# Data integration for 3D structure determination.
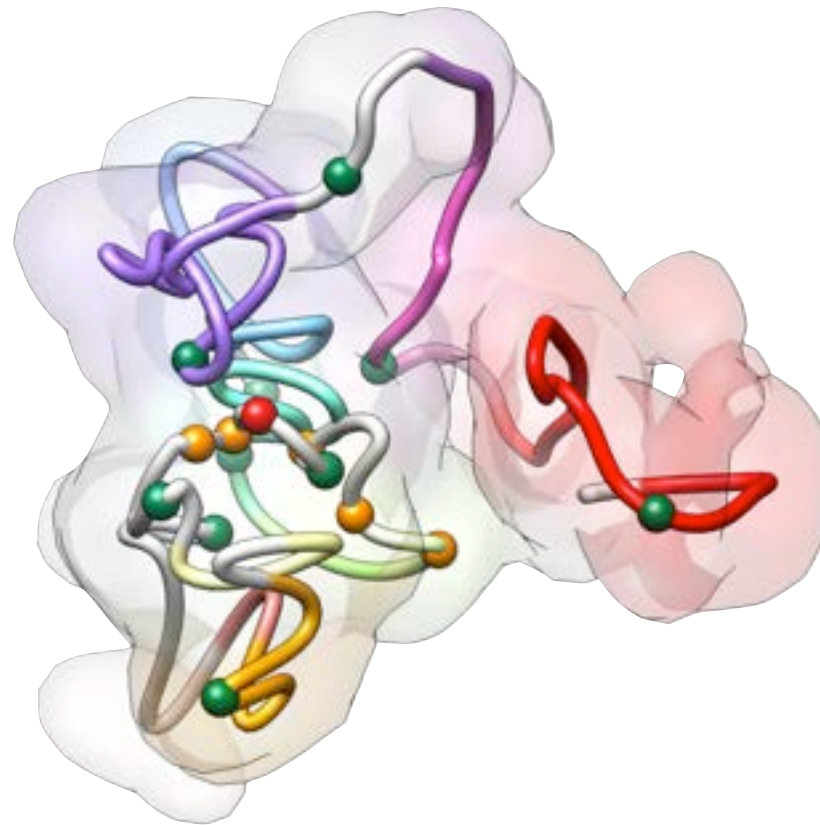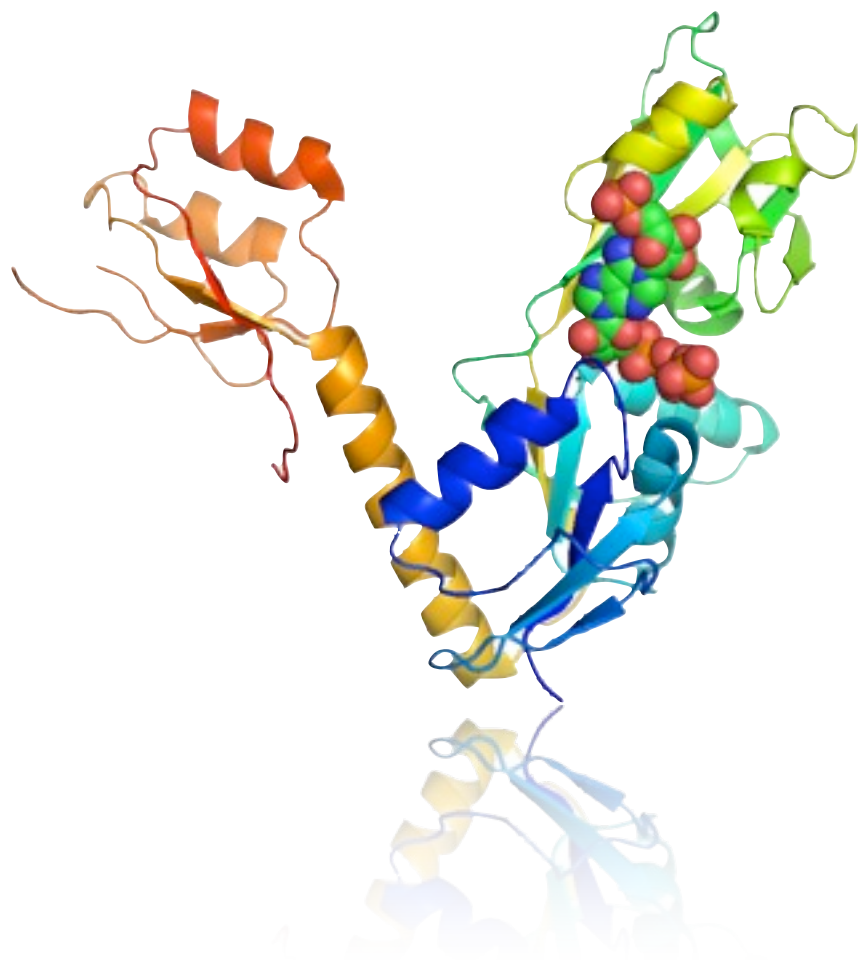
**Marc A. Marti-Renom**
*Genome Biology Group (CNAG)*
*Structural Genomics Group (CRG)*
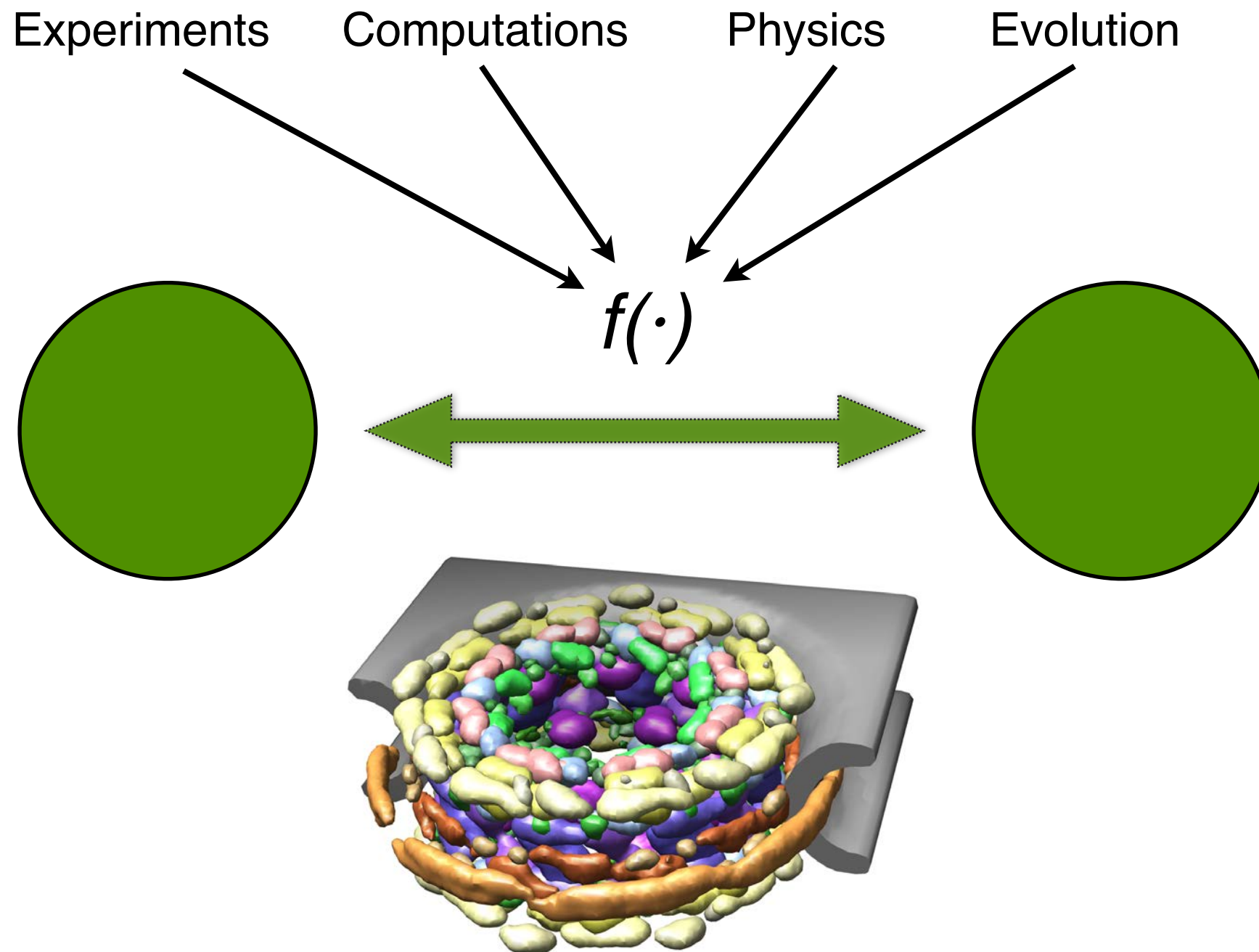
cnag

CRG
Centre
for Genomic
Regulation

# Structural Genomics Group

http://www.marciuslab.org

# Integrative Modeling Platform

http://www.integrativemodeling.org

Experiments     Computations     Physics     Evolution

$f(\cdot)$



*From: Russel, D. et al. PLOS Biology 10, e1001244 (2012).*

# 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
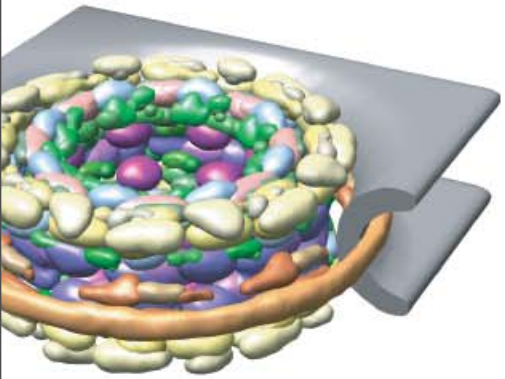
# Advantages

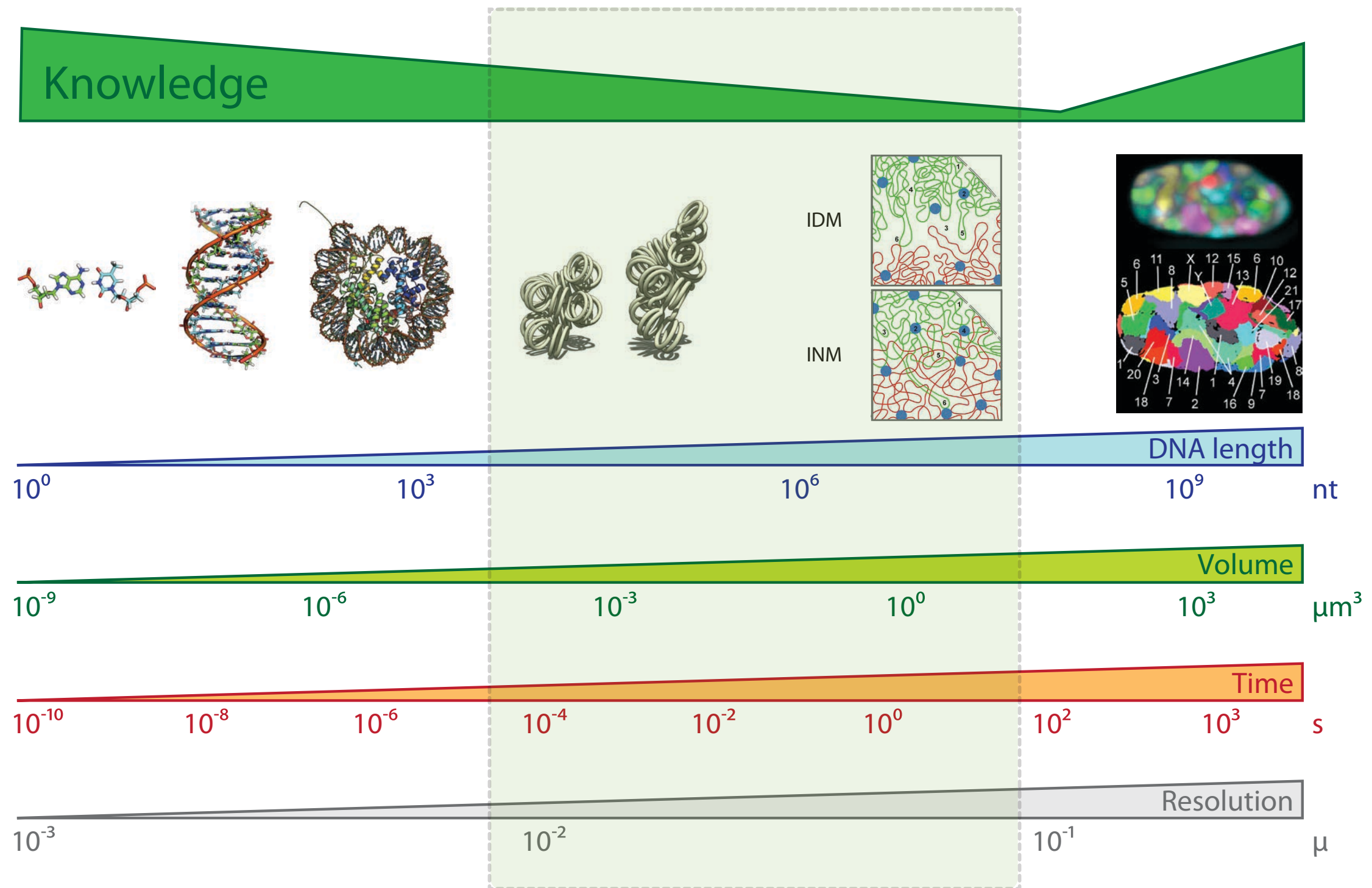**Using New Information.** Integrative modeling makes it easy to take advantage of new information and new types of information, resulting in a low barrier for using incremental information that is generally not applied to structure characterization. Even when a single data type is relatively uninformative, multiple types can give a surprisingly complete picture of an assembly [9,10].

**Maximizing Accuracy, Precision and Completeness.** Integrative models fit multiple types of information, and can thus be more accurate, precise, and complete than models based on the individual sources.

**Understanding and Assessing the Models.** By exhaustively sampling the space of models fitting the information, integrative modeling can find all models fitting the information, not only one. A full sampling of the models of a structure can improve the understanding of its function [49]. Because the data are encoded in scoring functions and the full set of models can be found, integrative modeling facilitates assessing the input information and output models in terms of precision and accuracy.

**Planning Experiments.** Integrative modeling provides feedback to guide future experiments, by computationally testing the impact of hypothetical datasets. As a result, experiments can be chosen to best improve our knowledge of the assembly.

**Understanding and Assessing Experimental Accuracy.** Data errors present a challenge for all methods of model building. Integrative modeling can detect inconsistent data as no models will exist that fit all the data. In addition, integrative modeling facilitates the application of more sophisticated methods for error estimation, such as Inferential Structure Determination [16].

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

cnag

# Data Integration

# Data Integration

# Data Integration

# Resolution Gap

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

# Complex genome organization

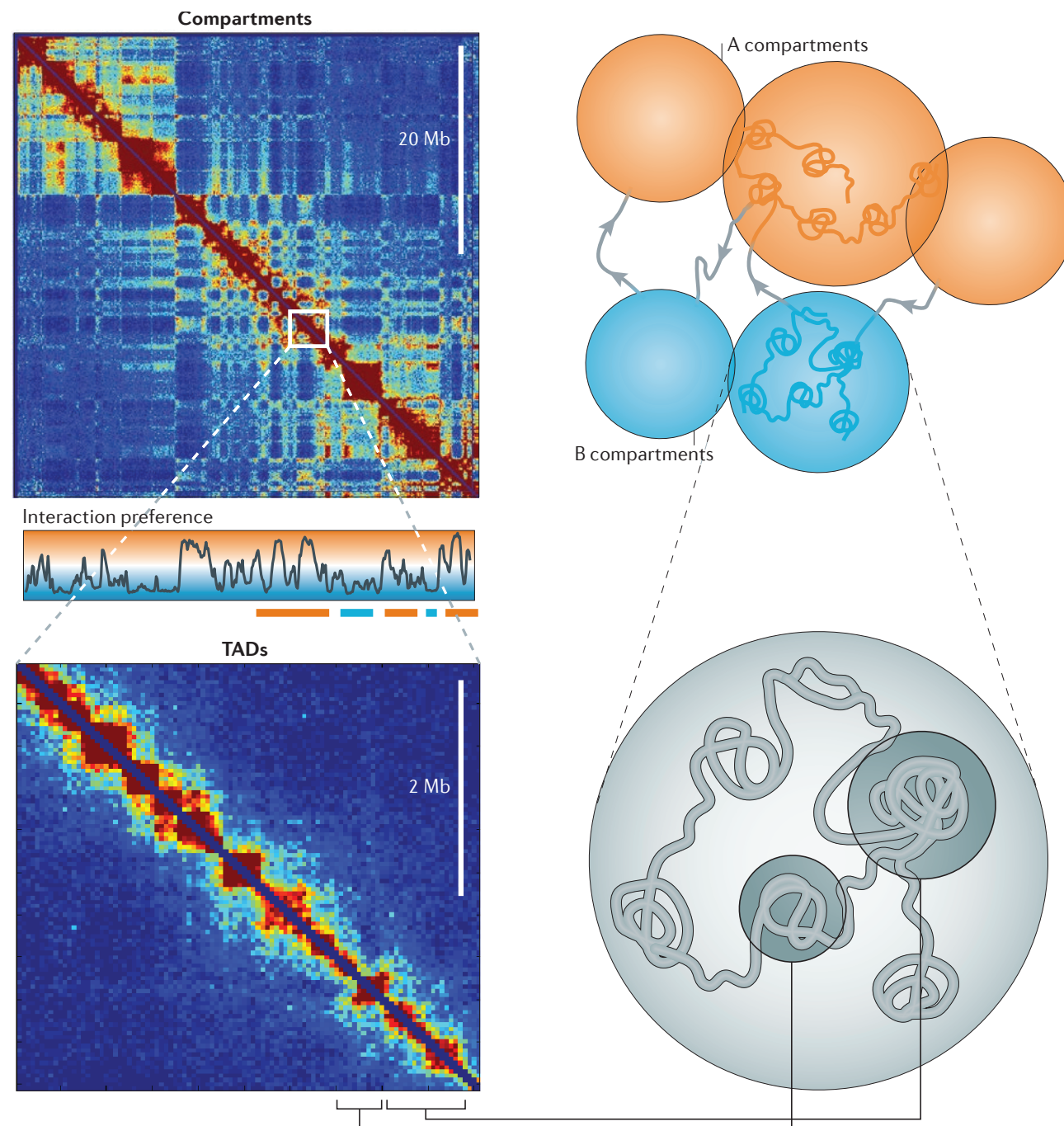Takizawa, T., Meaburn, K. J. & Misteli, T. The meaning of gene positioning. Cell 135, 9–13 (2008).

# Complex genome organization

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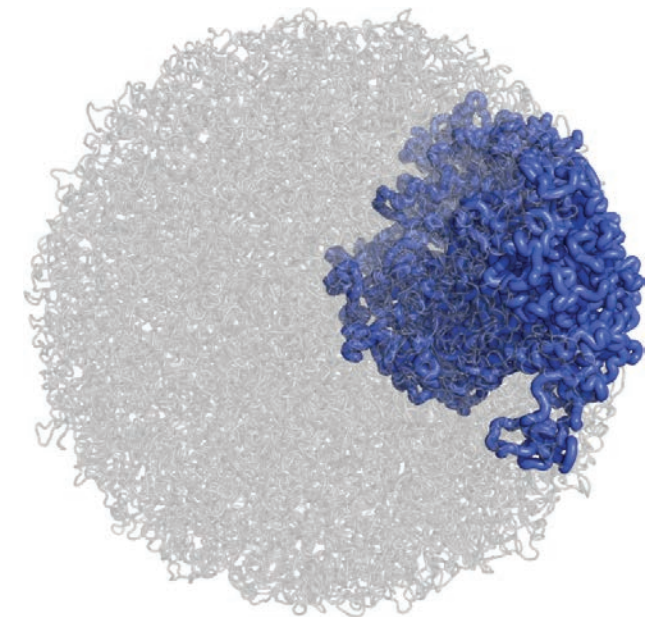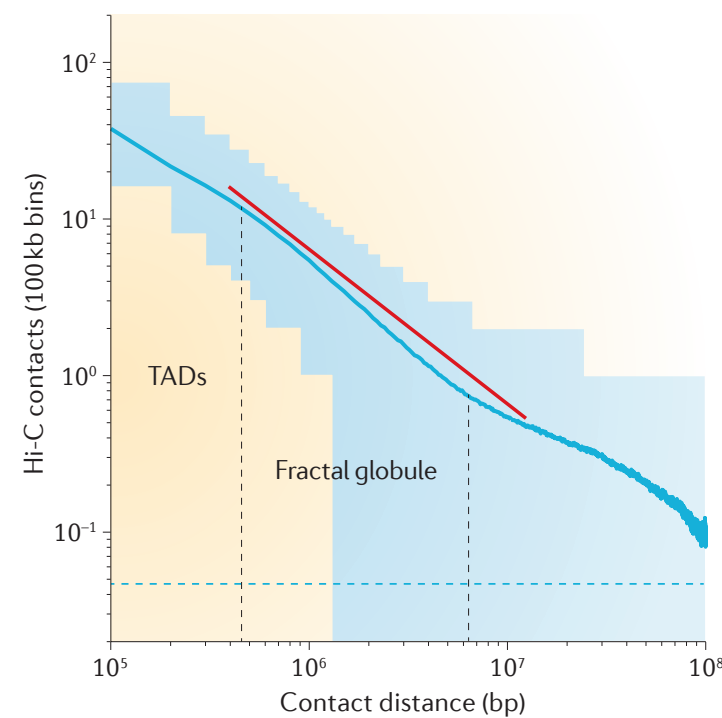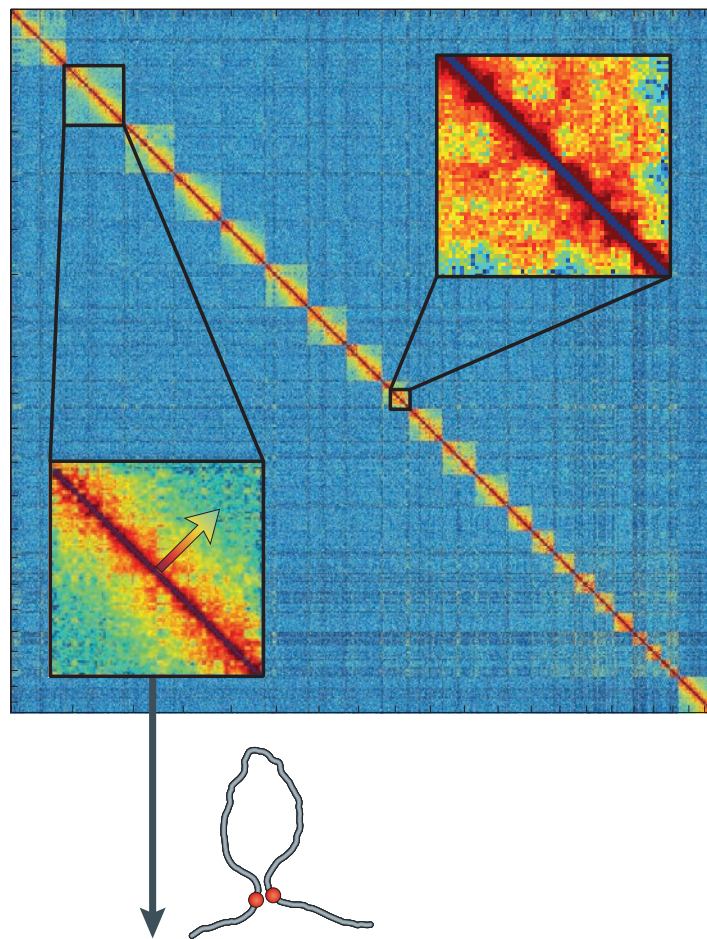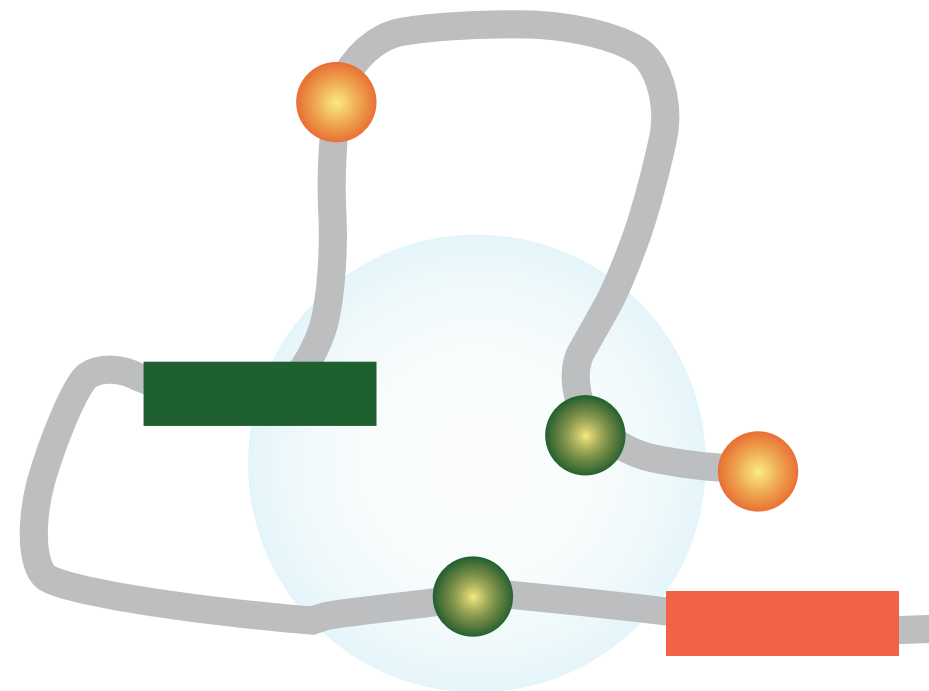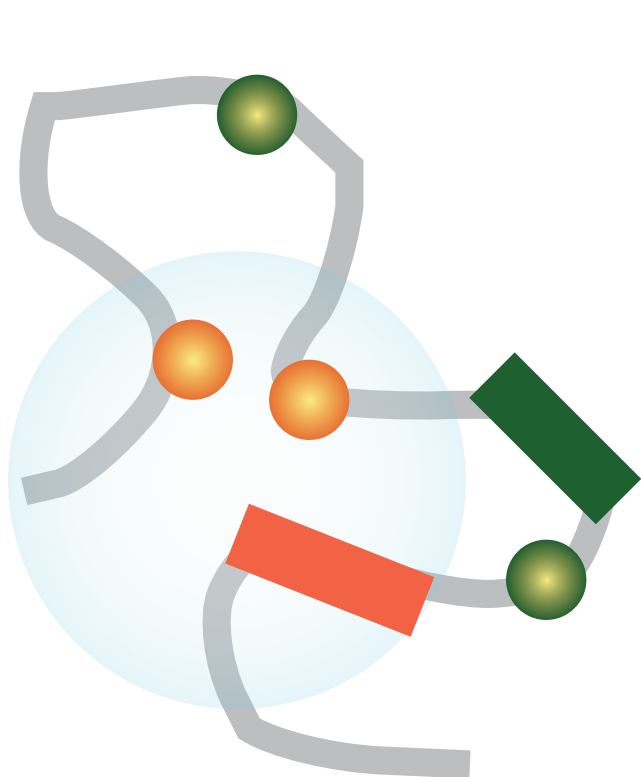
# Complex genome organization

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).

# Complex genome organization

Lieberman-Aiden, E. et al. Comprehensive mapping of long-range interactions reveals folding principles of the human genome. Science (New York, NY) 326, 289–293 (2009).

# Complex genome organization

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

cnag

# Chromosome Conformation Capture

Tuesday, October 15, 13

| | 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.

cnag CRG

# 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

# Representation

## Harmonic

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

## Harmonic Lower Bound

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

## Harmonic Upper Bound

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

# Scoring



GM12878

70 fragments
1,520 restraints

Harmonic          Harmonic Lower Bound          Harmonic Upper Bound

K562

70 fragments
1,049 restraints

# Optimization

# Clustering

# Not just *one* solution



GM12878

K562

# Consistency

# Regulatory elements
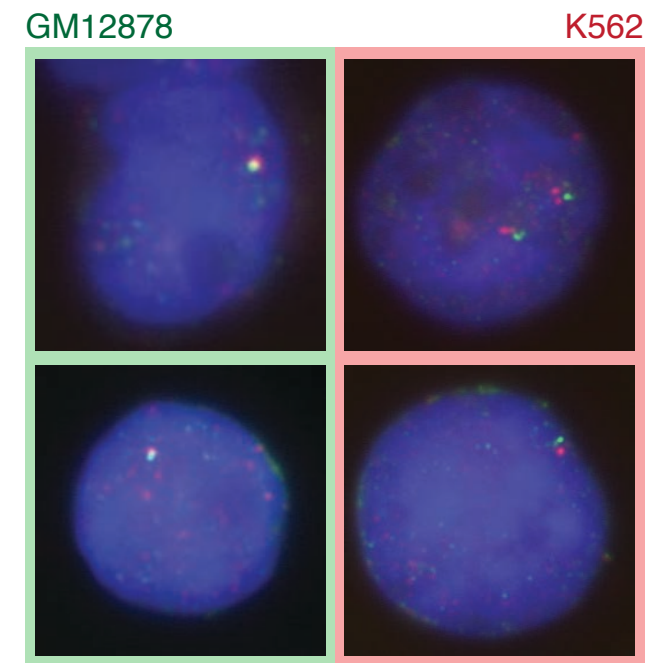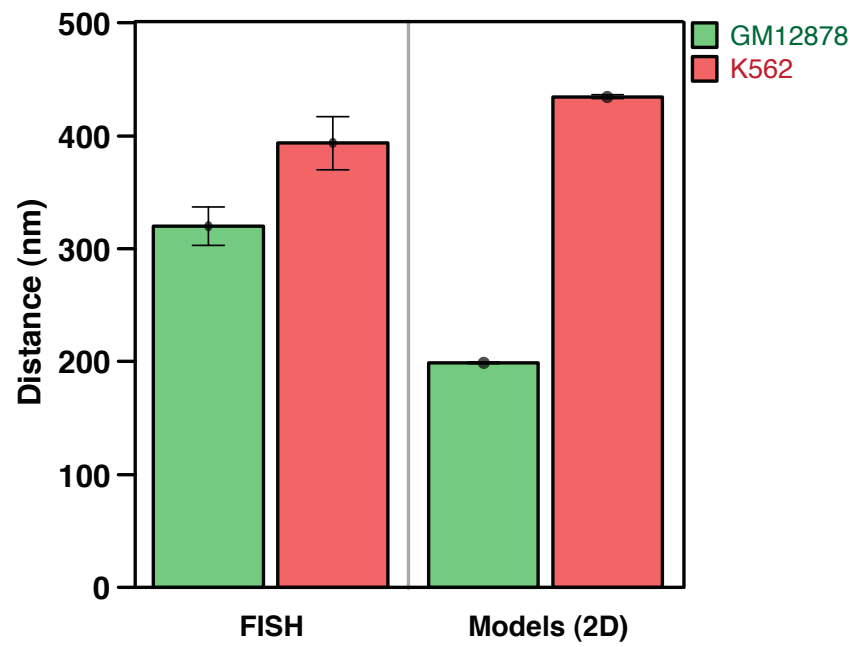
# Compactness
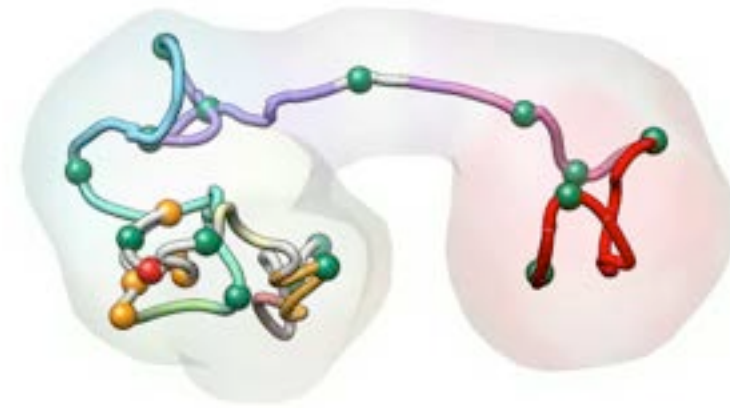
# Multi-loops

# Expression

# FISH validation



GM12878
Cluster #1
2780 model

K562
Cluster #2
314 model

# 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)*

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
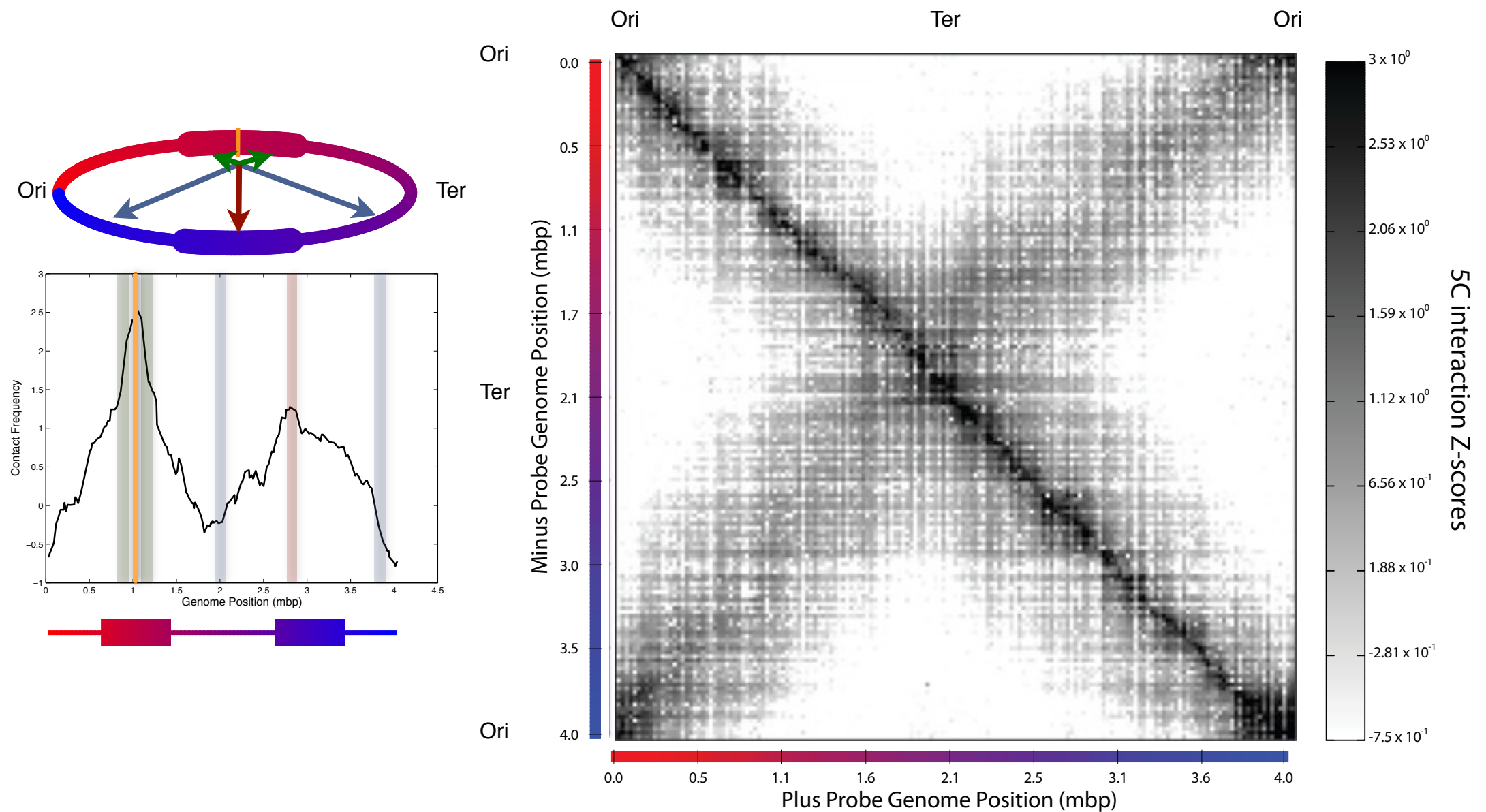
# The 3D architecture of Caulobacter Crescentus

4,016,942 bp & 3,767 genes



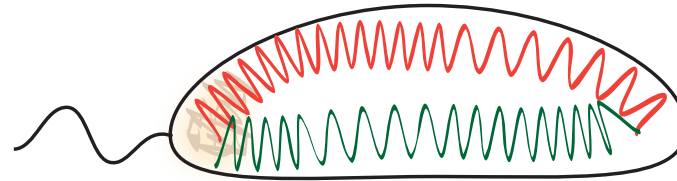Origin

= + Strand
= - Strand

Terminus

169 5C primers on + strand
170 5C primers on – strand
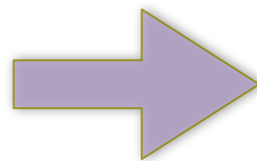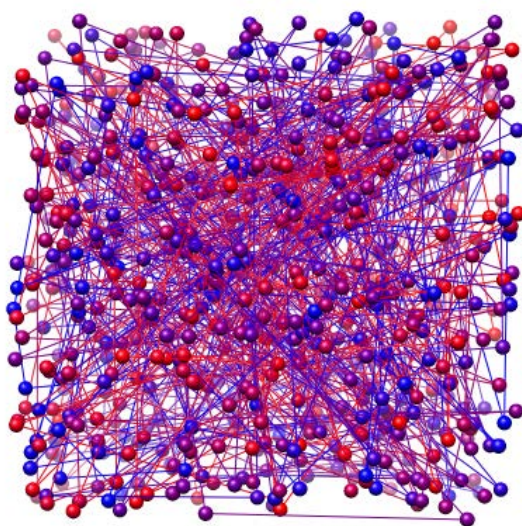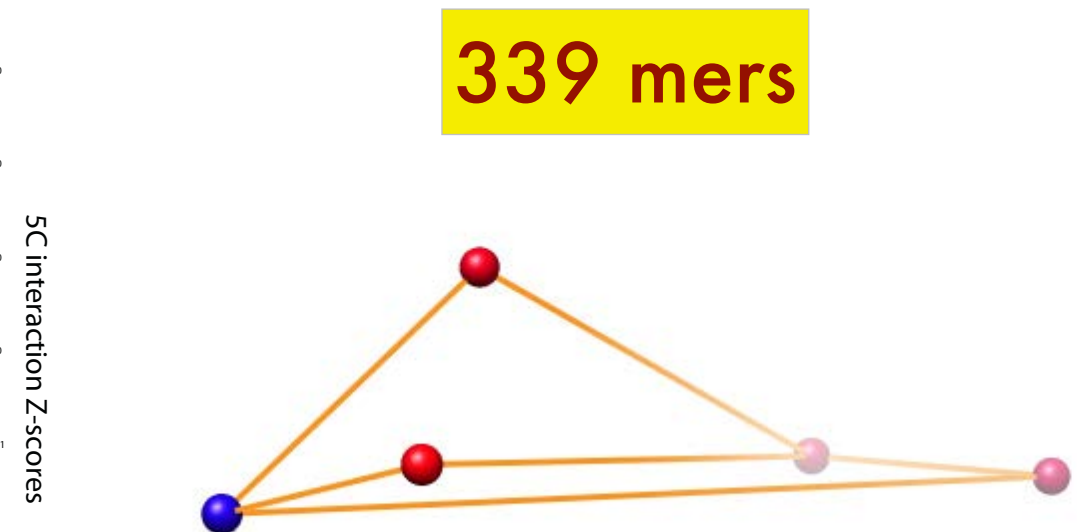28,730 chromatin interactions
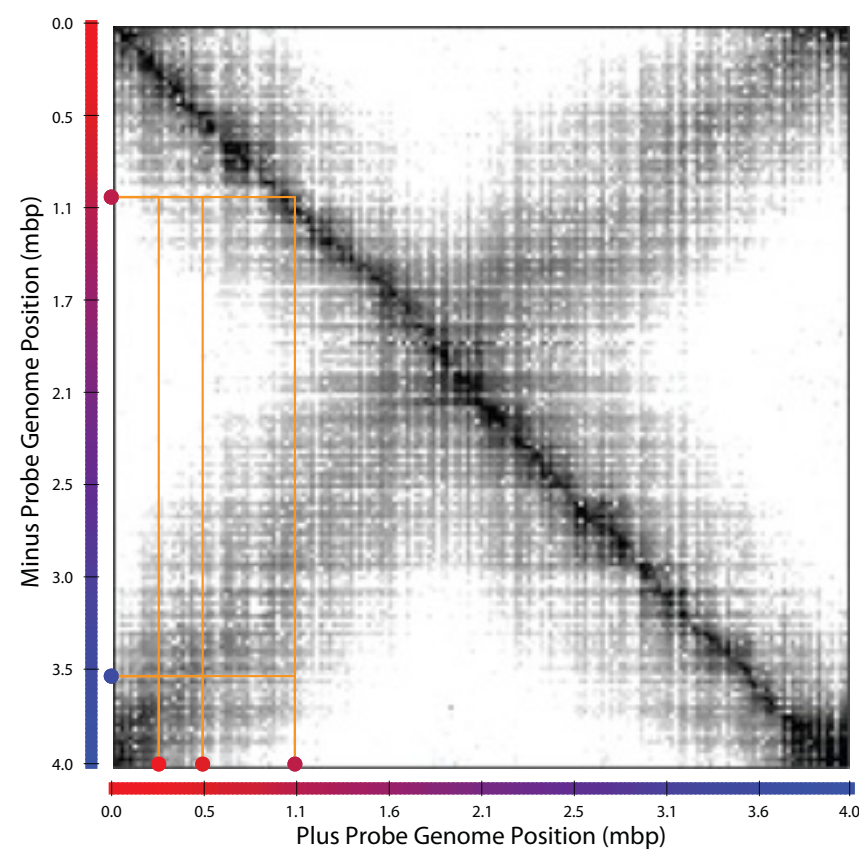
~13Kb

cnag CRG

# 5C interaction matrix

## ELLIPSOID for Caulobacter cresentus

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



339 mers

# Genome organization in Caulobacter crescentus

# Moving the parS sites 400 Kb away from Ori



Wild-type

ET166

# 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.

# 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

# Take home message



Data collection

Data interpretation
Representation
Scoring

Modeling
Sampling

Model analysis