#### **Course outline**

**Theory Practice** 

- Day 1 Introduction to structure determination Chromatin structure and Hi-C data Introduction to linux and python (FACULTATIVE) The Integrative Modeling Platform and Chimera
- Day 2 The Integrative Modeling Platform applied to chromatin TADbit introduction and installation Topologically Associated Domains detection and analysis

Day 3 The TADbit documentation: examples and code snippets 3D modeling of real Hi-C data Analysis of the results



#### **3D structure determination**

Davide Baù & François Serra

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



#### **Structural Genomics Group**

http://www.marciuslab.org





## Data groups



Experimental observations





Statistical rules



Laws of physics



# The importance of the 3D structure

The biochemical function of a molecule is defined by its interactions

The biological function is in large part a consequence of these interations

The 3D structure is more informative than sequence alone

Evolution tends to conserve function and function depends more directly on structure than on sequence







### Structure prediction vs determination



Thursday, April 23, 2009



#### Data integration









# The four stages of integrative modeling

Stage 1: Gathering experimental and statistical Information

**Stage 2: Choosing How To Represent And Evaluate Models** 

**Stage 3: Finding Models That Score Well** 

**Stage 4: Analyzing Resulting Models and Information** 











# Advantages of integrative modeling

- It facilitates the use of new information
- It maximizes accuracy, precision and completeness of the models
- It facilitates assessing the input information and output models
- It helps in understanding and assessing experimental accuracy

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





# Integrative Modeling Platform

http://www.integrativemodeling.org



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











## The simulating annealing procedure





### En example of nergy optimization







# Integrative Modeling Platform

http://www.integrativemodeling.org



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







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







#### PROTEINS



#### COMPLEXES



#### GENOMES



### Proteins

Single data type



**Amino Acids** 









#### **Complexes** Multiple data types





#### S. cerevisiae ribosome



Fitting of comparative models into 15Å cryo-electron density map.

43 proteins could be modeled on 20-56% seq.id. to a known structure.

The modeled fraction of the proteins ranges from 34-99%.

C. Spahn, R. Beckmann, N. Eswar, P. Penczek, A. Sali, G. Blobel, J. Frank. Cell 107, 361-372, 2001.



### The nuclear pore complex





#### Integrative Modeling of the NPC

F. Alber et al. Natute (2007) Vol 450





Representation

 $\theta$ 

436 proteins!

τ	$N^1_{ au}$	$N_{\tau}^2$	К	$\{B_j^\kappa\}$	n <sub>ĸ</sub>	r	τ	$N_{\tau}^{1}$	$N_{\tau}^2$	К	$\{B_j^\kappa\}$	$n_{\kappa}$	r
Nup192	1	1	1,2,5	00	2	3.0	Nup1	0	1	1,5	000000000	9	1.5
			3	-	1	-				2	<b>00</b> 0000000	2	1.5
Nup188	1	1	1,2,5	99	2	3.0				3	-	1	-
			3	-	1	-				4	ംക്കാരം	7	1.5
Nup170	1	1	1,2,5	99	2	2.9	Nsp1	2	2	1,5	*****	12	1.3
			3	-	1	-				2		3	1.3
Nup157	1	1	1,2,5	889	3	2.5				3	-	1	-
			3	-	1	-				4		9	1.3
Nup133	1	1	1,2,5	<b></b>	2	2.7	Gle1	1	0	1,2,5	<b></b>	2	2.1
			3	-	1	-				3	-	1	-
Nup120	1	1	1,2,5	<b></b>	2	2.6	Nup60	0	1	1,5	<b></b>	4	1.6
			3	-	1	-				2,3	<b>0</b> 000	1	1.6
Nup85	1	1	1,2,5	000	3	2.0				4	ಾಲ	3	1.6
			3	-	1	-	Nup59	1	1	1,5		4	1.6
Nup84	1	1	1,2,5		3	2.0				2		2	1.6
			3	-	1	-				3	-	1	-
Nup145C	1	1	1,2,5	<b></b>	2	2.3				4	<b>00</b> 00	2	1.6
•			3	-	1	-				1,5	888	3	1.8
Seh1	1	1	1,2,3,5	9	1	2.2	Nup57	1	1	2,3		1	1.8
Sec13	1	1	1,2,3,5	٩	1	2.1				4	<b>ee</b>	2	1.8
Gle2	1	1	1,2,3,5	٢	1	2.3	Nup53	1	1	1,5	<b></b>	3	1.7
Nic96	2	2	1,2,5	<b>33</b>	2	2.4				2,3	000	1	1.7
			3	-	1	-				4	<b>99</b>	2	1.7
Nup82	1	1	1,2,5	<b></b>	2	2.3	Nup145N	0	2	1,5	333333	6	1.5
			3	-	1	-				2,3	000000	1	1.5

Alber, F., Dokudovskaya, S., Veenhoff, L. M., Zhang, W., Kipper, J., Devos, D., Suprapto, A., et al. (2007). Nature, 450(7170), 695–701



K



Data generation		Data interpretation							
Method	Experiments	Restraint	R <sub>c</sub>	Ro	R <sub>A</sub>	Functional form of activated feature restraint			
fractionation	30 nup sequences	Protein excluded volume restraint	-	-	1,864 1,863/2	Protein-protein:   Violated for $f < f_o$ . $f$ is the distance between two beads, $f_o$ is the sum of the bead radii, and $\sigma$ is 0.01 nm.   Applied to all pairs of particles in representation $\kappa$ =1: $B^{m} = \left\{ B_j^{m-1}(\theta, s, \tau, i) \right\}$			
Bioinformatics and Membrane	30 nup sequences	Surface localization restraint	-	-	48	$\begin{array}{l} \textbf{Membrane-surface location:}\\ \textbf{Violated if } f \neq f_{o}, f \text{ is the distance between a protein particle and the closest point on the NE surface (half-torus), f_{o} = 0 nm, and \sigma \text{ is } 0.2 nm. Applied to particles:}\\ B^{m} = \left\{B_{j}^{r-6}(\theta, s, \tau, i) \mid \tau \in (\text{Ndcl}, \text{Poml52}, \text{Pom34})\right\} \end{array}$			
	30 Nup sequences and immuno-EM (see below)		-	-	64	$\label{eq:pore-side volume location:} \begin{aligned} & \text{Pore-side volume location:} \\ & \text{Violated if } f < f_o, f \text{ is the distance between a protein particle and the closest point on the} \\ & \text{NE surface (half-torus), } f_o = 0 \text{ nm, and } \sigma \text{ is } 0.2 \text{ nm. Applied to particles:} \\ & B^{m} = \left\{ B_j^{r-s}(\theta,s,\tau,i) \mid r \in (\text{Ndc1,Pom152,Pom34}) \right\} \end{aligned}$			
			-	-	80	$\label{eq:period} \begin{array}{l} \textbf{Periouclear volume location:}\\ \textbf{Violated if } f > f_{o}, f \text{ is the distance between a protein particle and the closest point on the NE surface (half-forus), f_{o} = 0 nm, and \sigma \text{ is } 0.2 nm. Applied to particles:}\\ B^{see} = \left\{B_{j}^{s-7}(\theta,s,\tau,i)\tau\in(\text{Pom152})\right\} \end{array}$			
Hydrodynamics experiments	1 S-value	Complex shape restraint	1	164	1	$\label{eq:complex_diameter} \begin{array}{l} \textbf{Complex_diameter} \\ \mbox{Violated if } f < f_o. f is the distance between two protein particles representing the largest diameter of the largest complex, f_o is the complex maximal diameter D=19.2-R, where R is the sum of both particle radii, and \sigma is 0.01 nm. Applied to particles of proteins in composite C45: B^{me} = \left\{ B_j^{n-1}(\theta,s,\tau,i) \mid \tau \in C_{s1} \right\}$			
	30 S-values	Protein chain restraint	-	-	1,680	Protein chain   Violated if $f \neq f_o$ . $f$ is the distance between two consecutive particles in a protein, $f_o$ is the sum of the particle radii, and $\sigma$ is 0.01 nm. Applied to particles: $B = \left\{ B_j^{\kappa}(\theta, s, \tau, i)   \kappa = 1 \right\}$			
Immuno-Electron microscopy		Protein localization restraint	-		456	<b>Z-axial position</b> Violated for $f < f_o$ , $f$ is the absolute Cartesian Z-coordinate of a protein particle, $f_o$ is the lower bound defined for protein type $r$ , and $\sigma$ is 0.1 nm. Applied to particles: $B = \left\{ B_j^c (\theta, s, r, i)   \kappa = l, j = l \right\}$			
	particles				456	Violated for $f > f_o$ , $f$ is the absolute Cartesian Z-coordinate of a protein particle, $f_o$ is the upper bound defined for protein type $r$ , and $\sigma$ is 0.1 nm. Applied to particles: $B = \left\{ B_i^r (\theta, s, \tau, i)   \kappa = 1, j = 1 \right\}$			
	10,940 gold			-	456	<b>Radial position</b> Violated for $f < f_o$ . <i>f</i> is the radial distance between a protein particle and the Z-axis in a plane parallel to the X and Y axes, $f_o$ is its lower bound defined for protein type $\tau$ , and $\sigma$ is 0.1 nm. Applied to particles: $B = \{B_r^c(\theta, s, \tau, i)   \kappa = 1, j = 1\}$			
	· ·				456	Violated for $f > f_o$ . <i>f</i> is the radial distance between a protein particle and the Z-axis in a plane parallel to the X and Y axes, $f_o$ is its upper bound defined for protein type $r$ , and $\sigma$ is 0.1 nm. Applied to particles: $B = \left\{ B_i^r \left( \theta, s, \tau, i \right)   \kappa = 1, j = 1 \right\}$			
Overlay assays	13 contacts	Protein interaction restraint	20	112	20	Protein contact   Violated for $f > f_o$ . $f$ is the distance between two protein particles, $f_o$ is the sum of the particle radii multiplied by a tolerance factor of 1.3, and $\sigma$ is 0.01 nm. Applied to particle: $B = \{B_j^{\kappa}(\theta, s, \tau, i)   \kappa \in (2, 4, 9), \theta \in (1, 2, 3)\}$			
Affinity purification	4 complexes	Competitive binding restraint	1	132	4	Protein contact   Violated for $f > f_o$ . $f$ is the distance between two protein particles, $f_o$ is the sum of the particle radii multiplied by a tolerance factor of 1.3, and $\sigma$ is 0.01 nm. Applied to : $B = \left\{ B_j^r(\theta, s, \tau, i)   \theta \in (1, 2, 3), \kappa \in (2, 4, 6), \tau = (Nup 82, Nic 96, Nup 49, Nup 57) \right\}$			
	64 complexes	Protein proximity restraint	692	25,348	692	<b>Protein proximity</b> Violated for $f > f_o$ . $f$ is the distance between two protein particles, $f_o$ is the maximal diameter of a composite complex, and $\sigma$ is 0.01 nm. Applied to particles: $B = \left\{ B_j^{\kappa}(\theta, s, \tau, i)   \theta \in (1, 2, 3), \kappa \in (2, 4, 9) \right\}$			



#### Optimization







## The structure of the nuclear pore complex



www.nature.com/nature



#### Genomes

Limited data types







#### Main approaches

#### Light microscopy (FISH)



#### Cell/molecular biology (3C-based methods)





## The resolution gap







#### Simple genomes



#### Complex genomes











#### Job Dekker



Dostie et al. Genome Res (2006) vol. 16 (10) pp. 1299-309



#### **3C-like technologies**



Hakim and Misteli Cell (2012) vol. 148, March 2



### **3C-like technologies**

	3C 5C		4C Hi-C		ChIP-loop	ChIA-PET	
Principle	Contacts between two defined regions3,17All against all4,18		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 Commonly < 1Mb Genome-		Genome-wide	Commonly < 1Mb	Genome-wide	
Detection	Locus-specific PCR	Locus-specific PCR HT-sequencing		HT-sequencing	Locus-specific qPCR	HT-sequencing	
Limitations	Low throughput and coverage	Limited coverage	Limited to one viewpoint		Rely on one chromatin-a disregarding other conta	ssociated factor, cts	
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 and Misteli Cell (2012) vol. 148, March 2



# Take home message



