### Data integration for 3D structure determination.

#### Marc A. Marti-Renom

Genome Biology Group (CNAG) Structural Genomics Group (CRG)









## **Structural Genomics Group**

http://www.marciuslab.org





## Integrative Modeling Platform

http://www.integrativemodeling.org



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



# **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).







## anization

Lieberman-Aiden, E. et al. Comprehensive mapping of long-range interactions reveals folding principles of the human genome. Science (New York, NY) 326, 289–298 (2009).













#### **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.



## - some some pture

	3C	5C	4C	Hi-C	ChIP-loop	ChIA-PET
Principle	Contacts between two defined regions <sup>3,17</sup>	All against all <sup>4,18</sup>	All contacts with a point of interest <sup>14</sup>	All against all <sup>10</sup>	Contacts between two defined regions associated with a given protein <sup>8</sup>	All contacts associated with a given protein <sup>6</sup>
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 <sup>7</sup> or melting curve analysis <sup>1</sup>		Circular chromosome conformation capture <sup>20</sup> , open- ended chromosome conformation capture <sup>19</sup> , inverse 3C <sup>12</sup> , associated chromosome trap (ACT) <sup>11</sup> , affinity enrichment of bait- ligated junctions <sup>2</sup>	Yeast <sup>5,15</sup> , tethered conformation capture <sup>9</sup>		ChIA-PET combined 3C-ChIP-cloning (6C), <sup>16</sup> enhanced 4C (e4C) <sup>13</sup>

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



## Modeling 3D Genomes

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



![](_page_18_Picture_3.jpeg)

## Examples...

![](_page_19_Figure_1.jpeg)

![](_page_19_Picture_2.jpeg)

## Human &-globin domain

![](_page_20_Picture_1.jpeg)

![](_page_20_Picture_2.jpeg)

#### Human $\alpha$ -globin domain

ENm008 genomic structure and environment

![](_page_21_Figure_2.jpeg)

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 DNasel 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

![](_page_21_Picture_5.jpeg)

#### Human $\alpha$ -globin domain

ENm008 genomic structure and environment

![](_page_22_Figure_2.jpeg)

![](_page_22_Picture_3.jpeg)

## Representation

![](_page_23_Figure_1.jpeg)

![](_page_23_Picture_2.jpeg)

## Scoring

![](_page_24_Figure_1.jpeg)

![](_page_24_Picture_2.jpeg)

## Optimization

![](_page_25_Figure_1.jpeg)

![](_page_25_Picture_2.jpeg)

## Clustering

![](_page_26_Figure_1.jpeg)

![](_page_26_Picture_2.jpeg)

## Not just one solution

![](_page_27_Figure_1.jpeg)

![](_page_28_Picture_0.jpeg)

![](_page_28_Figure_1.jpeg)

![](_page_28_Picture_2.jpeg)

## The "Chromatin Globule" model

![](_page_29_Figure_1.jpeg)

![](_page_29_Picture_2.jpeg)

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.

![](_page_29_Picture_4.jpeg)

## Caulobacter crescentus genome

![](_page_30_Picture_1.jpeg)

![](_page_30_Picture_2.jpeg)

### The 3D architecture of Caulobacter Crescentus

4,016,942 bp & 3,767 genes

![](_page_31_Figure_2.jpeg)

![](_page_31_Picture_3.jpeg)

#### **5C interaction matrix**

**ELLIPSOID** for Caulobacter cresentus

![](_page_32_Picture_2.jpeg)

![](_page_32_Figure_3.jpeg)

![](_page_32_Picture_4.jpeg)

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

![](_page_33_Figure_1.jpeg)

![](_page_33_Figure_2.jpeg)

![](_page_33_Picture_3.jpeg)

![](_page_33_Picture_4.jpeg)

#### Genome organization in Caulobacter crescentus

![](_page_34_Picture_1.jpeg)

![](_page_34_Picture_2.jpeg)

#### Moving the parS sites 400 Kb away from Ori

![](_page_35_Figure_1.jpeg)

![](_page_35_Picture_2.jpeg)

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

![](_page_36_Figure_1.jpeg)

![](_page_36_Figure_2.jpeg)

#### Genome architecture in Caulobacter

![](_page_37_Picture_1.jpeg)

![](_page_37_Picture_2.jpeg)

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

![](_page_37_Picture_4.jpeg)

#### From Sequence to Function 5C + IMP

![](_page_38_Figure_1.jpeg)

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

![](_page_38_Picture_3.jpeg)

#### Bacteria has also TADs (CIDs)

Le, T. B. K., Imakaev, M. V., Mirny, L. A., & Laub, M. T. (2013). High-Resolution Mapping of the Spatial Organization of a Bacterial Chromosome. Science (New York, NY), 1242059

![](_page_39_Figure_2.jpeg)

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

#### **On TADs and hormones**

![](_page_40_Picture_1.jpeg)

![](_page_40_Picture_2.jpeg)

Davide Baù

![](_page_40_Picture_4.jpeg)

François le Dily

![](_page_40_Picture_6.jpeg)

## Progesterone-regulated transcription in breast cancer

![](_page_41_Figure_1.jpeg)

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?** 

![](_page_41_Picture_5.jpeg)

## Experimental design

![](_page_42_Figure_1.jpeg)

![](_page_42_Picture_2.jpeg)

## Are there TADs? how robust?

![](_page_43_Figure_1.jpeg)

![](_page_43_Picture_2.jpeg)

## Are TADs homogeneous?

![](_page_44_Figure_1.jpeg)

![](_page_44_Picture_2.jpeg)

## **Do TADs respond differently to Pg treatment?**

![](_page_45_Figure_1.jpeg)

![](_page_45_Figure_2.jpeg)

![](_page_45_Picture_3.jpeg)

## Do TADs respond differently to Pg treatment?

![](_page_46_Figure_1.jpeg)

Pg induced fold change per TAD (6h)

![](_page_46_Figure_3.jpeg)

## Modeling 3D TADs

![](_page_47_Figure_1.jpeg)

![](_page_47_Figure_2.jpeg)

61 genomic regions containing 209 TADs covering 267Mb

![](_page_47_Picture_4.jpeg)

## How TADs respond structurally to Pg?

![](_page_48_Figure_1.jpeg)

![](_page_48_Figure_2.jpeg)

![](_page_48_Picture_3.jpeg)

## Model for TAD regulation

![](_page_49_Figure_1.jpeg)

![](_page_49_Picture_2.jpeg)

#### **PLoS CB Outlook**

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

![](_page_50_Picture_2.jpeg)

access to unprecedented details of how genomes organize within the interphase nucleus. Development of new computational approaches to leverage this data has already resulted in the first three-dimensional structures of genomic domains and genomes. Such approaches expand our knowledge of the chromatin folding princiexpand our knowledge of the chromatin foroing princi-ples, which has been classically studied using polymer physics and molecular simulations. Our outlook describes computational approaches for integrating experimental data with polymer physics, thereby bridging the resolu-tion gap for structural determination of genomes and genomic domains.

#### This is an "Editors' Outlook" article for PLoS Computational Biology

Recent experimental and computational advances are resulting in an increasingly accurate and detailed characterization of how genomes are organized in the three-dimensional (3D) space of the nucleus (Figure 1) [1]. At the lowest level of chromatin organization, naked DNA is packed into nucleosomes, which forms the so-called chromatin fiber composed of DNA and proteins. However, this initial packing, which reduces the length of the DNA by about seven times, is not sufficient to explain the higher-order folding of chromosomes during interphase and metaphase. It is now accepted that chromosomes and genes are non-randomly and dynamically positioned in the cell nucleus during the interphase, which challenges the classical representa-tion of genomes as linear static sequences. Moreover, compartmentalization, chromatin organization, and spatial location of genes are associated with gene expression and the functional status of the cell. Despite the importance of 3D genomic architecture. we have a limited understanding of the molecular mechanisms that determine the higher-order organization of genomes and its relation to function. Computational biology plays an important role in the plethora of new technologies aimed at addressing this knowledge gap [2]. Indeed, Thomas Cremer, a pioneer in studying nuclear organization using light microscopy, recently high lighted the importance of computational science in complementing and leveraging experimental observations of genome organization [2]. Therefore, computational approaches to integrate experimental observations with chromatin physics are needed to determine the architecture (3D) and dynamics (4D) of genomes. We present two complementary approaches to address this challenge: (i) the first approach aims at developing simple polymer models of chromatin and determining relevant interactions (both

PLoS Computational Biology | www.ploscompbiol.org

its organization. These approaches are reminiscent of the protein-folding field where the first strategy was used for characterizing protein "foldability" and the second was implemented for modeling the structure of proteins using nuclear magnetic resonance and other experimental constraints. In fact, our outlook consistently returns to the many connections between the two fields.

#### What Does Technology Show Us?

Today, it is possible to quantitatively study structural features of genomes at diverse scales that range from a few specific loci, through chromosomes, to entire genomes (Table 1) [3]. Broadly, there are two main approaches for studying genomic organization light microscopy and cell/molecular biology (Figure 2). Light microcopy [4], both with fixed and living cells, can provide images of a few loci within individual cells [5,6], as well as their dynamics as a function of time [7] and cell state [8]. On a larger scale, light microscopy combined with whole-chromosome staining reveals chromosomal territories during interphase and their reorganization upon cell division. Immunofluorescence with fluorescent antibodies in combination with RNA, and DNA fluorescence in situ hybridization (FISH) has been used to determine the colocalization of loci and nuclear substructures.

Using cellular and molecular biology, novel chromosome conformation capture (3C)-based methods such 3C [9], 3C-onchip or circular 3C (the so-called 4C) [10,11], 3C carbon copy (5C) [12], and Hi-C [13] quantitatively measure frequencies of spatial contacts between genomic loci averaged over a large

Citation: Marti-Renom MA, Mirny LA (2011) Bridging the Resolution Gap in Structural Modeling of 3D Genome Organization. PLoS Comput Biol 7(7): e1002125. doi:10.1371/journal.pcbi.1002125 Editor: Philip E. Bourne, University of California San Diego, United States of America

#### Published July 14, 2011

Copyright: © 2011 Marti-Renom, Mirny. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Funding: MAM-R acknowledges support from the Spanish Ministry of Science and Innovation (BFU2010-19310). LM is acknowledging support of the NCH-funded MIT Center for Physics Sciences in Oncology. The funders had no role in decision to publish, or preparation of the manuscript. Competing Interests: The authors have declared that no competing interests evist

\* E-mail: mmarti@cipf.es

July 2011 | Volume 7 | Issue 7 | e1002125

![](_page_50_Figure_16.jpeg)

000

0

0

![](_page_50_Picture_17.jpeg)

#### Take home message

![](_page_51_Figure_1.jpeg)

![](_page_51_Picture_2.jpeg)

## Acknowledgments

![](_page_52_Picture_1.jpeg)

![](_page_52_Picture_2.jpeg)

http://marciuslab.org http://integrativemodeling.org http://cnag.cat · http://crg.cat