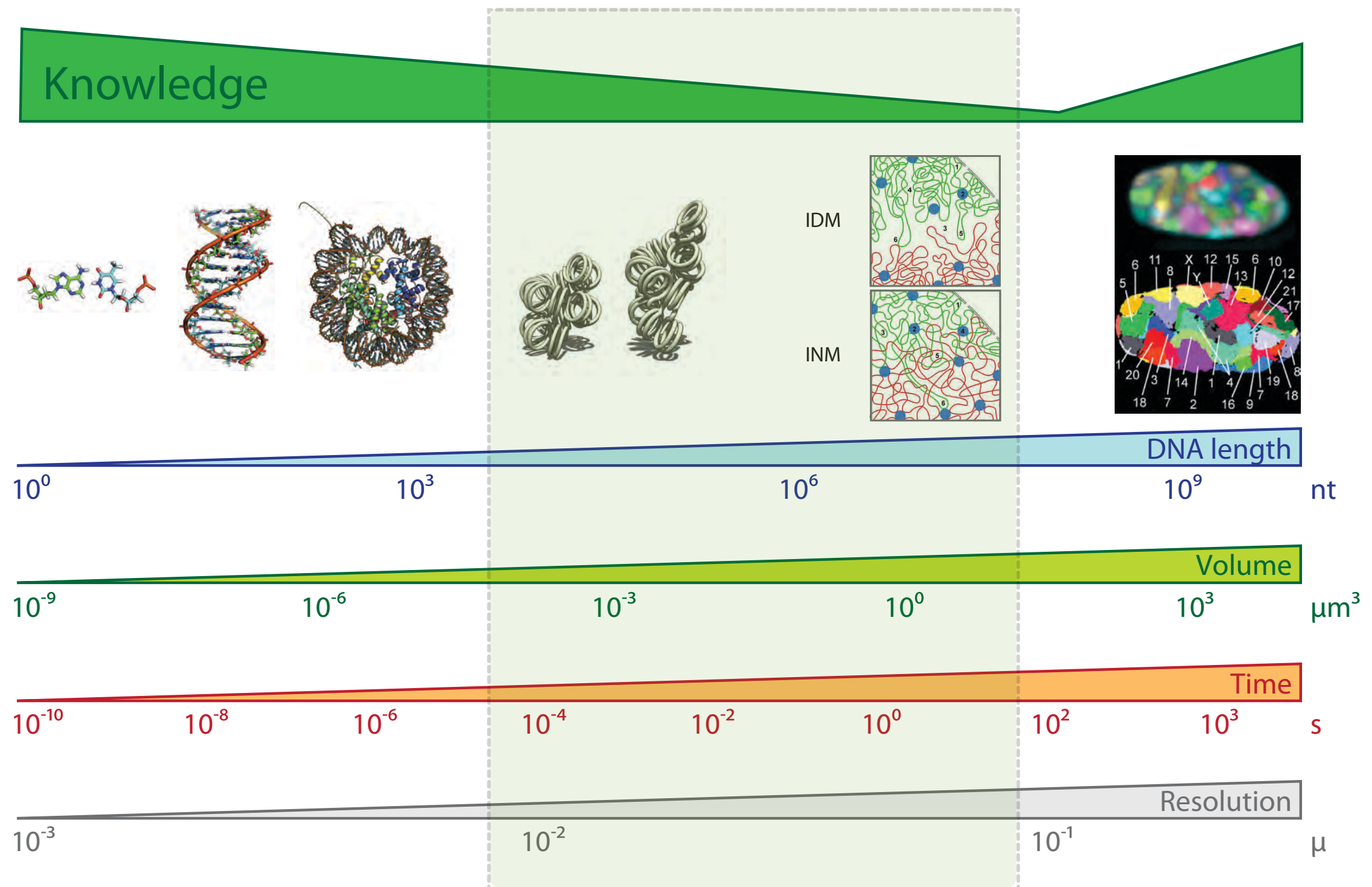
# Data Integration

# Data Integration

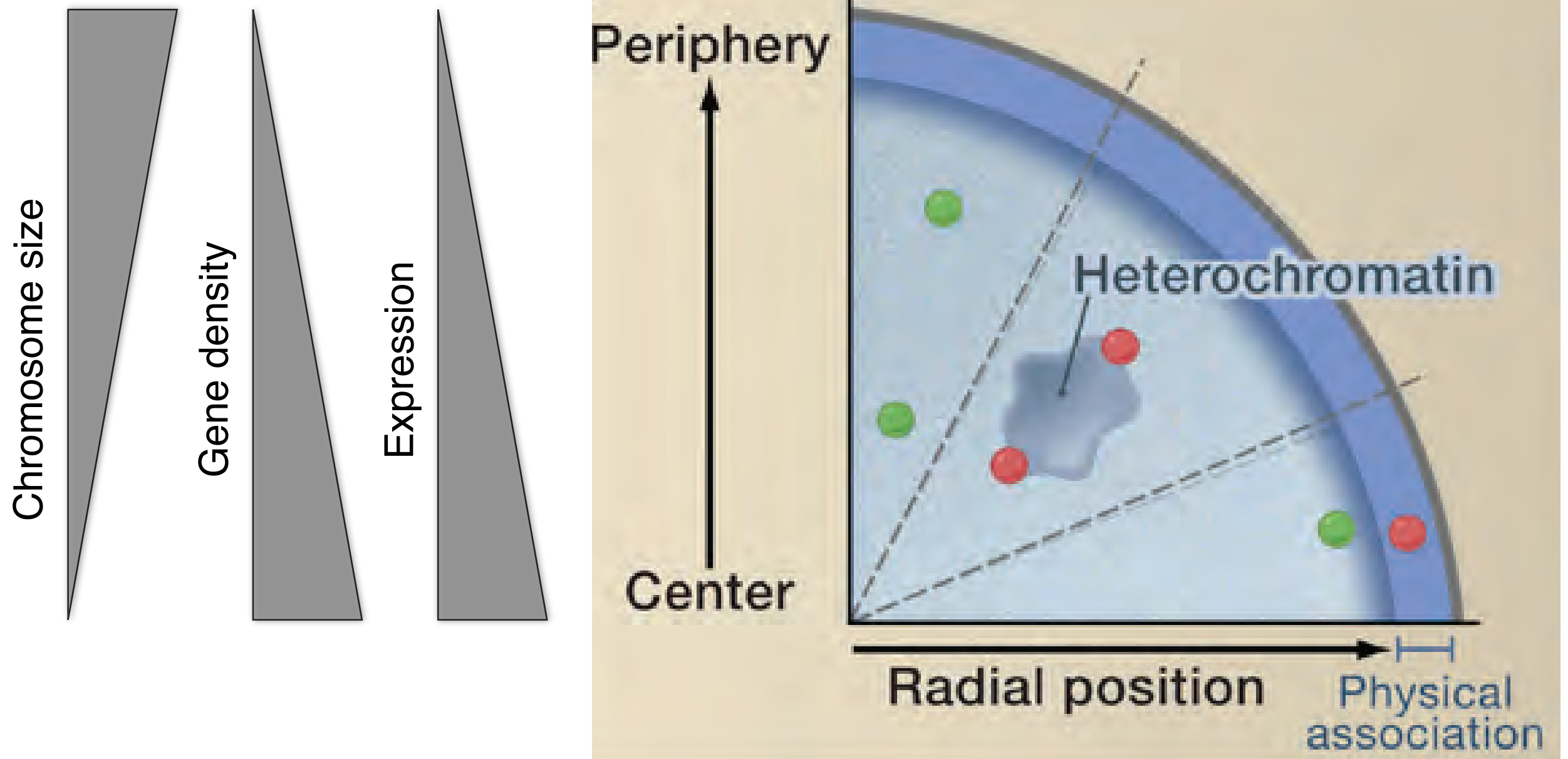# Data Integration
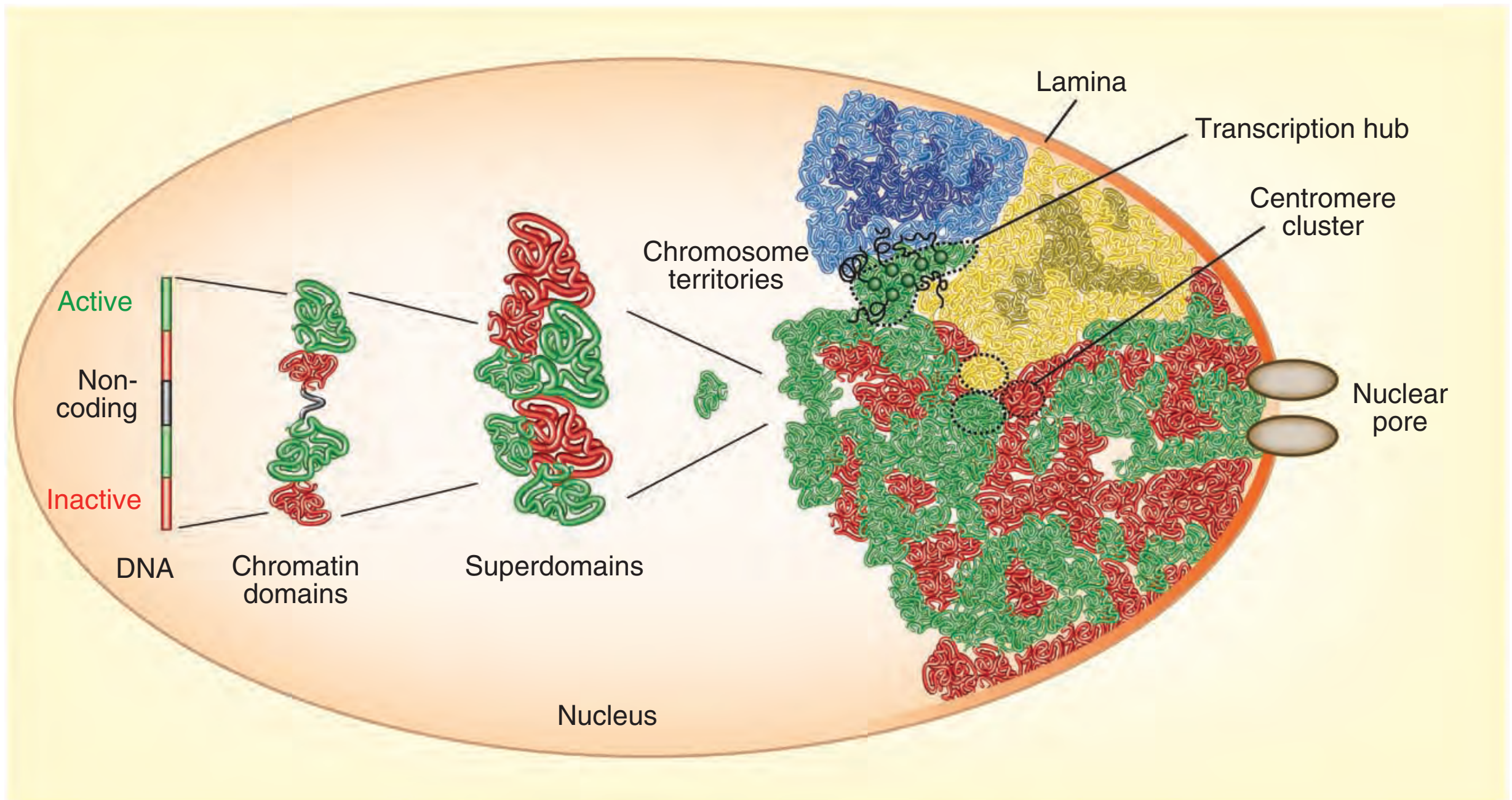
# Complex genome organization

# Resolution Gap

# Radial organization of the genome

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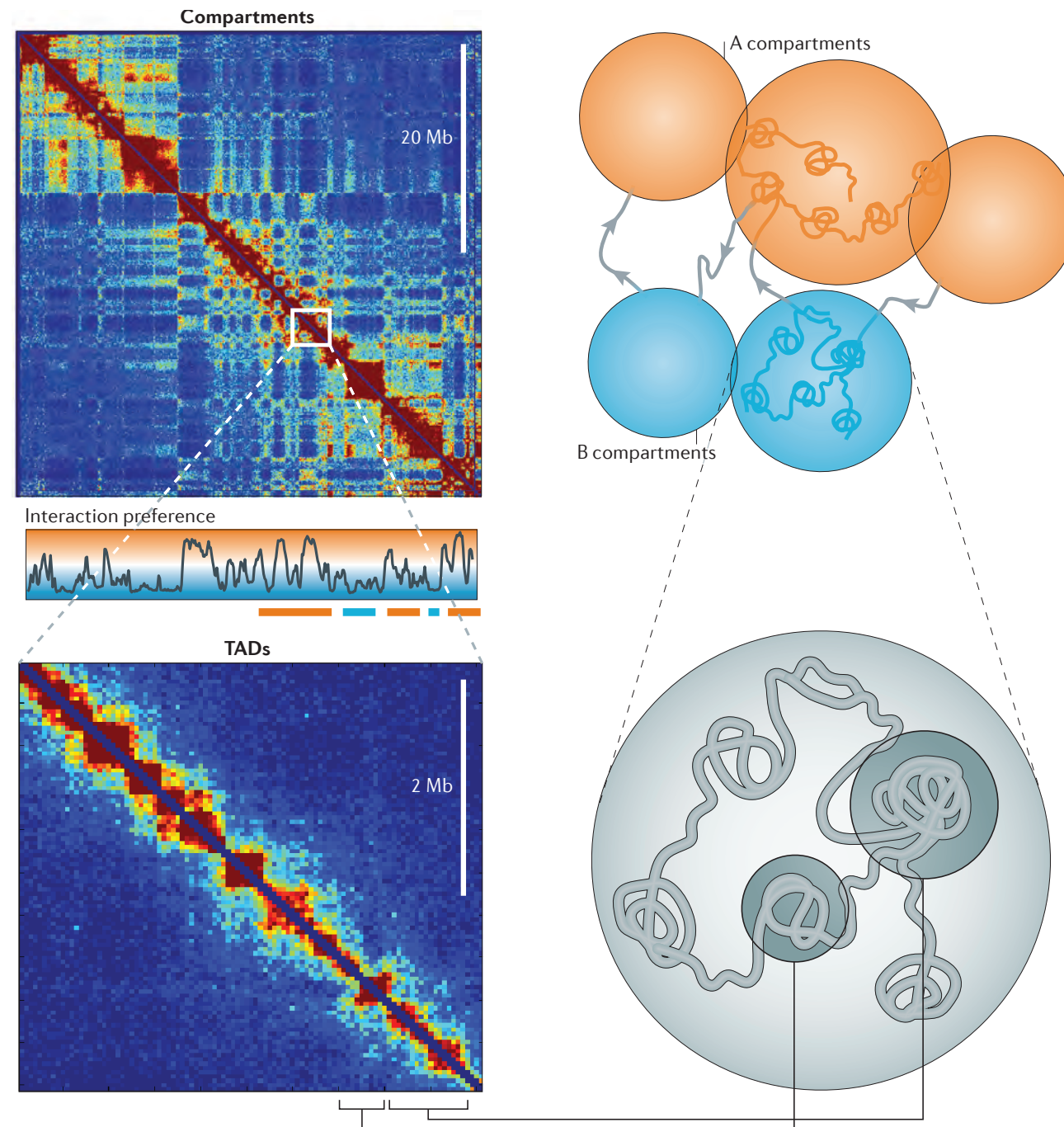
# Complex genome organization

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

# A/B Compartments and TADs

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).

# 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.

# Modeling Genomes

Experiments



Grow GM12878 and K562 cells

↓

Perform 3C analysis

↓

Perform 5C analysis with 30+25 primers

↓

Analyze 5C products by paired-end Solexa sequencing
(131,947 paired end reads per library)

Computation

Biomolecular structure determination
2D-NOESY data



Chromosome structure determination
5C data

# Chromosome Conformation Capture

cnag · CRG Centre for Genomic Regulation

| | 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] |

cnag  CRG Centre for Genomic Regulation

# Modeling 3D Genomes

Baù, D. & Marti-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).

# Human α-globin domain

## ENm008 genomic structure and environment



K562 cells:
α-globin genes active

# Representation

**Harmonic**

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

**Harmonic Lower Bound**

$$\begin{cases} if \ \ d_{i,j} \le d_{i,j}^0; & lbH_{i,j} = k\left(d_{i,j} - d_{i,j}^0\right)^2 \\[2mm] if \ \ d_{i,j} > d_{i,j}^0; & lbH_{i,j} = 0 \end{cases}$$

**Harmonic Upper Bound**

$$\begin{cases} if \ \ d_{i,j} \ge d_{i,j}^0; & ubH_{i,j} = k\left(d_{i,j} - d_{i,j}^0\right)^2 \\[2mm] if \ \ d_{i,j} < d_{i,j}^0; & ubH_{i,j} = 0 \end{cases}$$

# Scoring



GM12878

70 fragments
1,520 restraints

K562

70 fragments
1,049 restraints

Harmonic    Harmonic Lower Bound    Harmonic Upper Bound

# Optimization

# Clustering



Minimum IMP Objective function in cluster

Cluster #1
2780 model
910,280 IMP OF

Cluster #2
2668 model
910,400 IMP OF

Cluster #4
2270 model
915,890 IMP OF

Cluster #3
2282 model
915,890 IMP OF

# Not just one solution

# The "Chromatin Globule" model



Active genes
Inactive genes
CTCF sites
HS sites
Globule core

(a) chromatin subcompartment
subcompartment loop base spring (magnified)

Münkel et al. JMB (1999)

Osborne et al. Nat Genet (2004)

al. Science (2009)

cnag    CRG Centre for Genomic Regulation

# Caulobacter crescentus genome

# The 3D architecture of Caulobacter Crescentus

4,016,942 bp & 3,767 genes



Origin

= + Strand
= - Strand

Terminus

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

~13Kb

# 5C interaction matrix

ELLIPSOID for Caulobacter cresentus

# 3D model building with the 5C + IMP approach



339 mers

# Genome organization in Caulobacter crescentus

Arms are helical

Resolution

*dif* site 47±17Kb from Ter

Centromer-like

*parS* sites 25±17Kb from Ori



| Cluster 1 | Cluster 2 | Cluster 3 | Cluster 4 |

180°    180°    180°    180°

500 nm    500 nm    500 nm    500 nm

**MIRRORS!**

cnag   CRG Centre for Genomic Regulation

# Moving the parS sites 400 Kb away from Ori

# Moving the parS sites results in whole genome rotation!

# Genome architecture in Caulobacter

# From Sequence to Function
## 5C + IMP

**Technology**



**Hypothesis**

**Function!**

D. Baù and M.A. Marti-Renom **Chromosome Res** (2011) 19:25-35.

# Bacteria has also TADs (CIDs)

**Fig. 1. Partitioning of the *Caulobacter* chromosome into chromosomal interaction domains (CIDs). (A)**

# On TADs and hormones



Davide Baù

François le Dily

# Progesterone-regulated transcription in breast cancer



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

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

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

**Regulation in 3D?**

# Experimental design



ChIP-Seq
RNA-Seq
Hi-C

HiC libraries

- Pg

+ Pg

HindIII          NcoI

Chr.18 (Hind III)          Chr.18 (NcoI)
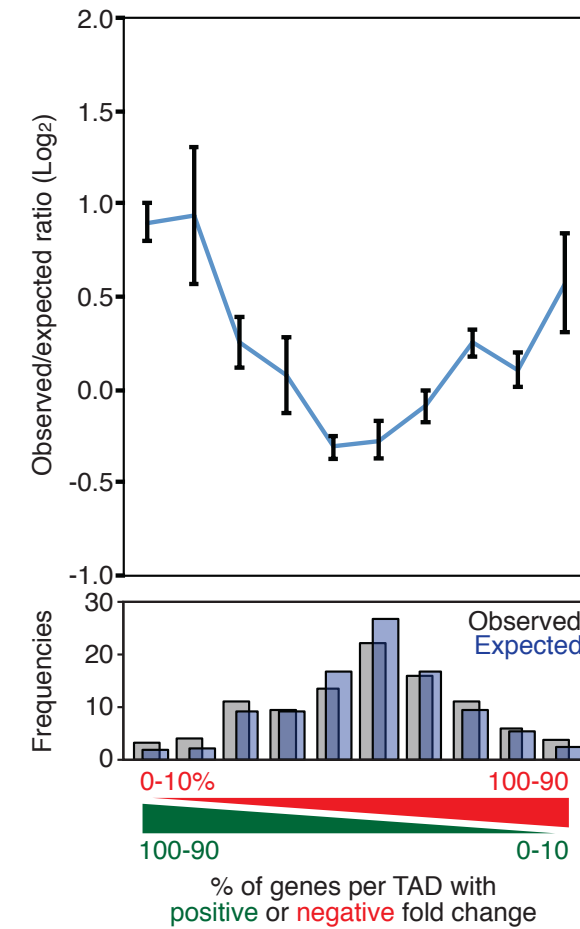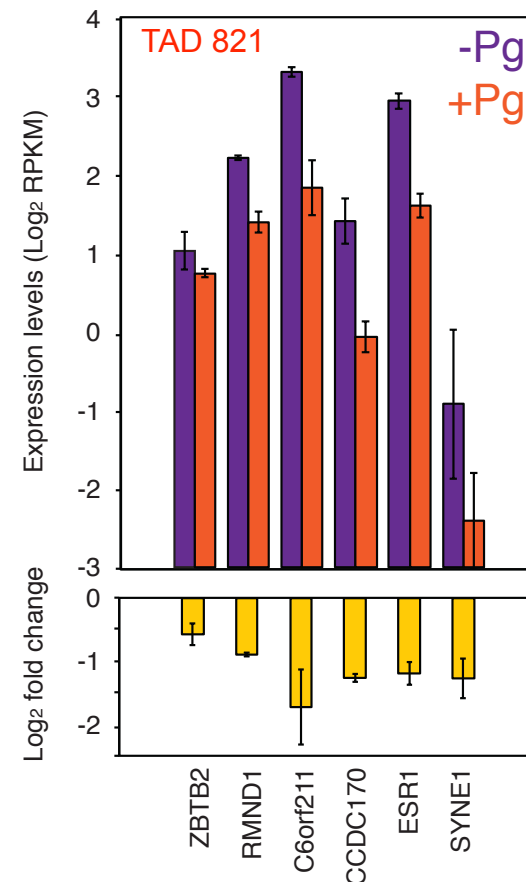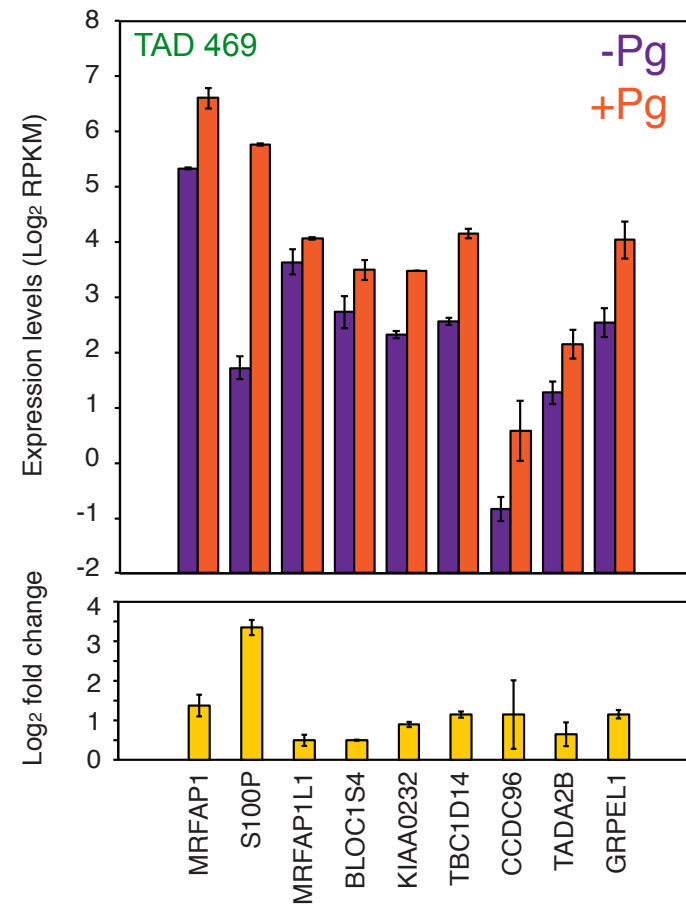
# Are there TADs? how robust?



>2,000 detected TADs

Chr.18
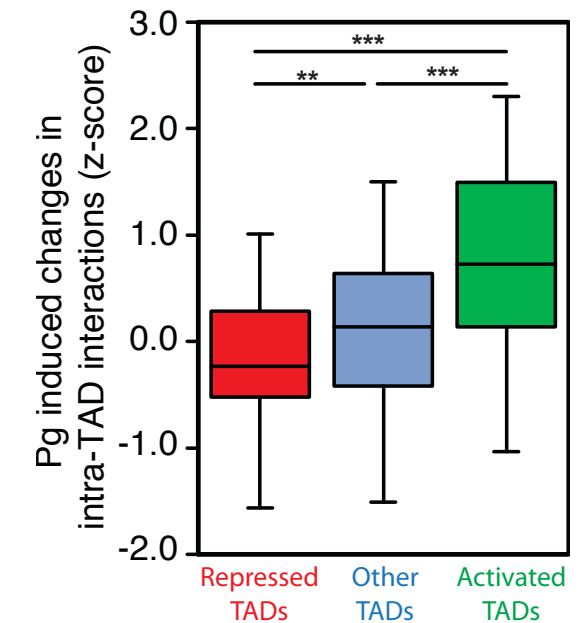
-Pg    +Pg
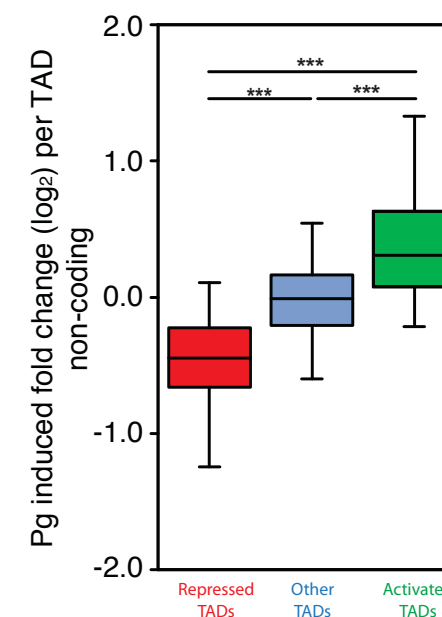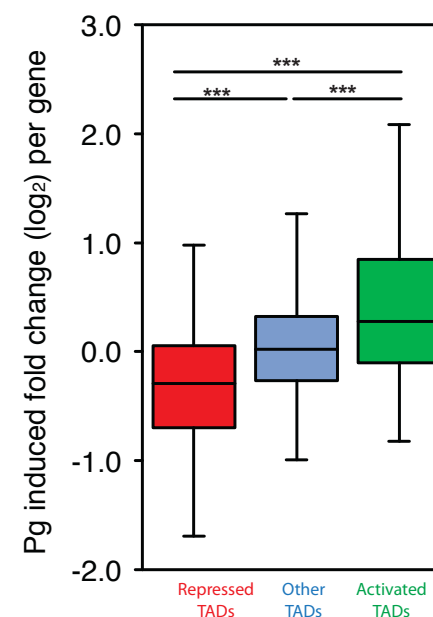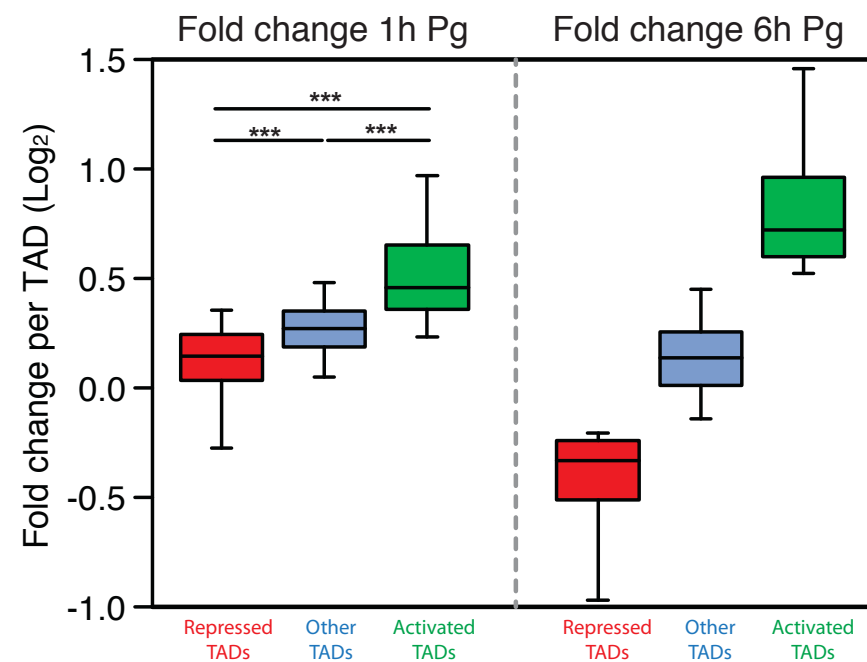
conserved
100kb
±200kb or more

# Are TADs homogeneous?

# Do TADs respond differently to Pg treatment?

# Do TADs respond differently to Pg treatment?
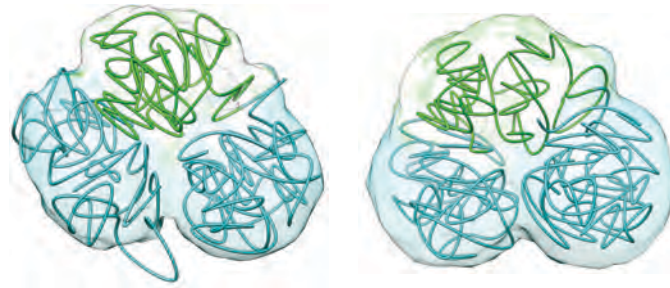
# Modeling 3D TADs



Chr1:26,800,000-28,700,000
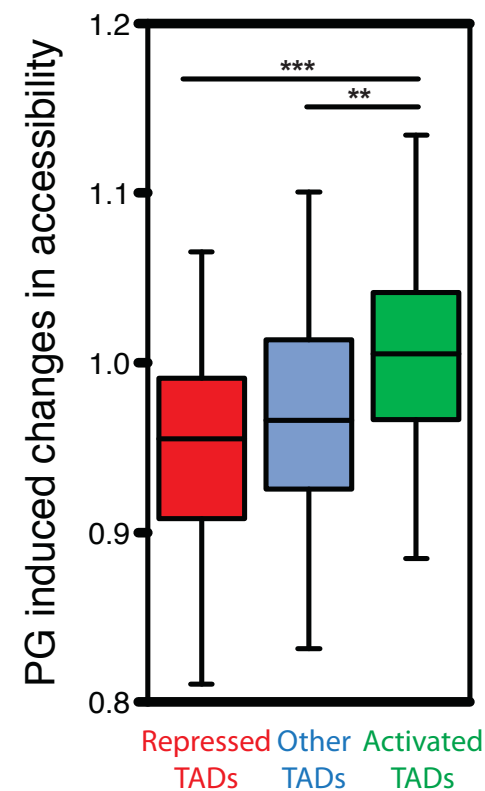
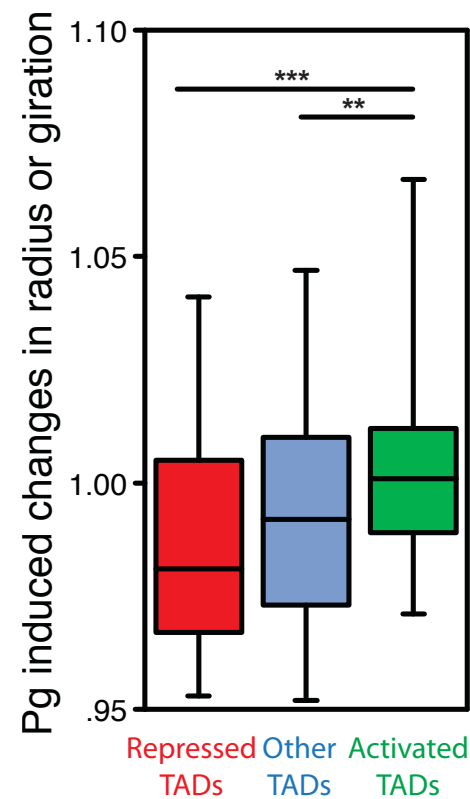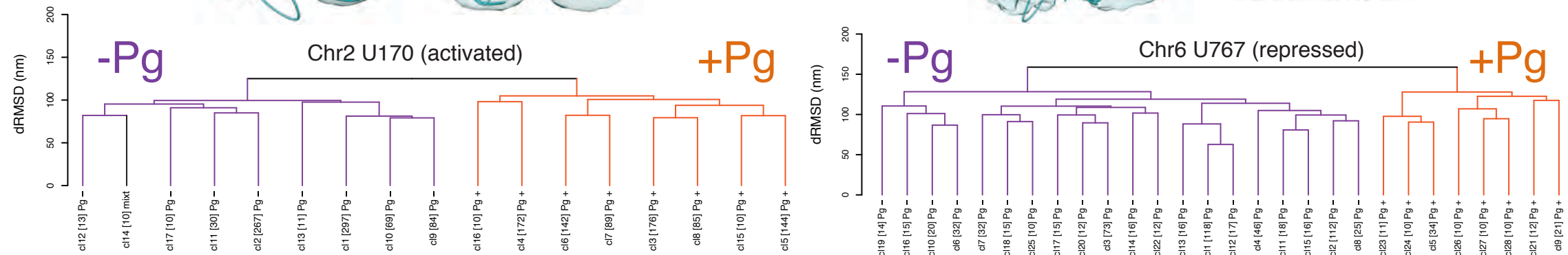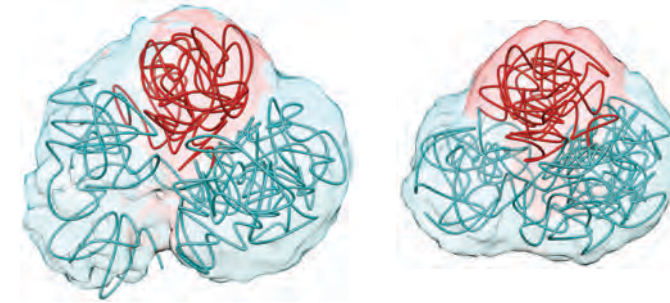61 genomic regions containing 209 TADs covering 267Mb

# How TADs respond structurally to Pg?
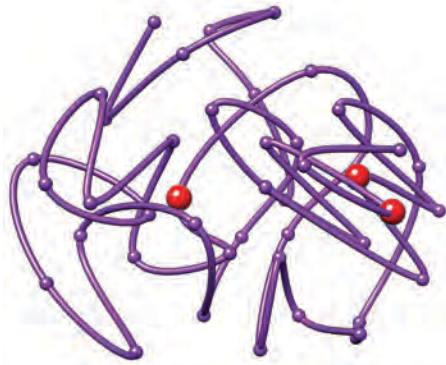


Chr2:9,600,000-13,200,000
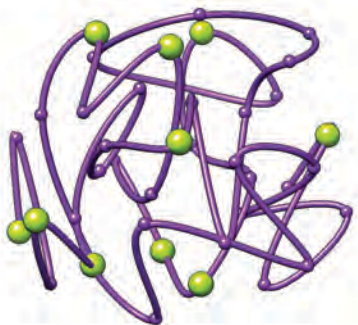
Chr6:71,800,000-76,500,000

# Model for TAD regulation

**Repressed TAD**
chr1 U41
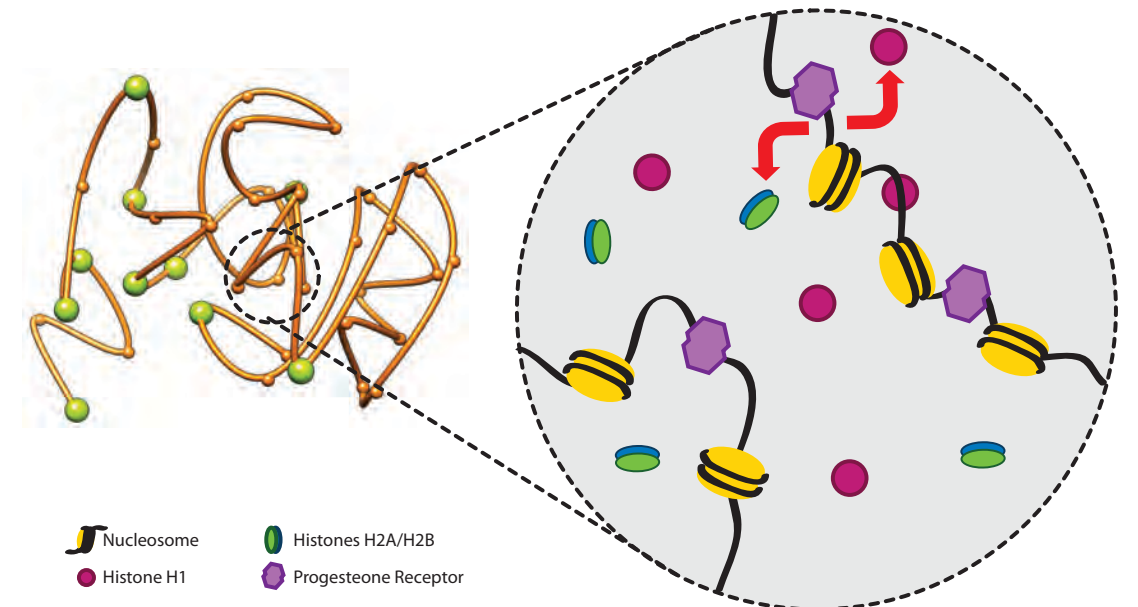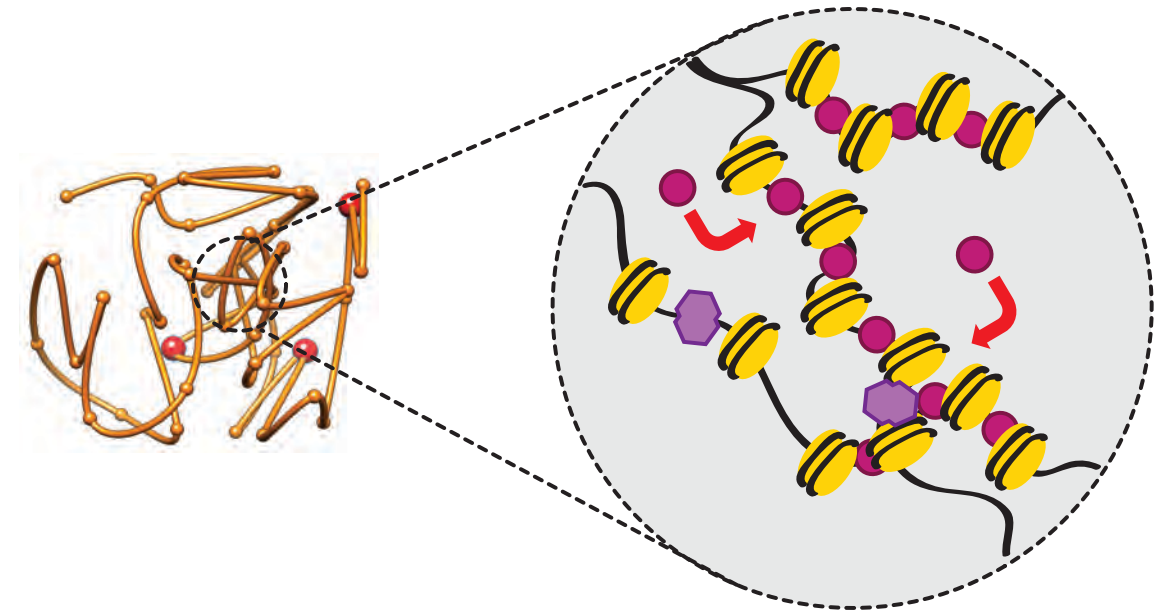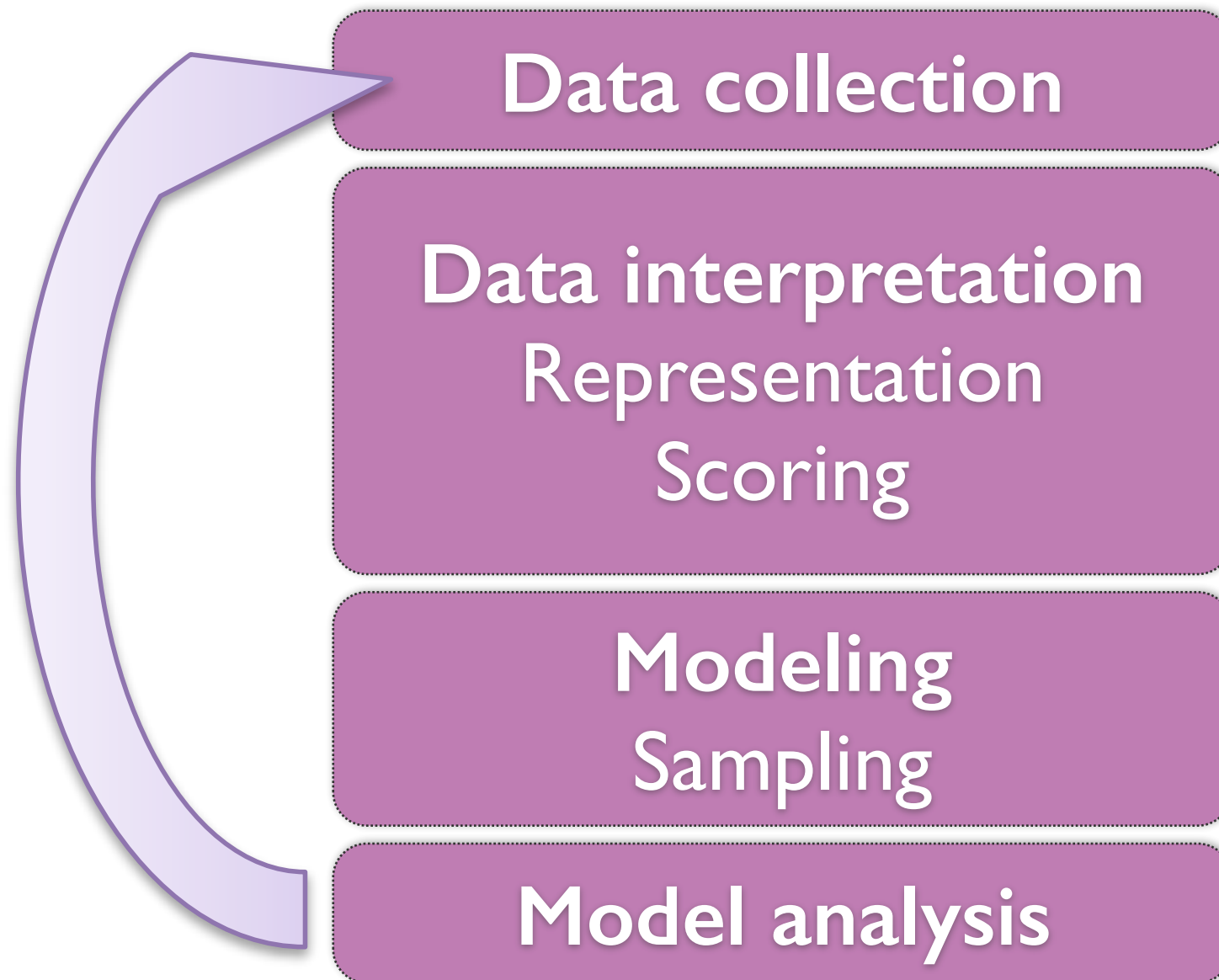
**Activated TAD**
chr2 U207

Structural transition **+Pg**

DHS HP1 H1.2 H2A MNAse H3K27me3 H3K9m3 H3K14ac H3K4me1 H3K36me2 H3K4me3

Nucleosome
Histones H2A/H2B
Histone H1
Progesteone Receptor

cnag CRG Centre for Genomic Regulation

# PLoS CB Outlook

Marti-Renom MA, Mirny LA (2011) PLoS Comput Biol 7(7): e1002125.



MURRE
Cell (2008) **133**:265-79

DOSTIE/BLANCHETTE
Genome Biol (2009) **10**: R37

DEKKER/LANDER/MIRNY
Science (2009) **326**:289-93

NOBLE
Nature (2010) **465**: 363-7

DEKKER/MARTI-RENOM
NSMB (2011) **18**:107-14

# Acknowledgments