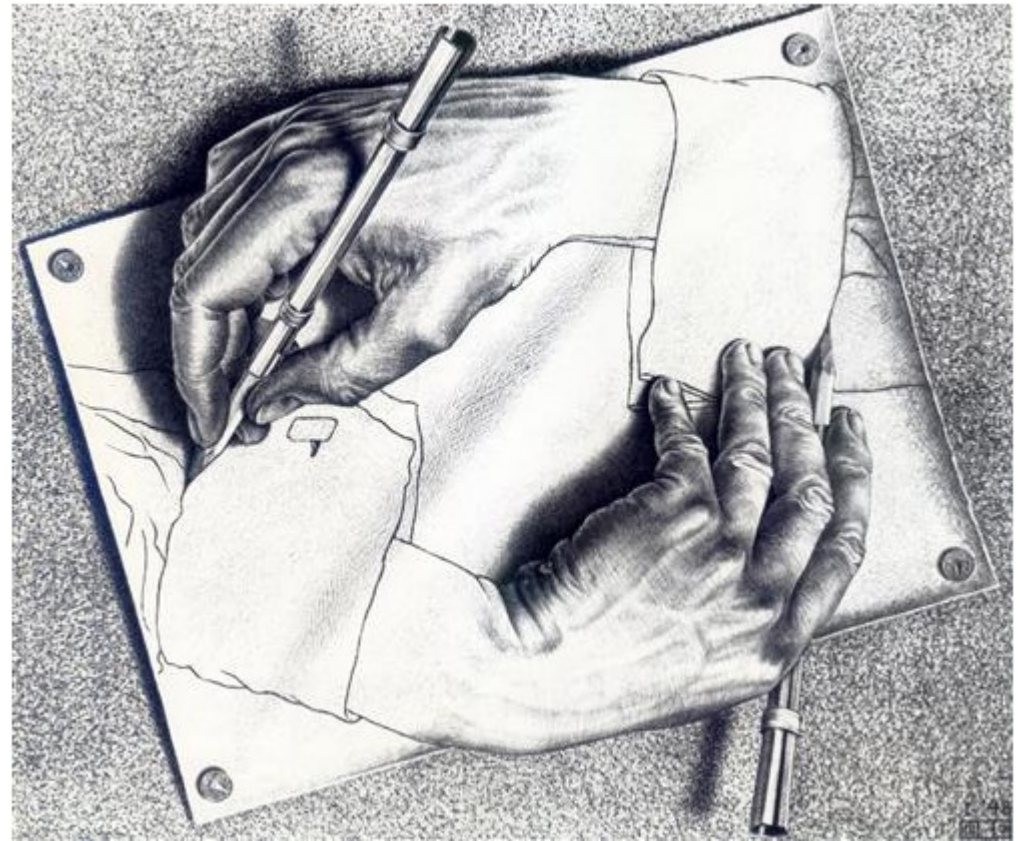


# Hi-C data analysis - Exploring the 3D structure of the chromatin by processing DNA sequences



# Outline

- ◆ Biological context
- ◆ More biological context
- ◆ Hi-C data processing
  - ◆ map
  - ◆ filter
  - ◆ count
  - ◆ normalize
  - ◆ segment
  - ◆ compare
- ◆ Conclusion, discussion, NETBIO lunch

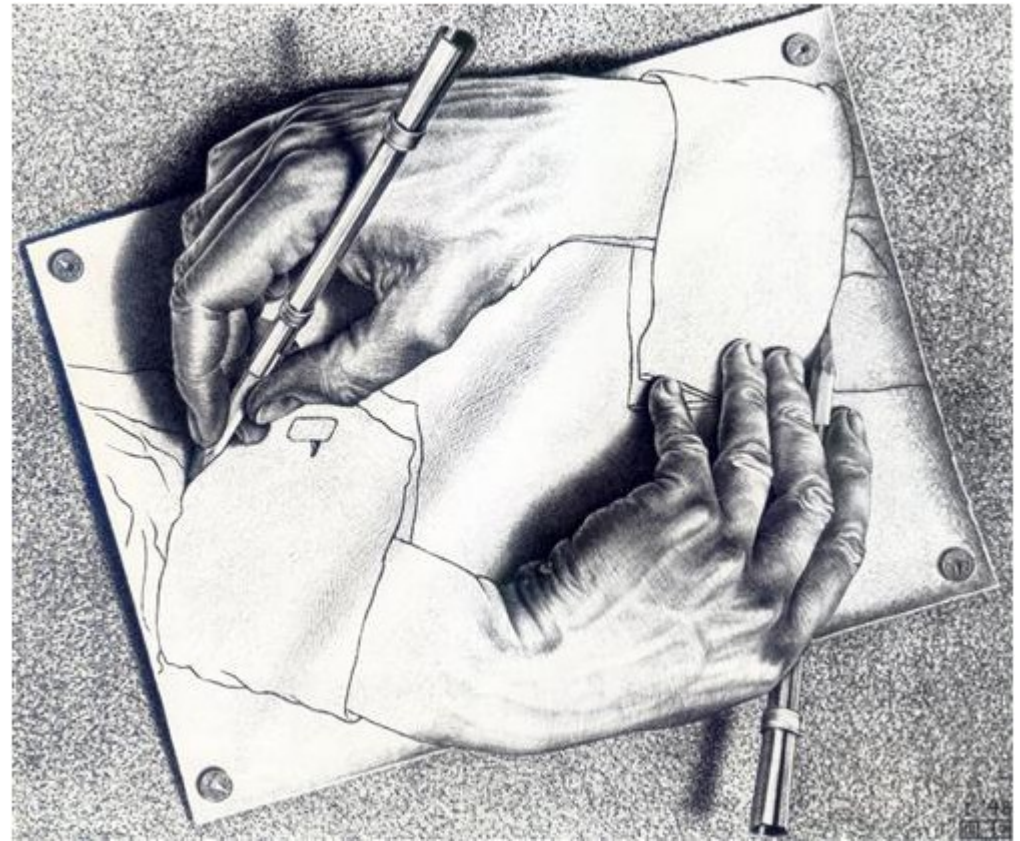


*M.C. Escher, 1948*

# Outline

## Biological context

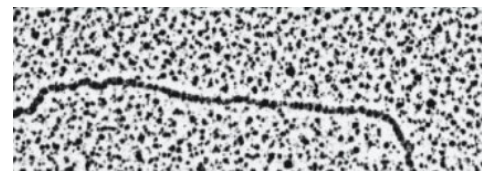
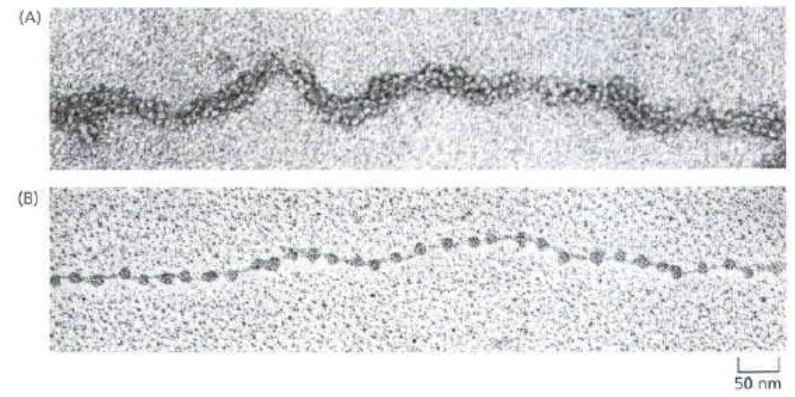
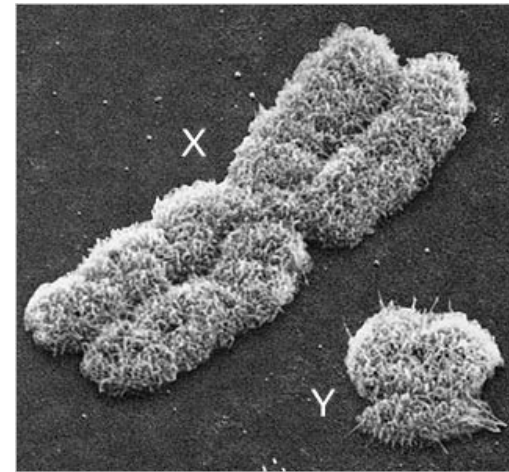
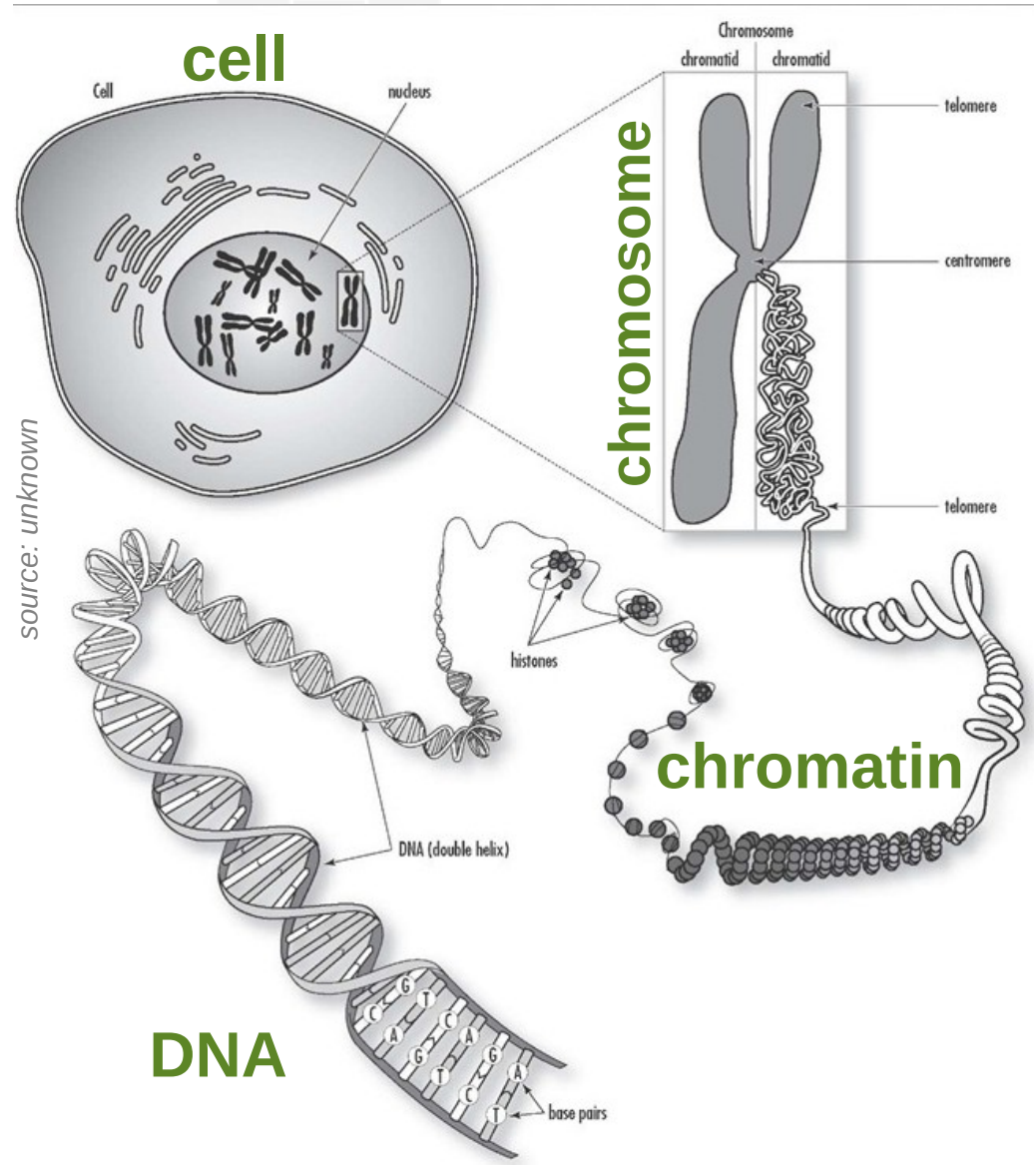
- ◆ More biological context
- ◆ Hi-C data processing
  - ◆ map
  - ◆ filter
  - ◆ count
  - ◆ normalize
  - ◆ segment
  - ◆ compare
- ◆ Conclusion, discussion, NETBIO lunch



*M.C. Escher, 1948*



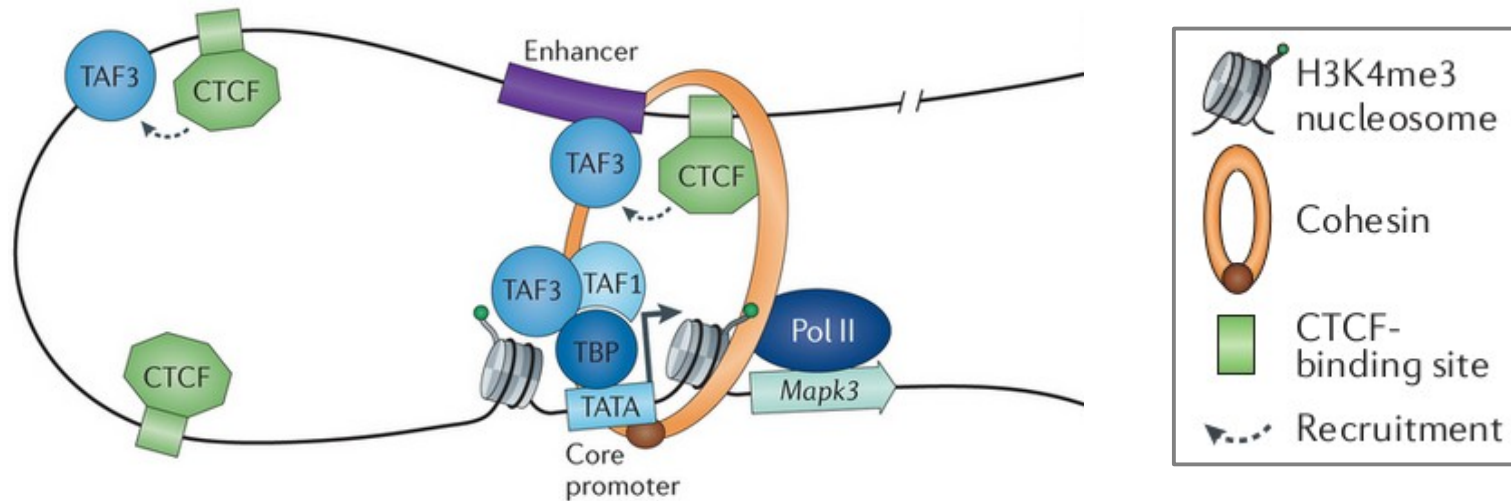
# Life, cell, chromosome & DNA



# Life, cell, chromosome & DNA



# From structure to function



*Ong & Corces, Nat. Rev. Genet., 2014*

**3D structure impacts gene expression**



# From structure to function

## Chromosomal Contact Permits Transcription between Coregulated Genes

Stephanie Fanucchi,<sup>1</sup> Youtaro Shibayama,<sup>1</sup> Shaun Burd,<sup>1</sup> Marc S. Weinberg,<sup>3,4</sup> and Musa M. Mhlanga<sup>1,2,\*</sup>

<sup>1</sup>Gene Expression and Biophysics Group, Synthetic Biology Emerging Research Area, Biosciences Unit, Council for Scientific and Industrial Research, Pretoria, Gauteng 0001, South Africa

<sup>2</sup>Unidade de Biologia e Expressão Genética, Instituto de Medicina Molecular, Faculdade de Medicina, Universidade de Lisboa, Lisboa, 1649-028 Portugal

<sup>3</sup>Antiviral Gene Therapy Research Unit, Department of Molecular Medicine and Haematology, School of Pathology, Faculty of Health Sciences, University of the Witwatersrand, Johannesburg, Gauteng 2193, South Africa

<sup>4</sup>Department of Molecular and Experimental Medicine, The Scripps Research Institute, La Jolla, CA 92037, USA

\*Correspondence: [yoda@mhlangalab.org](mailto:yoda@mhlangalab.org)

<http://dx.doi.org/10.1016/j.cell.2013.09.051>

### SUMMARY

Transcription of coregulated genes occurs in the context of long-range chromosomal contacts that

Nucleic Acids Research Advance Access published February 4, 2015

*Nucleic Acids Research* 2015, 43, 1093-1104  
doi:10.1093/nar/gkv046

2012). These highly sensitive assays capture nascent mRNA and have revealed the FISH foci in a fraction of the population (2010; Papantonis et al., 2010). This study

## Spatial re-organization of myogenic regulatory sequences temporally controls gene expression

Akihito Harada<sup>1</sup>, Chandrashekara Mallappa<sup>2</sup>, Seiji Okada<sup>1</sup>, John T. Butler<sup>2</sup>, P. Baker<sup>2,3</sup>, Jeanne B. Lawrence<sup>2</sup>, Yasuyuki Ohkawa<sup>1,2,\*</sup> and Anthony N. Im

<sup>1</sup>Department of Advanced Medical Initiatives, JST-CREST, Faculty of Medicine, Kyushu University, 812-8582, Japan, <sup>2</sup>Department of Cell and Developmental Biology, University of Massachusetts Lowell, Lowell, MA 01854, USA, <sup>3</sup>Department of Biology, University of Massachusetts Lowell, Lowell, MA 01854, USA

\*Correspondence: [yoda@mhlangalab.org](mailto:yoda@mhlangalab.org) and K. Paszkiewicz and the Exeter Sequencing Service facility for genome sequencing services. This work was also supported by the Wellcome Trust (grant number 098346/Z/05/A) and the European Union (FP7-240622).

### TRANSCRIPTION

## CTCF establishes discrete functional chromatin domains at the *Hox* clusters during differentiation

Varun Narendra,<sup>1,2</sup> Pedro P. Rocha,<sup>3</sup> Disi An,<sup>4</sup> Ramya Ravindran,<sup>1,2</sup> Esteban O. Mazzoni,<sup>4,\*</sup> Danny Reinberg<sup>1,2,\*</sup>

Polycomb and Trithorax group proteins encode the epigenetic identity by establishing inheritable domains of repressive and active chromatin within the *Hox* clusters. Here we demonstrate that the CCCTC-binding factor (CTCF) functions as

## Nuclear Aggregation of Olfactory Receptor Genes Governs Their Monogenic Expression

E. Josephine Clowney,<sup>1</sup> Mark A. LeGros,<sup>2,4</sup> Colleen P. Eirene C. Markenskoff-Papadimitriou,<sup>3</sup> Markko Myllys, and Stavros Lomvardas<sup>1,2,3,\*</sup>

<sup>1</sup>Program in Biomedical Sciences

<sup>2</sup>Physical Biosciences Division, Lawrence Berkeley National Laboratory



GENOME RESEARCH

Next Generation  
USA Congress  
27 - 28 October 2015

HOME | ABOUT | ARCHIVE | SUBMIT | SUBSCRIBE | ADVERTISE | AUTHOR INFO | CONTACT

Institution: SWETS SUBSCRIPTIONSERVICE Sign In via User Name

## Spatial enhancer clustering and regulation of enhancer-proximal genes by cohesin

Elizabeth Ing-Simmons<sup>1,2,7</sup>, Vlad C. Seitan<sup>1,7</sup>, Andre J. Faure<sup>3,8</sup>, Paul Flicek<sup>3,4</sup>,

## Disruptions of Topological Chromatin Domains Cause Pathogenic Rewiring of Gene-Enhancer Interactions

Daño G. Lupiáñez,<sup>1,2</sup> Katerina Kraft,<sup>1,2</sup> Verena Heinrich,<sup>2</sup> Peter Krawitz,<sup>1,2</sup> Francesco Brancati,<sup>3</sup> Eva Kl Denise Horn,<sup>2</sup> Hülya Kayserli,<sup>5</sup> John M. Opitz,<sup>6</sup> Renata Laxova,<sup>6</sup> Fernando Santos-Simarro,<sup>7,8</sup> Brigitte Gilbert-Dussardier,<sup>9</sup> Lars Wittler,<sup>10</sup> Marina Borschiwer,<sup>1</sup> Stefan A. Haas,<sup>11</sup> Marco Osterwalder,<sup>12</sup> Bernd Timmermann,<sup>13</sup> Jochen Hecht,<sup>1,14</sup> Malte Spielmann,<sup>1,2,14</sup> Axel Visel,<sup>12,15,16</sup> and Stefan Mundlos<sup>1</sup>

<sup>1</sup>Max Planck Institute for Molecular Genetics, RG Development & Disease, 14195 Berlin, Germany

<sup>2</sup>Institute for Medical and Human Genetics, Charité Universitätsmedizin Berlin, 13353 Berlin, Germany

<sup>3</sup>Medical Genetics Unit, Policlinico Tor Vergata University Hospital, 00133 Rome, Italy

Leading Edge  
Previews

### A CRISPR Connection between Chromatin Topology and Genetic Disorders

Bing Ren<sup>1,\*</sup> and Jesse R. Dixon<sup>1</sup>

<sup>1</sup>Ludwig Institute for Cancer Research, University of California, San Diego, School of Medicine, 9500 Gilman Drive, La Jolla, CA 92093-0653

\*Correspondence: [binren@ucsd.edu](mailto:binren@ucsd.edu)

<http://dx.doi.org/10.1016/j.cell.2015.04.047>

Structural variations are common in the human genome, but their contributions to human diseases are poorly understood. Lupiáñez et al. demonstrate that some structural variants can interrupt resulting in ectopic enhancer-promoter interactions, altered spatiotemporal patterns, and developmental disorders.

as insertions, deletions, duplications, inversions, translocations, and inversions of DNA segments, are suggesting that they are stable during development and are not easily disrupted by transcriptional activities of the cell. In the human cases (Figure 1), remarkably, mutant mice carrying these structural alterations accurately reproduce

## 3D structure impacts gene expression

# From structure to function

**Chromosomal Contact Permits Transcription between Coregulated Genes**

Stephanie Fanucchi,<sup>1</sup> Younsoo Shin,<sup>1,2</sup> Eshwar Bhardwaj,<sup>1,2</sup> Marc S. Weinberg,<sup>1,2</sup> and Maya M. Malkina,<sup>1,2</sup> <sup>1</sup>Genome Division and Division of Chromatin Biology, Center for Genome Sciences and Policy, Duke University, Durham, NC 27708, USA; <sup>2</sup>Department of Molecular Biology, Princeton University, Princeton, NJ 08542, USA

**SUMMARY** Transcription of coregulated genes occurs in the same nuclear compartment, and this compartment is enriched for transcription factors and coactivators. Here, we show that transcription of coregulated genes occurs in the same nuclear compartment, and this compartment is enriched for transcription factors and coactivators. Here, we show that transcription of coregulated genes occurs in the same nuclear compartment, and this compartment is enriched for transcription factors and coactivators.

**Spatial enhancer clustering and regulation of enhancer-proximal genes by cohesin**

Elizabeth Ing-Simmons<sup>1,2,7</sup>, Vlad C. Seitan<sup>1,7</sup>, Andre J. Fauro<sup>3,8</sup>, Paul Filice<sup>3,4</sup>,

**Disruptions of Topological Chromatin Domains Cause Pathogenic Rewiring of Gene-Enhancer Interactions**

Daño G. Lupáñez,<sup>1,2</sup> Katerina Kraft,<sup>1,2</sup> Verena Heinrich,<sup>1</sup> Peter Krawitz,<sup>1,2</sup> Francesco Brancati,<sup>3</sup> Eva K. Dennis-Horn,<sup>4</sup> Hülya Kayserili,<sup>5</sup> John M. Cople,<sup>6</sup> Ramona Laxova,<sup>7</sup> Fernando Santos-Simarro,<sup>1,2</sup> Brigitte Gilbert-Dussardier,<sup>8</sup> Lars Wittler,<sup>9</sup> Marina Borchow,<sup>10</sup> Stefan A. Haas,<sup>11</sup> Marco Caterwilder,<sup>12</sup> Bernd Timmermann,<sup>13</sup> Jochen Hecht,<sup>1,14</sup> Malte Spielmann,<sup>1,15</sup> Axel Visel,<sup>1,16,17</sup> and Stefan Mundlos<sup>1,18</sup>

**Spatial re-organization of myogenic regulatory sequences temporally controls gene expression**

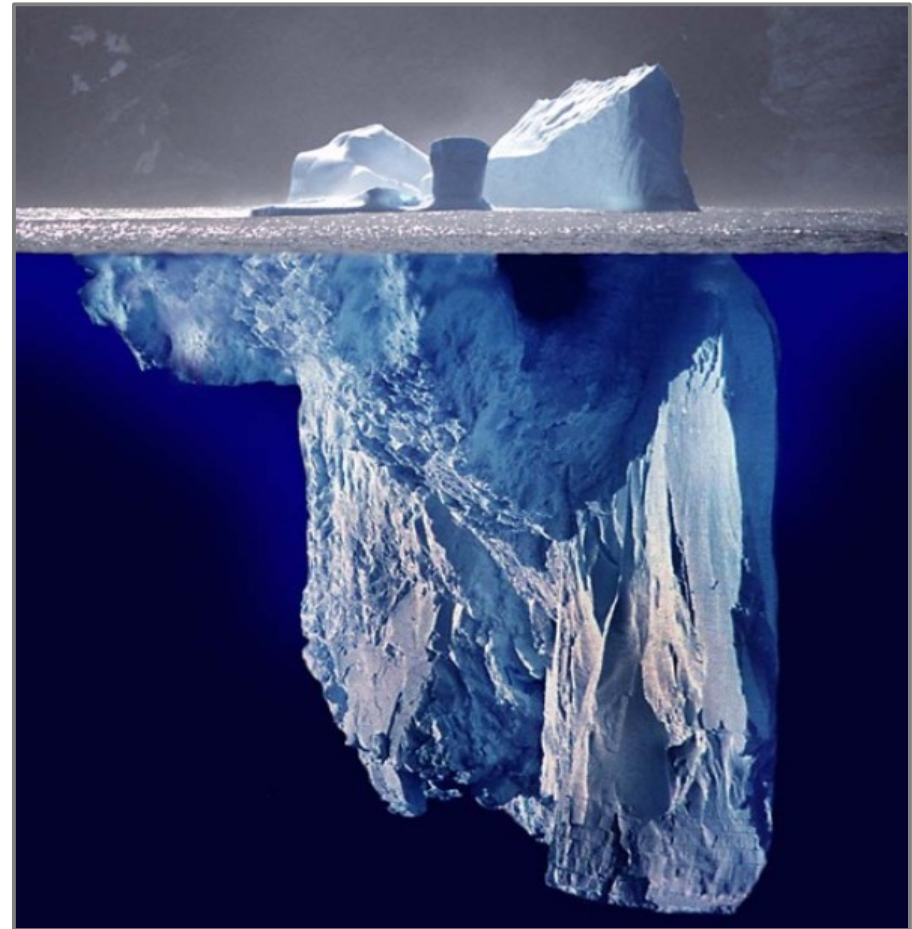
Akihito Harada,<sup>1</sup> Chandrasekhar Mallappa,<sup>2</sup> Saji Okada,<sup>3</sup> John T. Butler,<sup>4</sup> P. Baker,<sup>5</sup> Jeanne B. Lawrence,<sup>6</sup> Yasuyuki Ohkawara,<sup>7</sup> and Anthony M. Ingham<sup>8</sup>

**Nuclear Aggregation of Olfactory Receptor Genes Governs Their Monogenic Expression**

E. Josephine Clowney,<sup>1</sup> Mark A. LeGros,<sup>2,4</sup> Colleen P. Ereno,<sup>5</sup> C. Markoski-Papadimitriou,<sup>3</sup> Markko Mylly, and Stavros Lomvardas<sup>1,3,4</sup>

**CTCF establishes discrete functional chromatin domains at the *Hox* clusters during differentiation**

Yasuo Nakada,<sup>1,2</sup> Hideo O. Bando,<sup>3</sup> Daisuke Aki,<sup>4</sup> Romya Barik,<sup>5</sup> Jesse A. Shick,<sup>6</sup> Edoardo O. Mazzoni,<sup>7</sup> Dazhi Ratsberg<sup>8</sup>



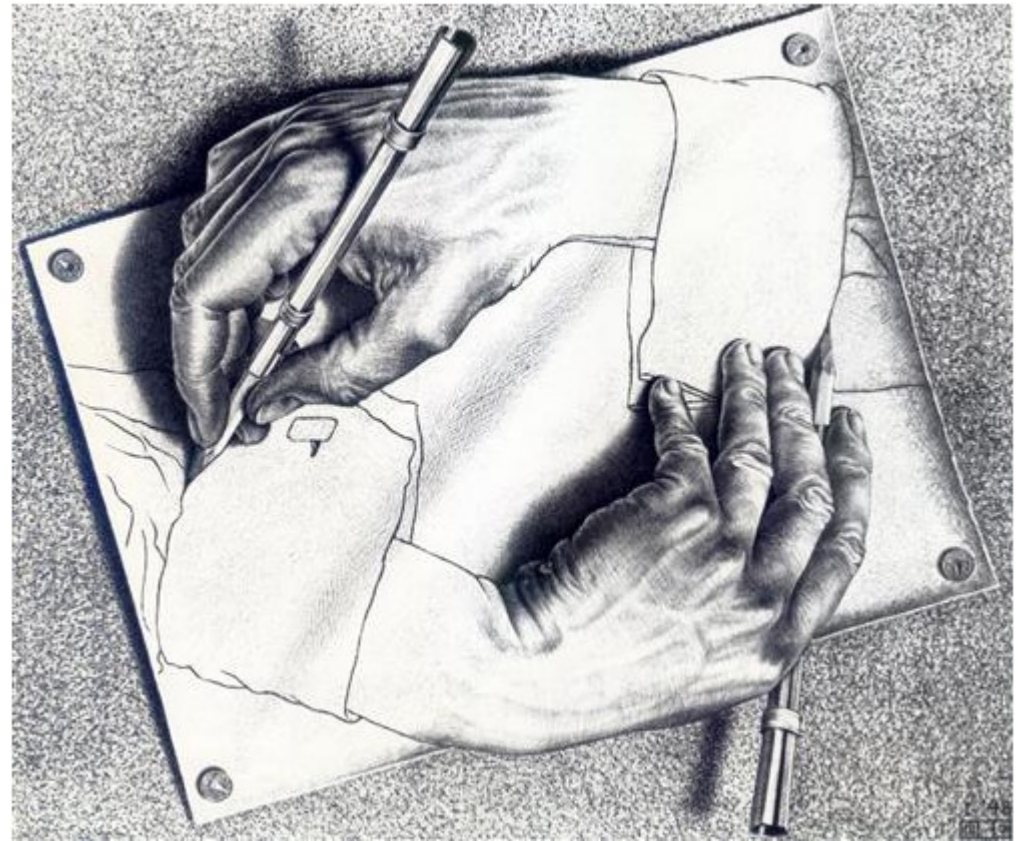
3D structure impacts gene expression



# Outline

## Biological context

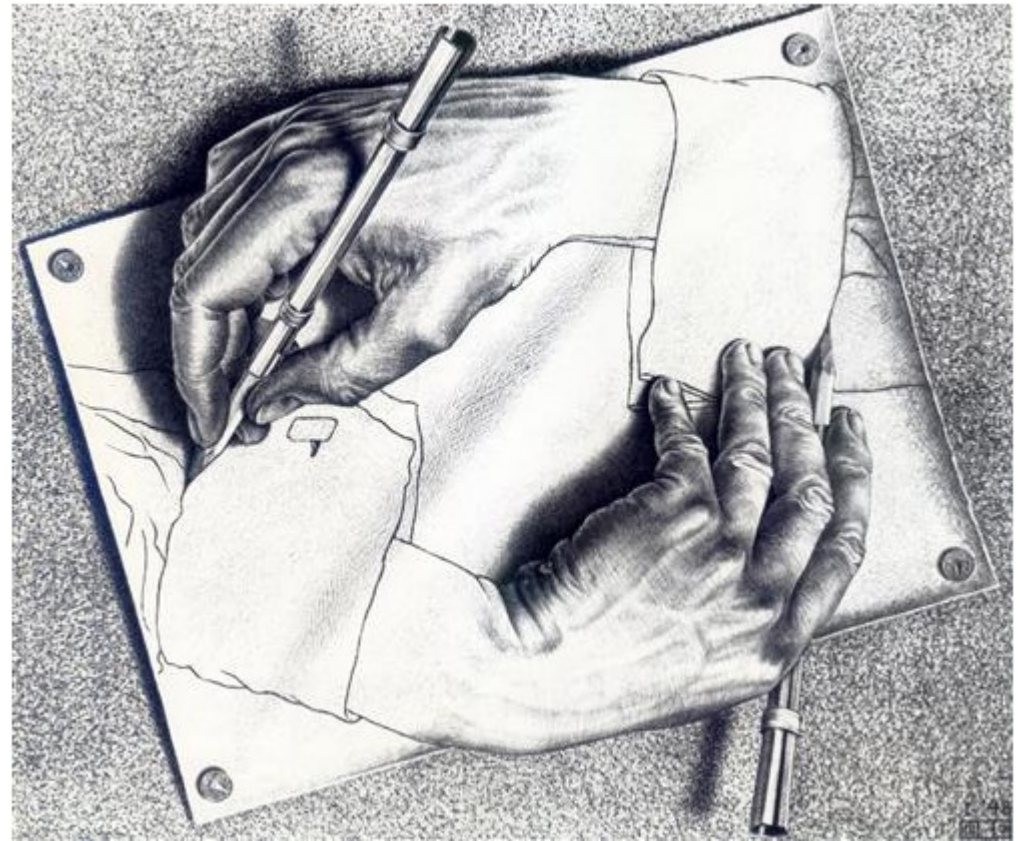
- ◆ More biological context
- ◆ Hi-C data processing
  - ◆ map
  - ◆ filter
  - ◆ count
  - ◆ normalize
  - ◆ segment
  - ◆ compare
- ◆ Conclusion, discussion, NETBIO lunch



*M.C. Escher, 1948*

# Outline

- ◆ Biological context
- ➔ More biological context
- ◆ Hi-C data processing
  - ◆ map
  - ◆ filter
  - ◆ count
  - ◆ normalize
  - ◆ segment
  - ◆ compare
- ◆ Conclusion, discussion, NETBIO lunch



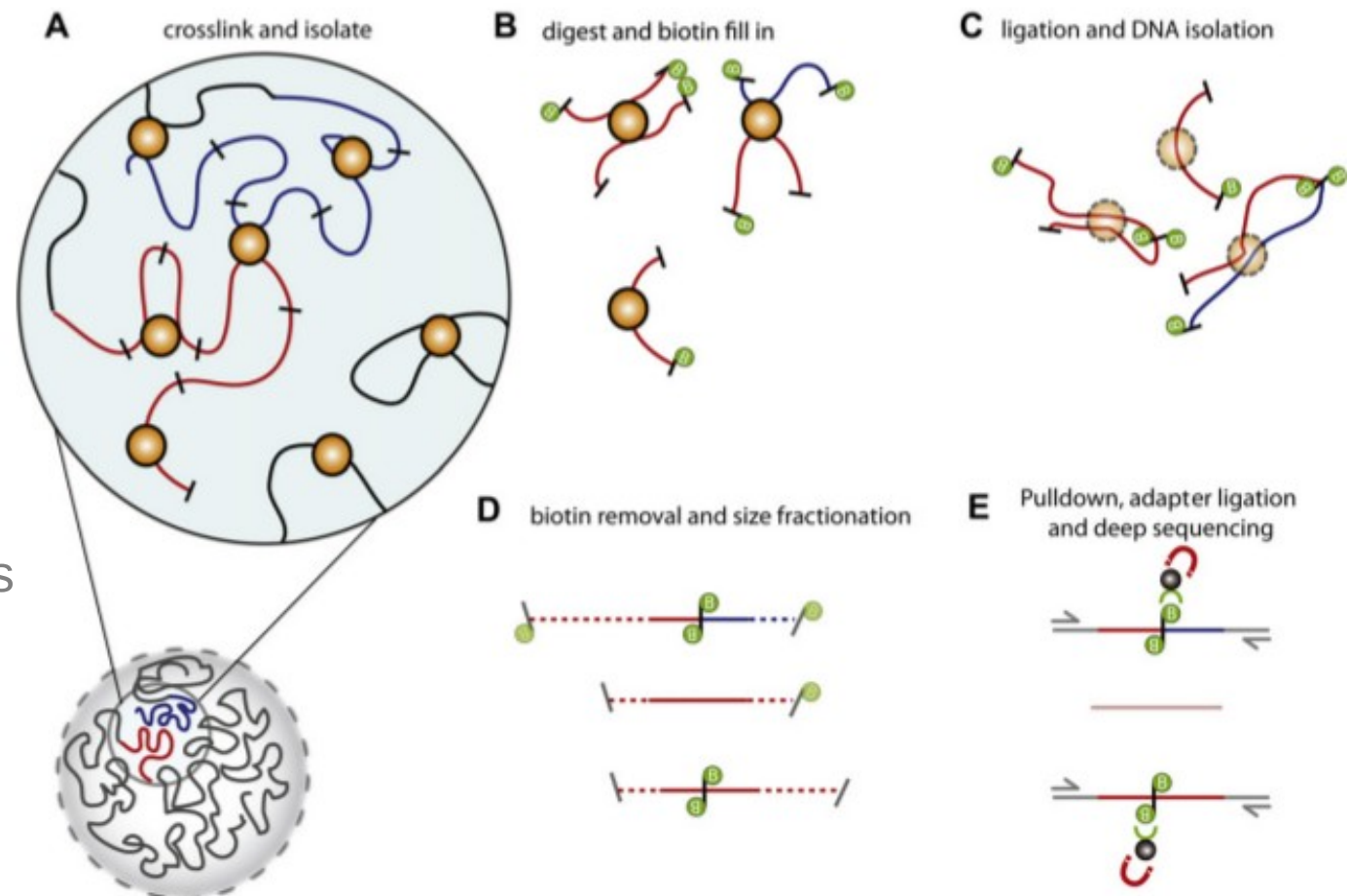
*M.C. Escher, 1948*

# Hi-C: the experiment

Hi-C: high-throughput chromatin conformation capture  
(Lieberman-Aiden et al, Science, 2009, Rao et al, Cell, 2014)

*J.-M. Belton et al./Methods 58 (2012) 268–276*

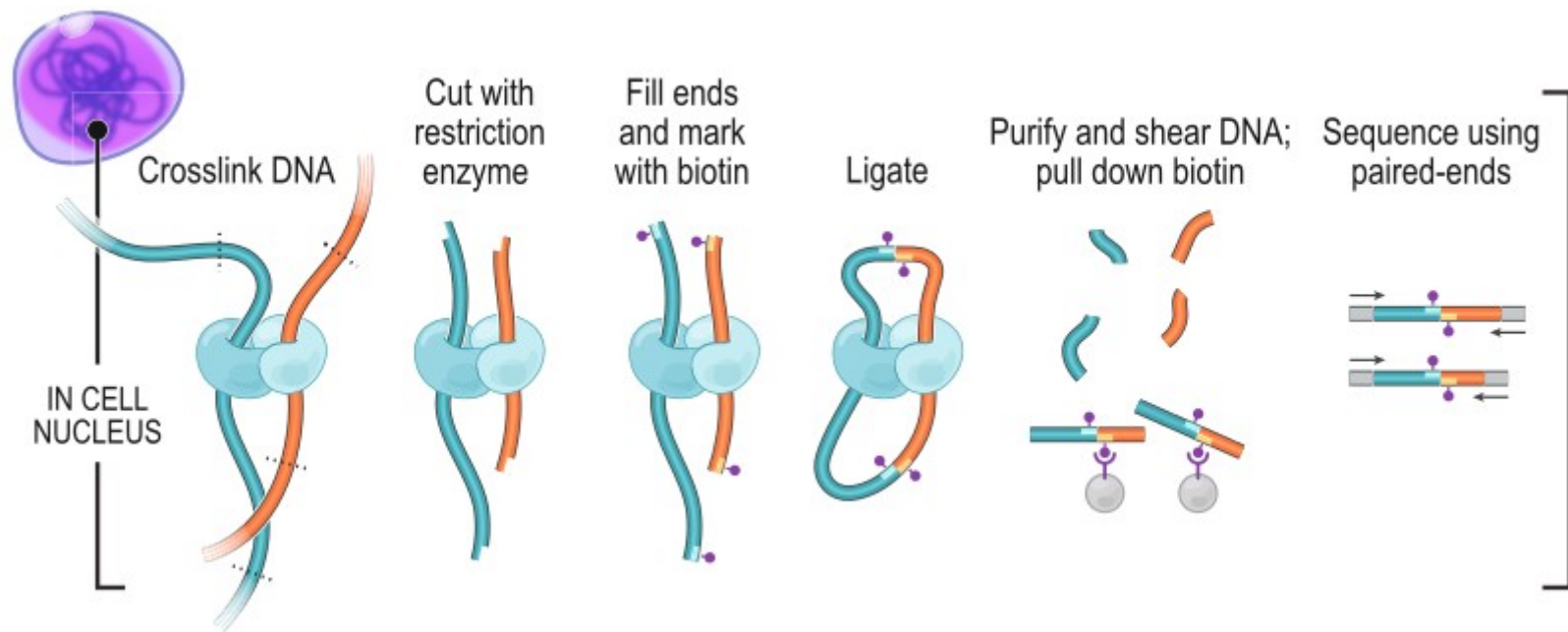
- ◆ crosslink DNA (“fixation”)
- ◆ cleave genome with restriction enzyme
- ◆ biotin-mark and ligate extremities
- ◆ fragment, select biotin-marked junctions
- ◆ sequence fragments (paired-ends)





# Hi-C: the experiment

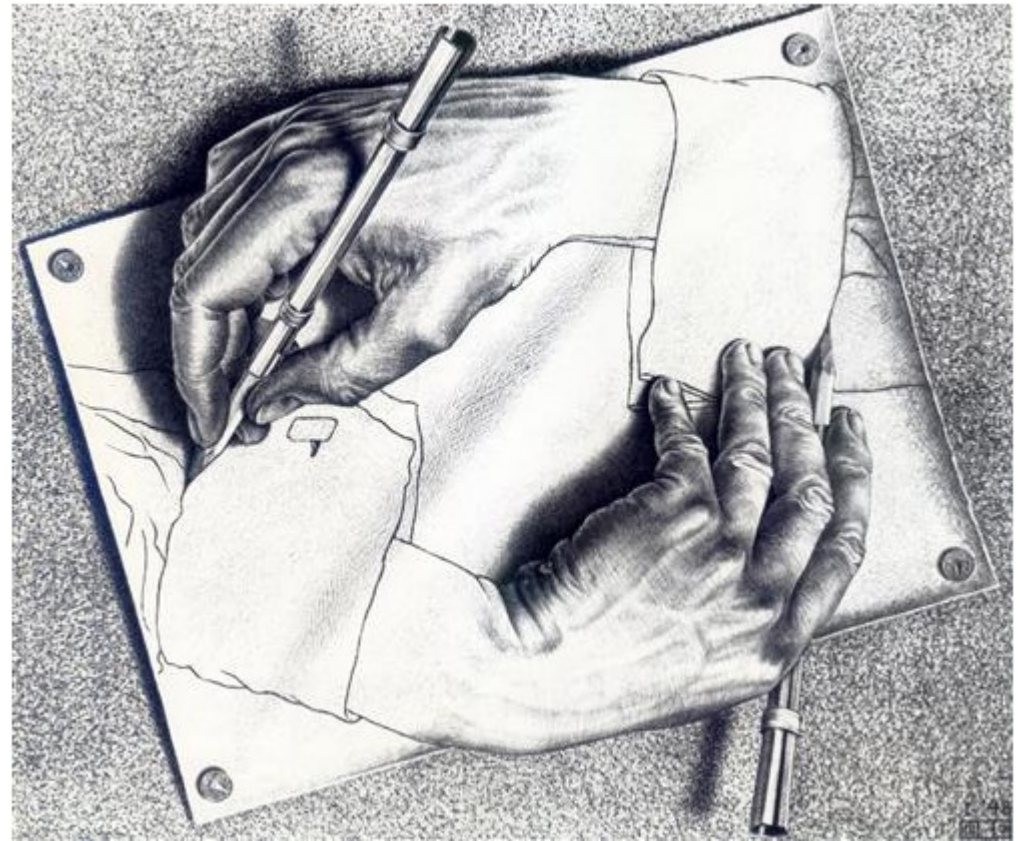
Hi-C: high-throughput chromatin conformation capture  
(Lieberman-Aiden et al, Science, 2009, Rao et al, Cell, 2014)



*Rao et al, Cell, 2014*

# Outline

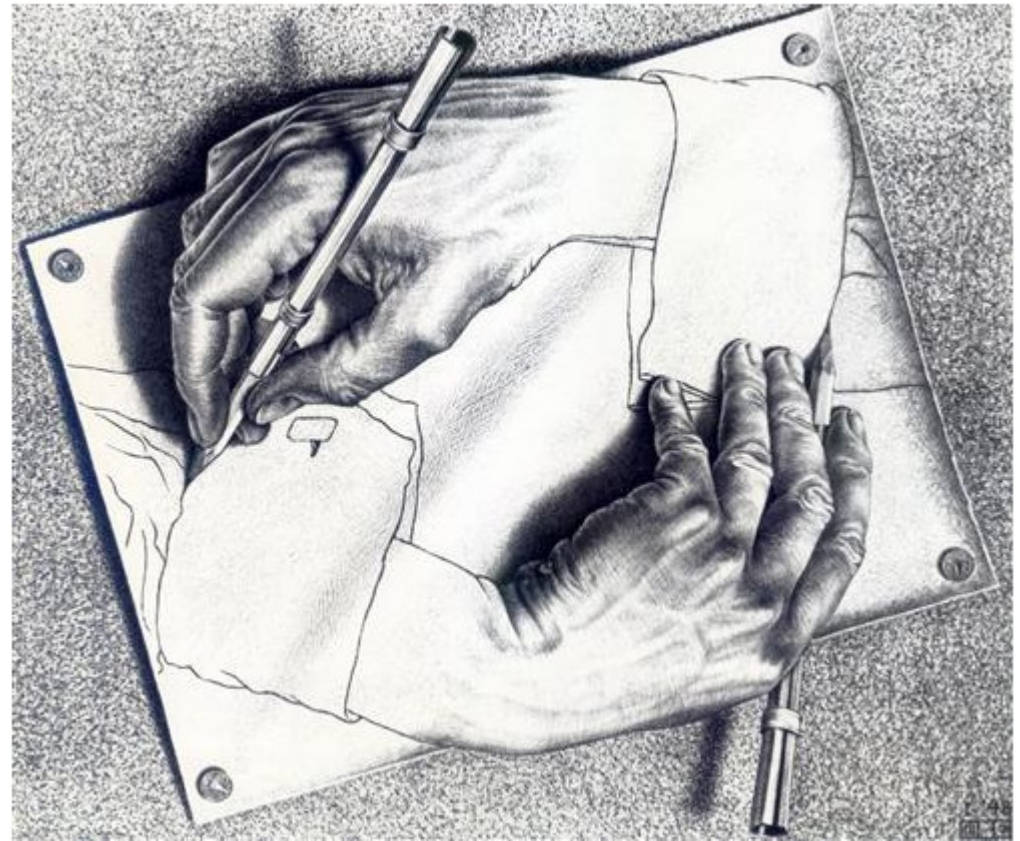
- ◆ Biological context
- ➔ More biological context
- ◆ Hi-C data processing
  - ◆ map
  - ◆ filter
  - ◆ count
  - ◆ normalize
  - ◆ segment
  - ◆ compare
- ◆ Conclusion, discussion, NETBIO lunch



*M.C. Escher, 1948*

# Outline

- ◆ Biological context
- ◆ More biological context
- ➔ Hi-C data processing
  - ◆ map
  - ◆ filter
  - ◆ count
  - ◆ normalize
  - ◆ segment
  - ◆ compare
- ◆ Conclusion, discussion, NETBIO lunch



*M.C. Escher, 1948*

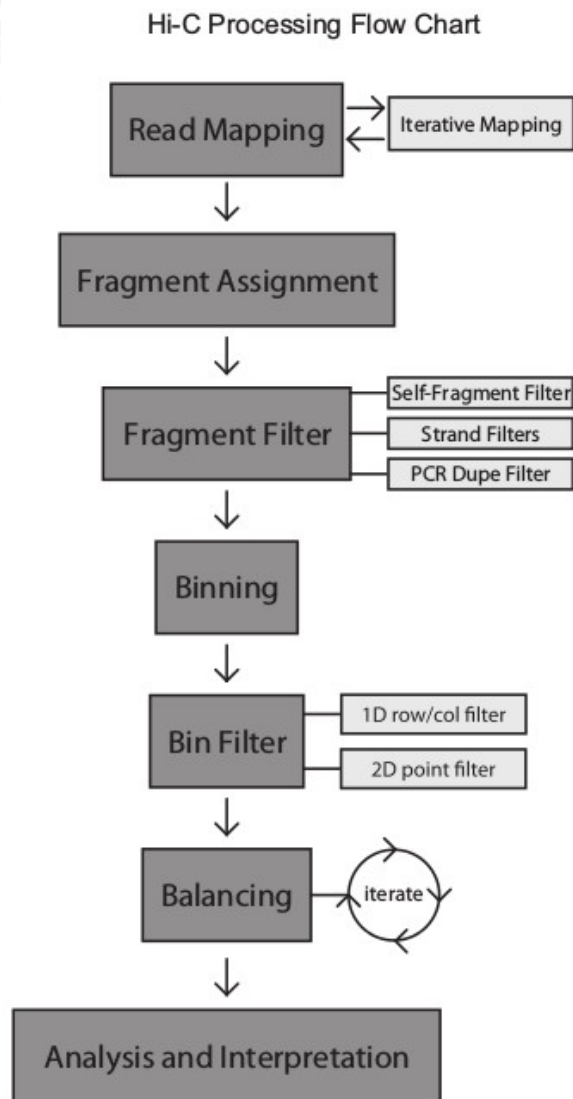




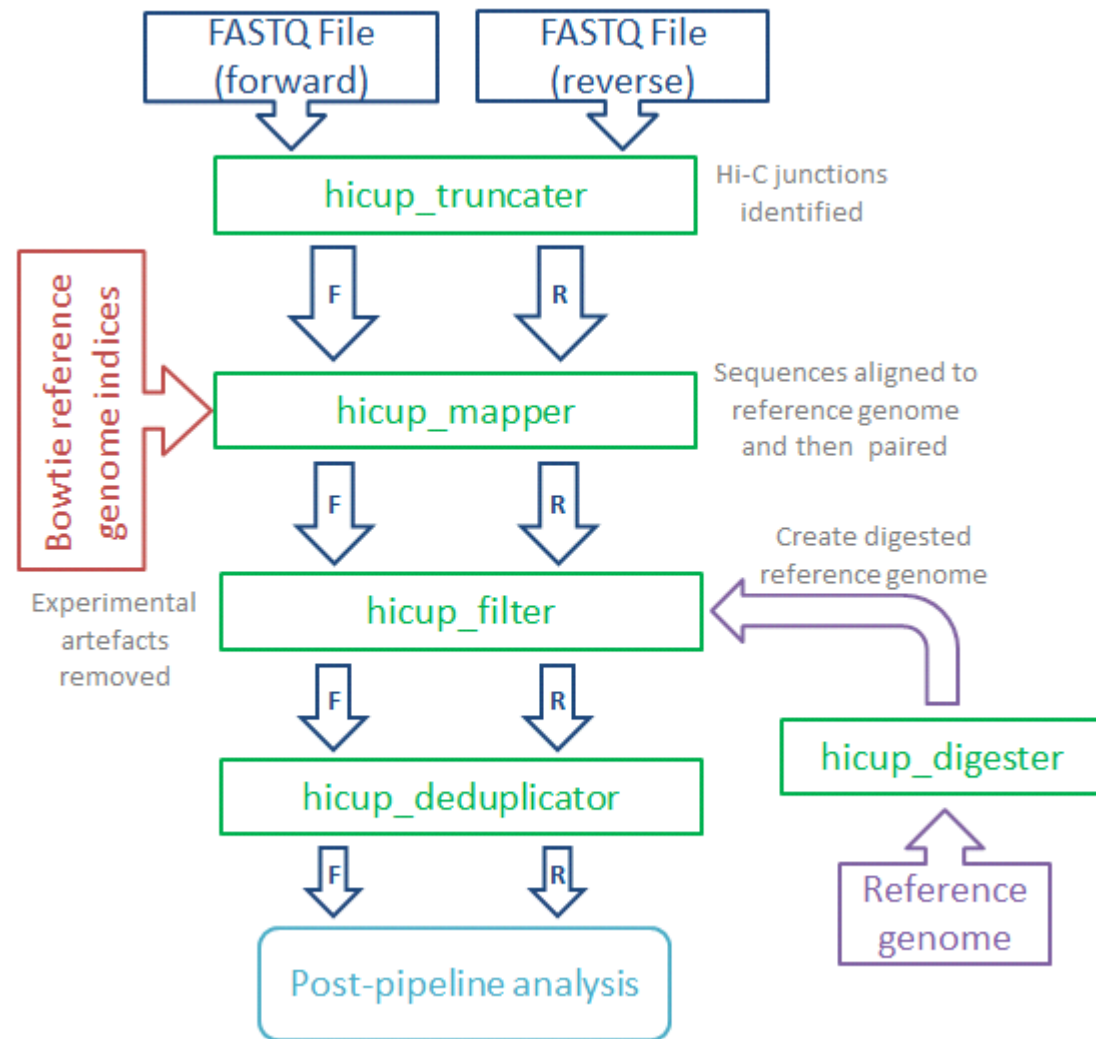
# Hi-C data analysis: overview

- ◆ clean and trim the reads
- ◆ map the reads on the genomic reference
- ◆ filter bogus configurations
- ◆ count the reads per genomic bin => contact matrix
- ◆ normalize the matrix
- ◆ identify topological domains, cis- and trans- interactions
- ◆ comparative/integrative analysis

# Hi-C data analysis: overview

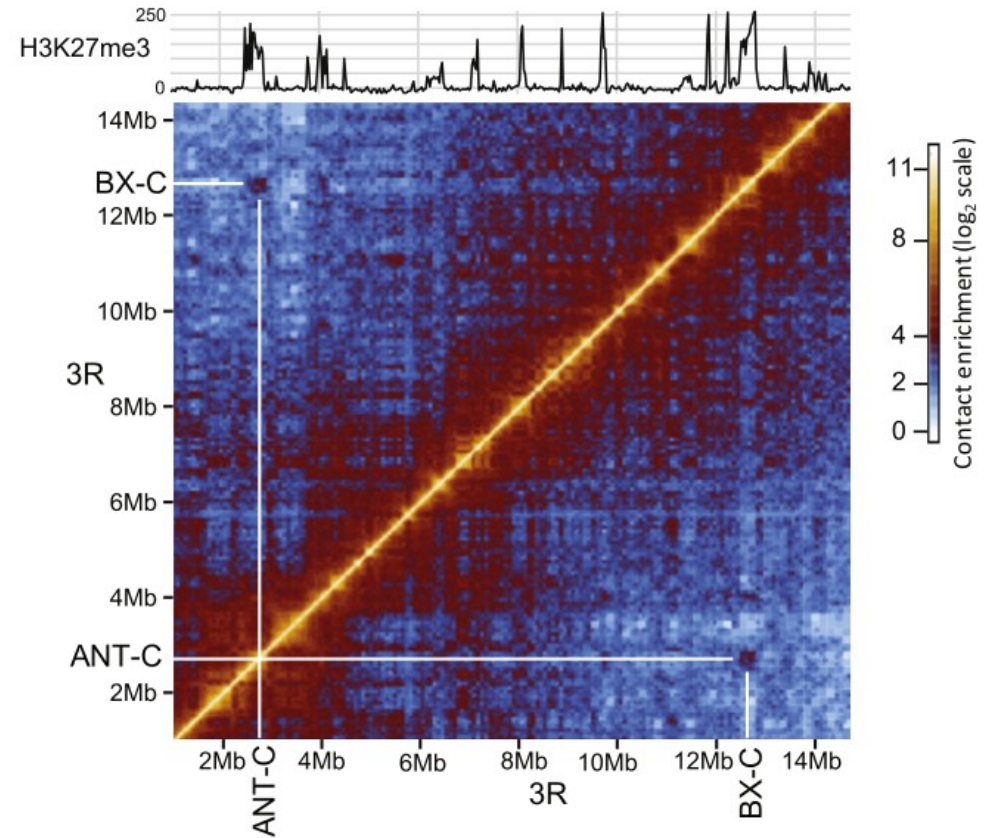
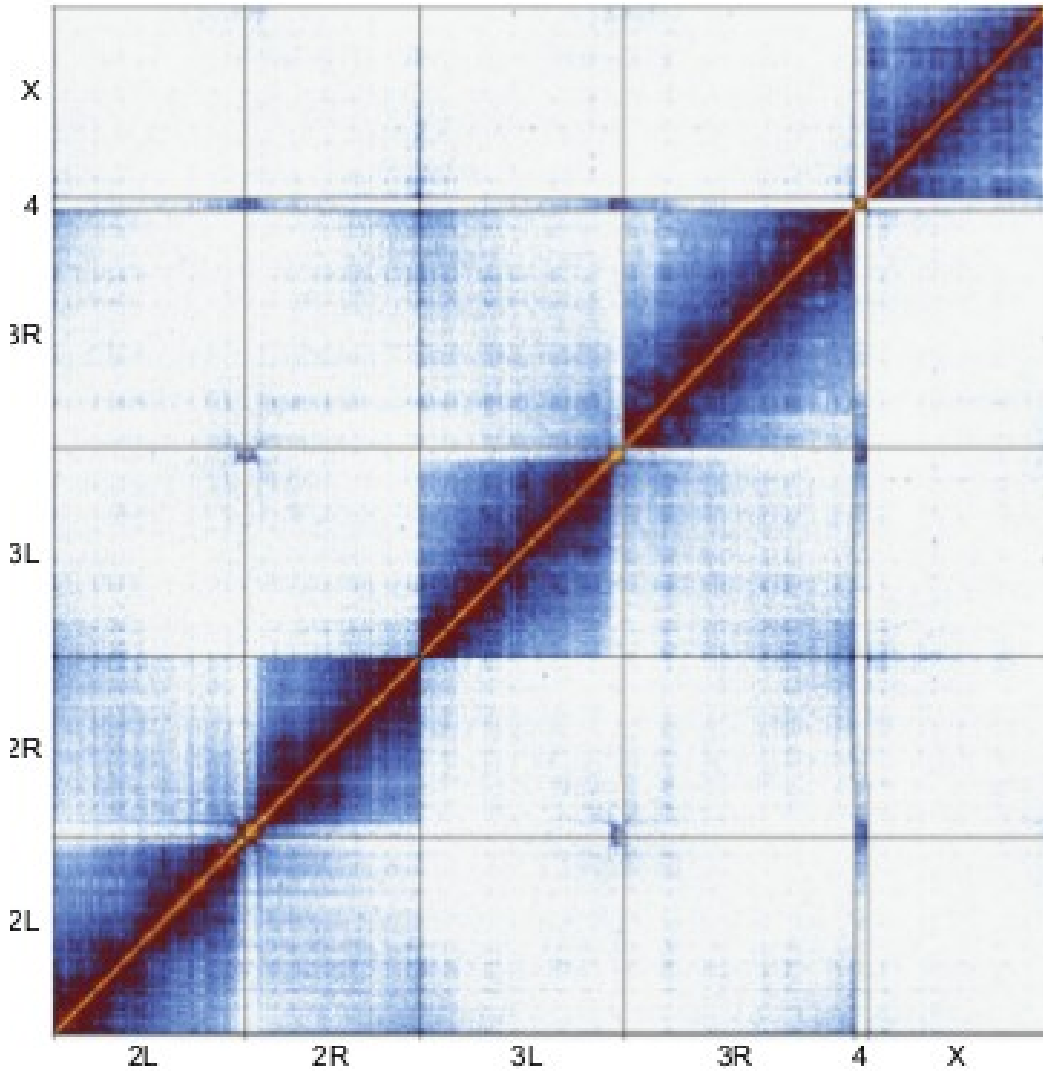


Lajoie et al, 2015



HiCUP, [www.bioinformatics.babraham.ac.uk/](http://www.bioinformatics.babraham.ac.uk/)

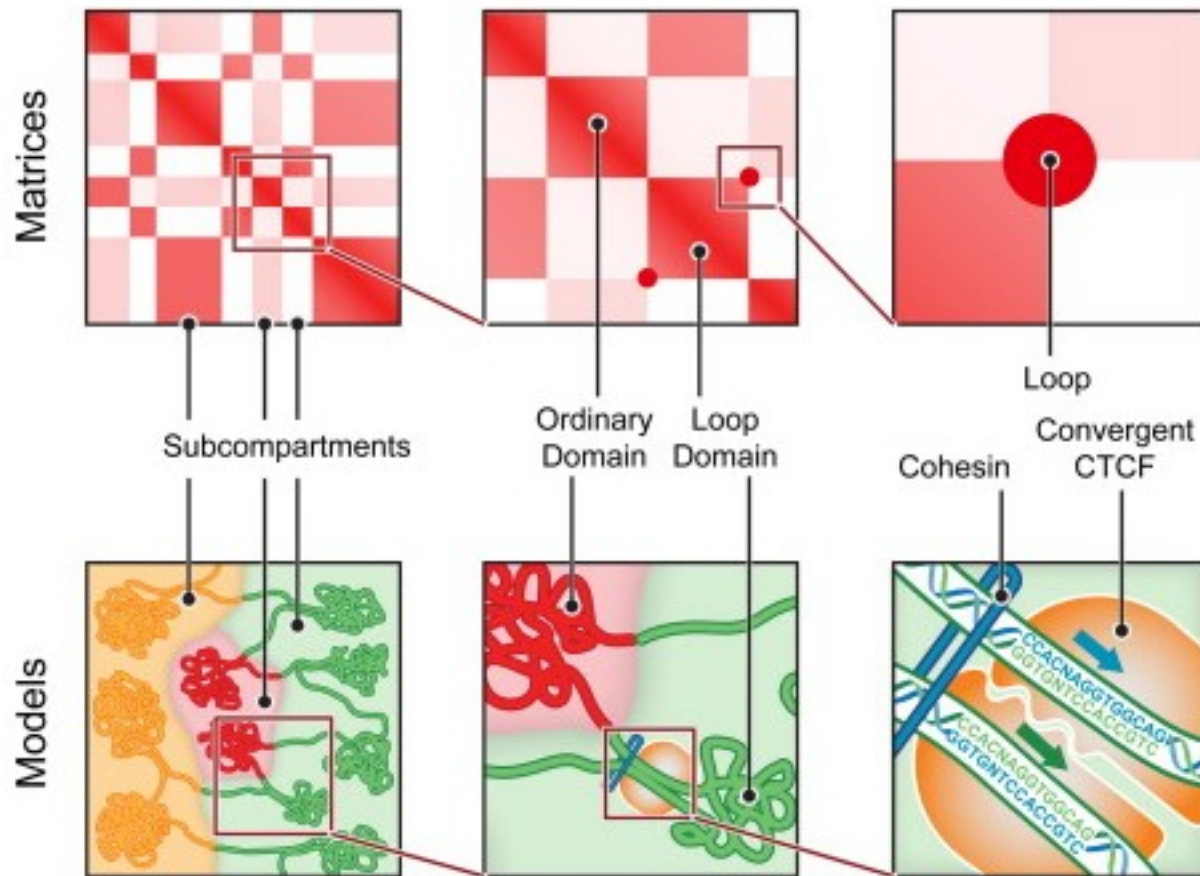
# Hi-C data analysis: the contact matrix



*Sexton et al 2012*



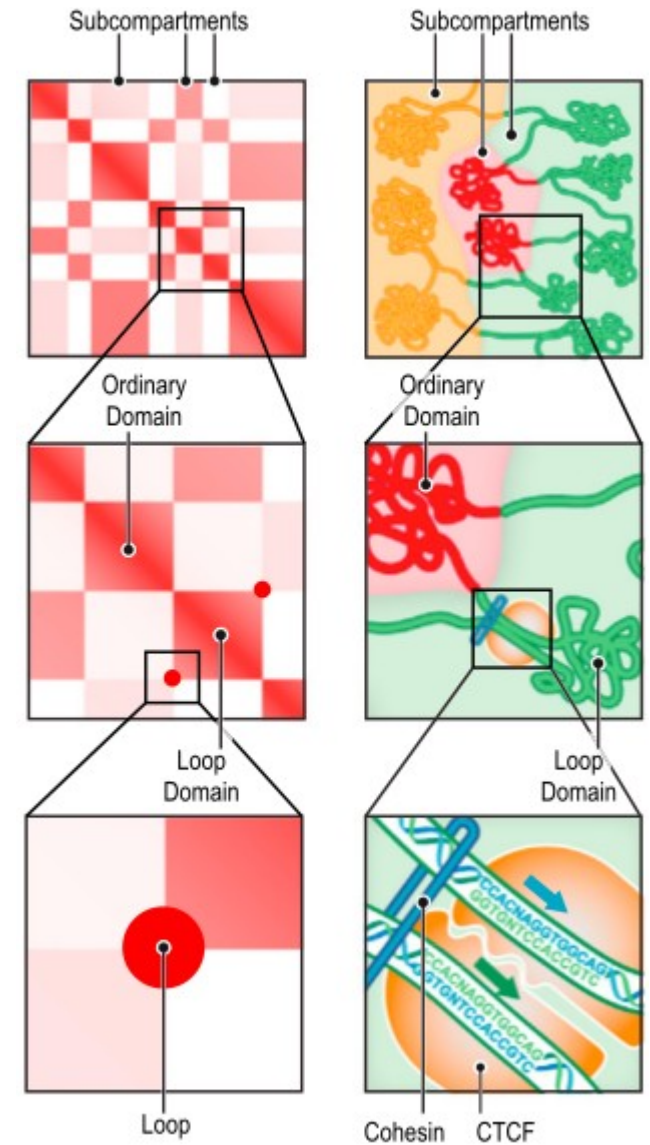
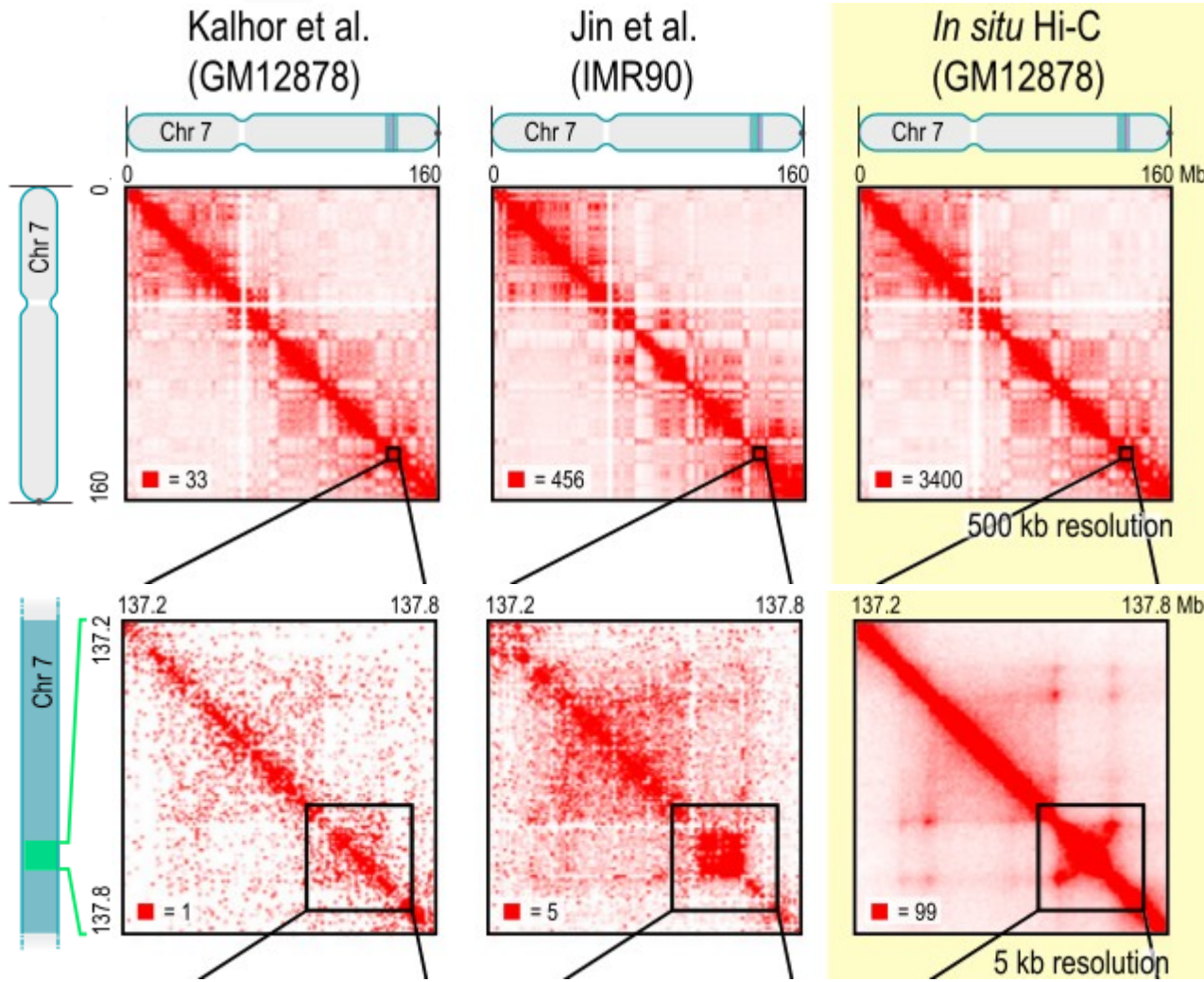
# Hi-C data analysis: the contact matrix



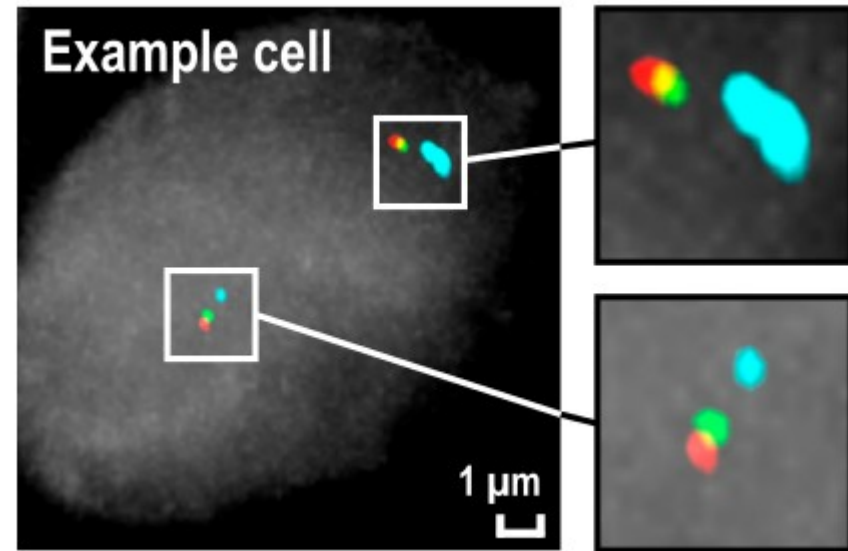
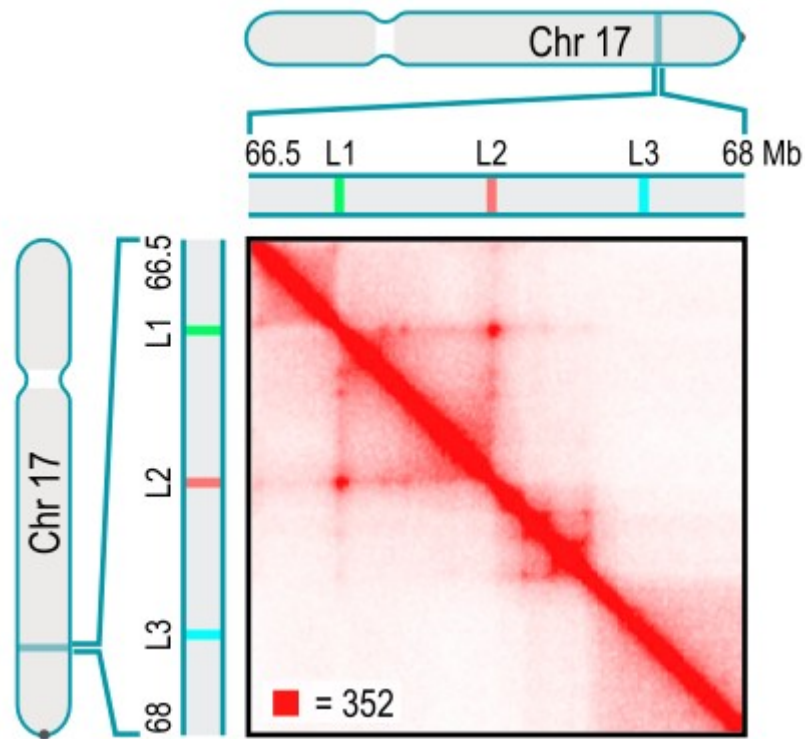
*Rao et al, Cell, 2014*

# Hi-C data analysis: the contact matrix

Rao et al, Cell, 2014



# Hi-C data analysis: the contact matrix



*Rao et al, Cell, 2014*





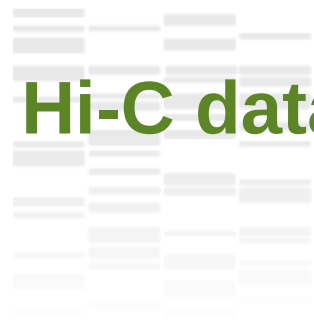
# Hi-C data analysis: matrix normalization

Number of reads per bin (coverage) depends on:

- ◆ GC%
- ◆ density of restriction sites
- ◆ repeats and “mappability”
- ◆ overall depth of coverage
- ◆ Others?

=> “Parametric” vs. “non-parametric” normalization

# Hi-C data analysis: matrix normalization



## Probabilistic modeling of Hi-C contact maps eliminates systematic biases to characterize global chromosomal architecture

Eitan Yaffe & Amos Tanay

Hi-C experiments measure the probability of physical proximity between pairs of chromosomal loci on a genomic scale. We report on several systematic biases that substantially

To fulfill this promise, 3C techniques and their derivations must become robust and quantitative. The complicated experimental procedure that includes fixation, digestion, ligation and amplification



NIH Public Access  
Author Manuscript

*Nat Methods*. Author manuscript; available

Published in final edited form as:

*Nat Methods*. 2012 October ; 9(10): . doi:10.10

NIH-PA Author Manuscript

### Iterative Correction of Hi-C Data Reveals Hallmarks of Chromosome Organization

Maxim Imakaev<sup>1,\*</sup>, Geoffrey Fudenberg<sup>2,\*</sup>, Rachel Patton McC  
Anton Goloborodko<sup>1</sup>, Bryan R. Lajoie<sup>3</sup>, Job Dekker<sup>3,#</sup>, and Le

<sup>1</sup>Department of Physics, MIT, Cambridge, MA

<sup>2</sup>Graduate Program in Bioscience, Harvard University, Cambridge

**BIOINFORMATICS APPLICATIONS NOTE** Vol. 28 no. 23 2012, pages 3131–3133  
doi:10.1093/bioinformatics/bts570

Genome analysis

Advance Access publication September 27, 2012

### HiCNorm: removing biases in Hi-C data via Poisson regression

Ming Hu<sup>1</sup>, Ke Deng<sup>1</sup>, Siddarth Selvaraj<sup>2,3</sup>, Zhaohui Qin<sup>4</sup>, Bing Ren<sup>2</sup> and Jun S. Liu<sup>1,\*</sup>

<sup>1</sup>Department of Statistics, Harvard University, Cambridge, MA 02138, USA, <sup>2</sup>Department of Cellular and Molecular Medicine, UCSD School of Medicine, La Jolla, CA 92093, USA, <sup>3</sup>Bioinformatics and Systems Biology Graduate Program, University of California, San Diego, La Jolla, CA 92093, USA and <sup>4</sup>Department of Biostatistics and Bioinformatics, Emory University, Atlanta, GA 30322 USA

Associate Editor: Alex Bateman

#### ABSTRACT

**Summary:** We propose a parametric model, HiCNorm, to remove systematic biases in the raw Hi-C contact maps, resulting in a matrix where the row and column sums are equal to one. However, the raw Hi-C contact matrix is of the massive size  $O(N^2)$ , where  $N$  is the number of genomic regions. Thus, it requires expensive computing resources such as large memory and long computation

biases, a non-parametric probabilistic model (referred to hereafter as the YT approach) was proposed that explicitly models the probability of observing a paired-end read spanning two fragment ends. This approach can remove the majority of sys-

Category

### HiCorrector: A fast, scalable and memory-efficient package normalizing large-scale Hi-C data

Wenyuan Li<sup>1</sup>, Ke Gong<sup>1</sup>, Qingjiao Li<sup>1</sup>, Frank Alber<sup>1\*</sup> and Xianghong Jasmine Z

<sup>1</sup> Molecular and Computational Biology, University of Southern California, Los Angeles, CA 90089, USA  
Associate Editor: Prof. Alfonso Valencia

#### ABSTRACT

**Summary:** Genome-wide proximity ligation assays, e.g. Hi-C and its variant TCC, have recently become important tools to study spatial genome organization. Removing biases from chromatin contact

and column sums of the matrix are equal to one. However, the raw Hi-C contact matrix is of the massive size  $O(N^2)$ , where  $N$  is the number of genomic regions. Thus, it requires expensive computing resources such as large memory and long computation

# Hi-C data analysis: matrix normalization

## A FAST ALGORITHM FOR MATRIX BALANCING

PHILIP A. KNIGHT\* AND DANIEL RUIZ<sup>†</sup>

**Abstract.** As long as a square nonnegative matrix  $A$  contains sufficient nonzero elements, then the matrix can be balanced, that is we can find a diagonal scaling of  $A$  that is doubly stochastic. A number of algorithms have been proposed to achieve the balancing, the most well known of these being Sinkhorn-Knopp. In this paper we derive new algorithms based on inner-outer iteration schemes. We show that Sinkhorn-Knopp belongs to this family, but other members can converge much more quickly. In particular, we show that while stationary iterative methods offer little or no improvement in many cases, a scheme using a preconditioned conjugate gradient method as the inner iteration can give quadratic convergence at low cost.

**Key words.** Matrix balancing, Sinkhorn-Knopp algorithm, doubly stochastic matrix, conjugate gradient iteration.

**AMS subject classifications.** 15A48, 15A51, 65F10, 65H10.

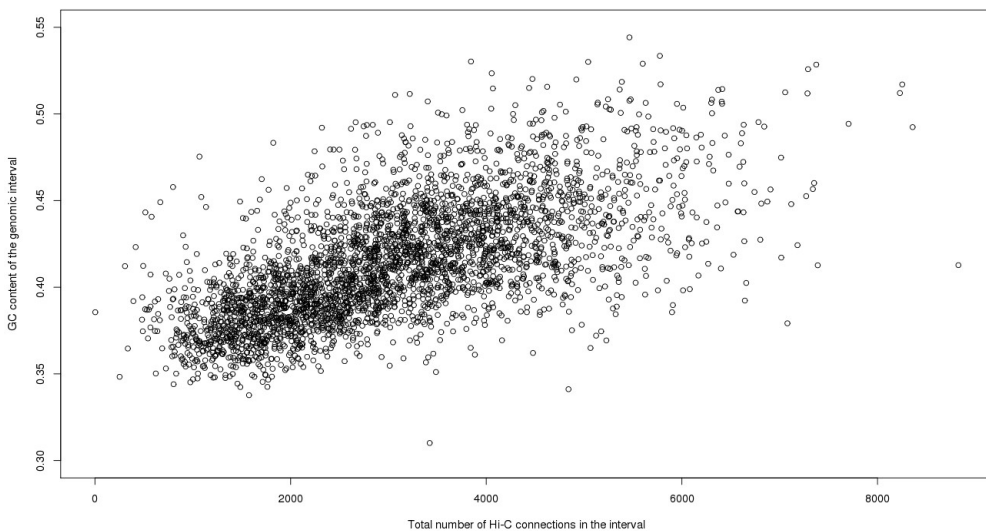
**1. Introduction.** For at least 70 years, scientists in a wide variety of disciplines have attempted to transform square nonnegative matrices into doubly stochastic form by applying diagonal scalings. That is, given  $A \in \mathbb{R}^{n \times n}$ ,  $A \geq 0$ , find diagonal matrices  $D_1$  and  $D_2$  so that  $P = D_1 A D_2$  is doubly stochastic. Motivations for achieving this balance include interpreting economic data [1], preconditioning sparse matrices [16], understanding traffic circulation [14], assigning seats fairly after elections [3], matching protein samples [4] and ordering nodes in a graph [12]. In all of these applications, one of the main methods considered is SK<sup>1</sup>. This is an iterative process that attempts to find  $D_1$  and  $D_2$  by alternately normalizing columns

*Knight & Ruiz, IMA J. Numer. Anal., 2013*

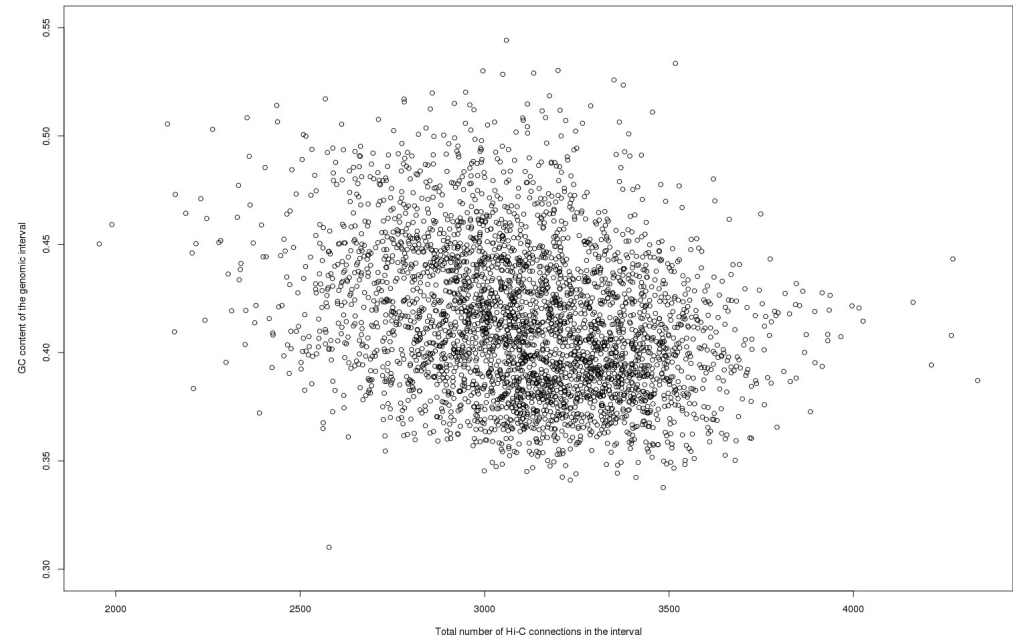


# Hi-C data analysis: matrix normalization

example with number of reads vs. GC%



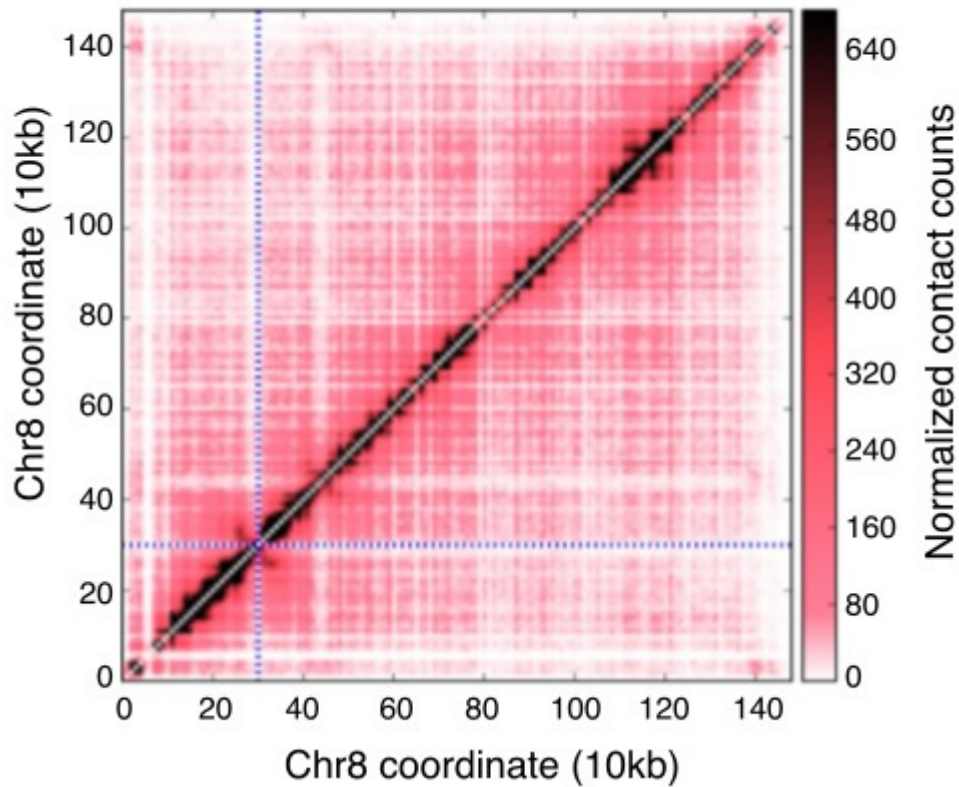
before normalization



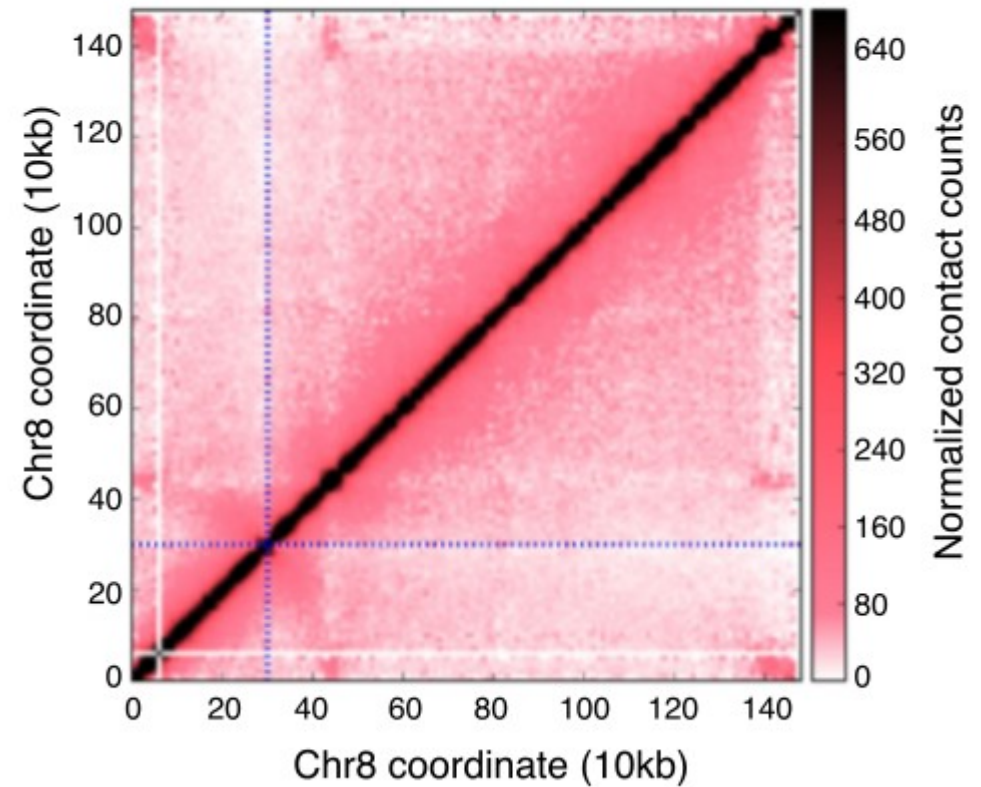
after normalization

# Hi-C data analysis: matrix normalization

**(a)** Raw contact map

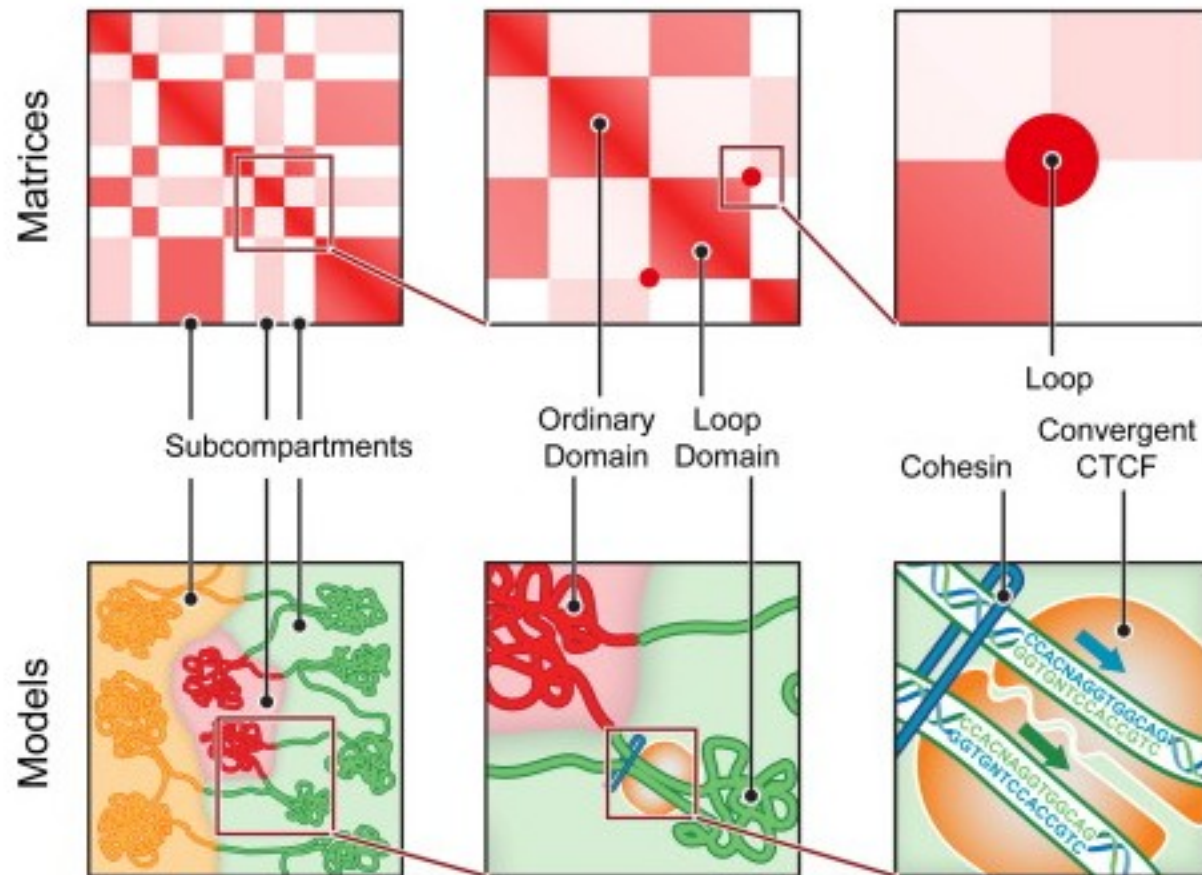


**(b)** Normalized contact map



*Ay & Noble, Genome Biology, 2015*

# Hi-C data analysis: finding topologically associated domains (TADs)

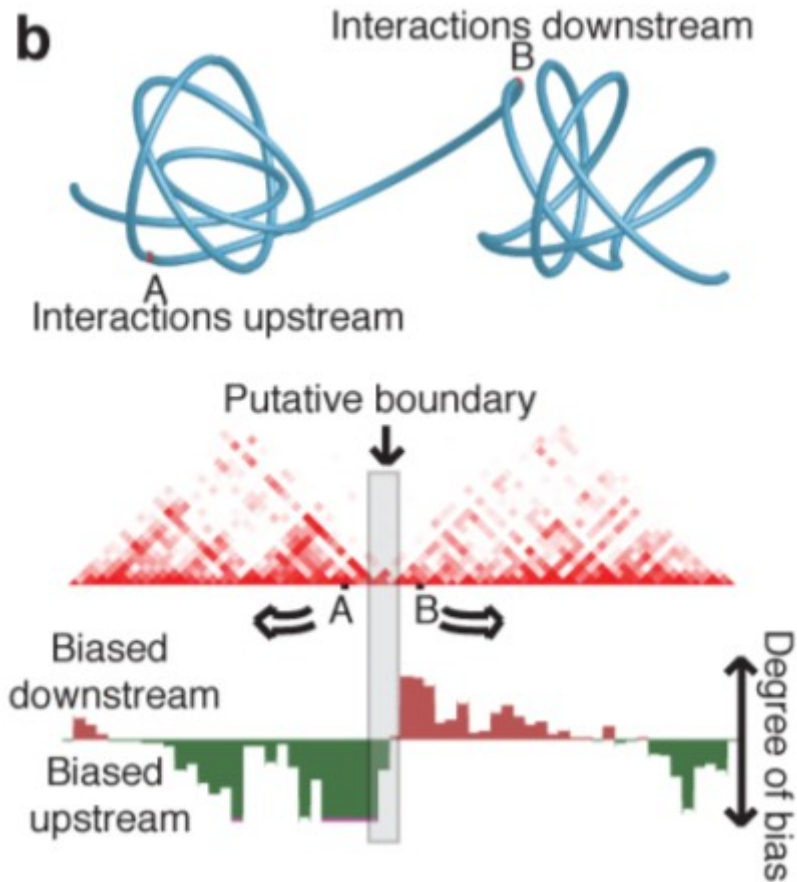


*Rao et al, Cell, 2014*

- ◆ methods: clustering, 2D-segmentation, etc



# Hi-C data analysis: finding topologically associated domains (TADs)



Dixon et al., Nature, 2012

NIH Public Access  
Author Manuscript  
Nature. Author manuscript; available in PMC 2012 November 17.

Published in final edited form as:  
*Nature.* ; 485(7398): 376–380. doi:10.1038/nature11082.

**Topological Domains in Mammalian Genomes Identified by Analysis of Chromatin Interactions**

**Jesse R. Dixon<sup>1,3,4</sup>, Siddarth Selvaraj<sup>1,5</sup>, Feng Yue<sup>1</sup>, Audrey Kim<sup>1</sup>, Yan Li<sup>1</sup>, Yin Shen<sup>1</sup>, Ming Hu<sup>6</sup>, Jun S. Liu<sup>6</sup>, and Bing Ren<sup>1,2,\*</sup>**

<sup>1</sup>Ludwig Institute for Cancer Research  
<sup>2</sup>University of California, San Diego School of Medicine, Department of Cellular and Molecular Medicine, Institute of Genomic Medicine, 9500 Gilman Drive, La Jolla, CA 92093  
<sup>3</sup>Medical Scientist Training Program, University of California, San Diego, La Jolla CA 92093  
<sup>4</sup>Biomedical Sciences Graduate Program, University of California, San Diego, La Jolla CA 92093  
<sup>5</sup>Bioinformatics and Systems Biology Graduate Program, University of California, San Diego, La Jolla CA 92093  
<sup>6</sup>Department of Statistics, Harvard University, 1 Oxford Street, Cambridge, MA 02138

**Abstract**

The spatial organization of the genome is intimately linked to its biological function, yet our understanding of higher order genomic structure is coarse, fragmented and incomplete. In the nucleus of eukaryotic cells, interphase chromosomes occupy distinct chromosome territories (CT).

NIH-PA Author Manuscript  
 NIH-PA Author Manuscript

Directionality Index

$$DI = \left( \frac{B - A}{|B - A|} \right) \left( \frac{(A - E)^2}{E} + \frac{(B - E)^2}{E} \right)$$

DI HMM => TADs

# Hi-C data analysis: finding topologically associated domains (TADs)

## Identification of hierarchical chromatin domains

Caleb Weinreb<sup>1</sup>, and Benjamin J. Raphael<sup>1,2\*</sup>

<sup>1</sup>Center for Computational Molecular Biology, Brown University, Providence, RI

<sup>2</sup>Department of Computer Science, Brown University, Providence, RI

Associate Editor: Prof. Gunnar Ratsch

### ABSTRACT

**Motivation:** The 3D structure of the genome is an important regulator of many cellular processes including differentiation and gene regulation. Recently, technologies such as Hi-C that combine proximity ligation with high-throughput sequencing have revealed domains of self-interacting chromatin, called topologically associating domains (TADs), in many organisms. Current methods for identifying TADs using Hi-C data assume that TADs are non-overlapping, despite evidence for a nested structure in which TADs and sub-TADs form a complex hierarchy.

**Results:** We introduce a model for hierarchical decomposition of contact frequencies into a hierarchy of nested TADs. This model is based empirical distributions of contact frequencies within TADs, where positions that are far apart have a greater enrichment of contacts than

resulting in a contact matrix  $A$ , where  $A_{ij}$  is the number of contacts between bins  $i$  and  $j$ , normalized for experimental bias. Several methods have been developed for the identification of TADs from Hi-C data. These methods may be roughly classified into two categories: (1) methods that define a 1D test statistic from the contact matrix  $A_{ij}$ ; (2) methods that exploit the 2D structure of the contact matrix.

Dixon *et al.* (2012) compute a 1D “directionality index” (DI) from the contact matrix. This index defines whether contacts have an upstream bias, downstream bias or no bias. Next, they use a hidden Markov model (HMM) to partition the genome into regions defined by changes in the directionality index. Each transition into downstream bias marks the start of a domain and the next transition out of upstream bias marks its end. Sauria *et al.* (2014) introduce a 1D

defined below.

**DEFINITION 2.** Consider a TAD  $D$  and interval  $[i, j] \subseteq [D_L, D_R]$ . Let the error compensation  $\mathcal{E}_C(i, j, D)$  be

$$\mathcal{E}_C(i, j, D) = \sum_{l=i}^j \sum_{k=l}^j (\bar{A}_D(l, k) - A_{lk})^2. \quad (8)$$

Using the error compensation, we derive an expression for the score of a TAD tree in terms of its root TAD and sub-trees.

**PROPOSITION 1.** Let  $T$  be a TAD tree consisting of a root TAD  $D$  and a collection of non-overlapping sub-trees  $T_1, \dots, T_m$ , spanning the intervals  $[i_1, j_1], \dots, [i_m, j_m]$ . The score  $\mathcal{O}_\gamma(T)$  can be decomposed as

$$\mathcal{O}_\gamma(T) = \mathcal{O}_\gamma(D) + \sum_{x=1}^m (\mathcal{O}_\gamma(T_x) + \mathcal{E}_C(i_x, j_x, D)). \quad (9)$$

We now describe steps (1-3) above in greater detail. To perform step (1), recall that a TAD is defined by four parameters,  $(L_D, R_D, \delta_D, \beta_D)$ . Thus, in choosing the root TAD  $D$ , two parameters are given ahead of time ( $[L_D, R_D] = [i, j]$ ), meaning we only need to select optimal values for  $\delta_D$  and  $\beta_D$ . Next, for a given choice of  $\delta_D$  and  $\beta_D$ , we must choose a set of non-overlapping sub-trees, defined by sub-intervals  $[i_x, j_x]$  and multiplicities  $n_x$  (steps 2-3). To that end, let  $\mathcal{I}(i, j, N)$  be the collection of sets  $\{(i_x, j_x, n_x)\}$  that satisfy the following properties: (i)  $[i_x, j_x]$  are non-overlapping sub-intervals of  $[i, j]$ ; (ii)  $\sum n_x = N - 1$ ; (iii)  $i_x$  and  $j_x$  are valid boundaries. Using  $\mathcal{I}(i, j, N)$  as a search space, we evaluate  $\Phi(i, j, N, \delta)$  as follows.

**PROPOSITION 2.** For each interval  $[i, j]$  and positive integer  $N$ ,

$$\Phi(i, j, N, \delta) = \max_{\{(i_x, j_x, n_x) \in \mathcal{I}(i, j, N)\}} \left( \mathcal{O}_\gamma(D) + \sum_{x=1}^m (\mathcal{O}_\gamma(T_x) + \mathcal{E}_C(i_x, j_x, D)) \right)$$

**PROBLEM 2.** Given  $N \in \mathbb{N}$  and  $\gamma \in \mathbb{R}^+$ , find the TAD forest  $F$  that maximizes the objective  $\mathcal{O}_\gamma(F) = \gamma \mathcal{B}_{p,q}(F) - \mathcal{E}(F)$  such that  $|F| = N$ , and each  $D \in F$  is locally fitted and has valid boundaries.

Once again, our first step in solving Problem 2 will be to find optimal TAD trees over every interval.

**DEFINITION 4.** Given  $N \in \mathbb{N}$ , and the interval  $[i, j]$ , define  $\hat{\Phi}(i, j, N) := \max \mathcal{O}_\gamma(T)$  over all TAD trees  $T$  such that (i)  $T$  is rooted at the interval  $[i, j]$ , (ii)  $T$  contains  $N$  TADs ( $|T| = N$ ), and (iii) each  $D \in T$  is locally fitted has valid boundaries.

In contrast to  $\Phi(i, j, N, \delta)$ ,  $\hat{\Phi}(i, j, N)$  does not take  $\delta$  as an argument, since it maximizes over TAD trees whose  $\delta$  values are fixed by the requirement that they be locally fitted. This leads to the following proposition, which shows how to evaluate  $\hat{\Phi}(i, j, N)$ .

**PROPOSITION 3.** For each interval  $[i, j]$  and positive integer  $N$ ,

$$\hat{\Phi}(i, j, N) = \mathcal{O}_\gamma(\hat{D}_{ij}) + \max_{\{(i_x, j_x, n_x) \in \mathcal{I}(i, j, N)\}} \left( \sum_x \mathcal{W}_x \right) \quad (12)$$

where

$$\mathcal{W}_x = \begin{cases} \hat{\Phi}(i_x, j_x, n_x) + \mathcal{E}_C(i_x, j_x, \hat{D}_{ij}) & \text{if } \delta(i_x, j_x) \geq \delta(i, j) \\ -\infty & \text{otherwise.} \end{cases}$$

### 2.3 Algorithm

To evaluate equation (12), we must choose a set of non-overlapping intervals  $[i_x, j_x]$  and multiplicities  $n_x$  that maximize  $\sum_x \mathcal{W}_x$  and satisfy  $\sum n_x = N - 1$ . Similarly, to assemble a TAD forest from TAD trees, we will likewise be choosing a non-overlapping set of intervals (leaves of TAD trees) with multiplicities (number of TADs in each tree) such that the sum of their scores is maximized and the multiplicities sum to a predefined  $N$ . These tasks are both similar to the weighed interval scheduling problem (Kleinberg

Weireb & Raphael, Bioinformatics, 2015

# Hi-C data analysis: finding topologically associated domains (TADs)

## Identification of hierarchical chromatin domains

Caleb Weinreb<sup>1</sup>, and Benjamin J. Raphael<sup>1,2\*</sup>

<sup>1</sup>Center for Computational Molecular Biology, Brown University, Providence, RI

<sup>2</sup>Department of Computer Science, Brown University, Providence, RI

Associate Editor: Prof. Gunnar Ratsch

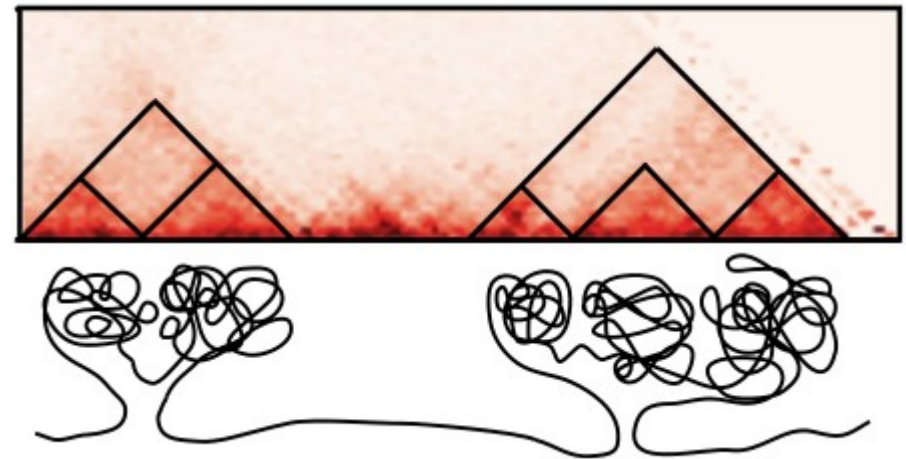
### ABSTRACT

**Motivation:** The 3D structure of the genome is an important regulator of many cellular processes including differentiation and gene regulation. Recently, technologies such as Hi-C that combine proximity ligation with high-throughput sequencing have revealed domains of self-interacting chromatin, called topologically associating domains (TADs), in many organisms. Current methods for identifying TADs using Hi-C data assume that TADs are non-overlapping, despite evidence for a nested structure in which TADs and sub-TADs form a complex hierarchy.

**Results:** We introduce a model for hierarchical decomposition of contact frequencies into a hierarchy of nested TADs. This model is based empirical distributions of contact frequencies within TADs, where positions that are far apart have a greater enrichment of contacts than

resulting in a contact matrix  $A$ , where  $A_{ij}$  is the number of contacts between bins  $i$  and  $j$ , normalized for experimental bias. Several methods have been developed for the identification of TADs from Hi-C data. These methods may be roughly classified into two categories: (1) methods that define a 1D test statistic from the contact matrix  $A_{ij}$ ; (2) methods that exploit the 2D structure of the contact matrix.

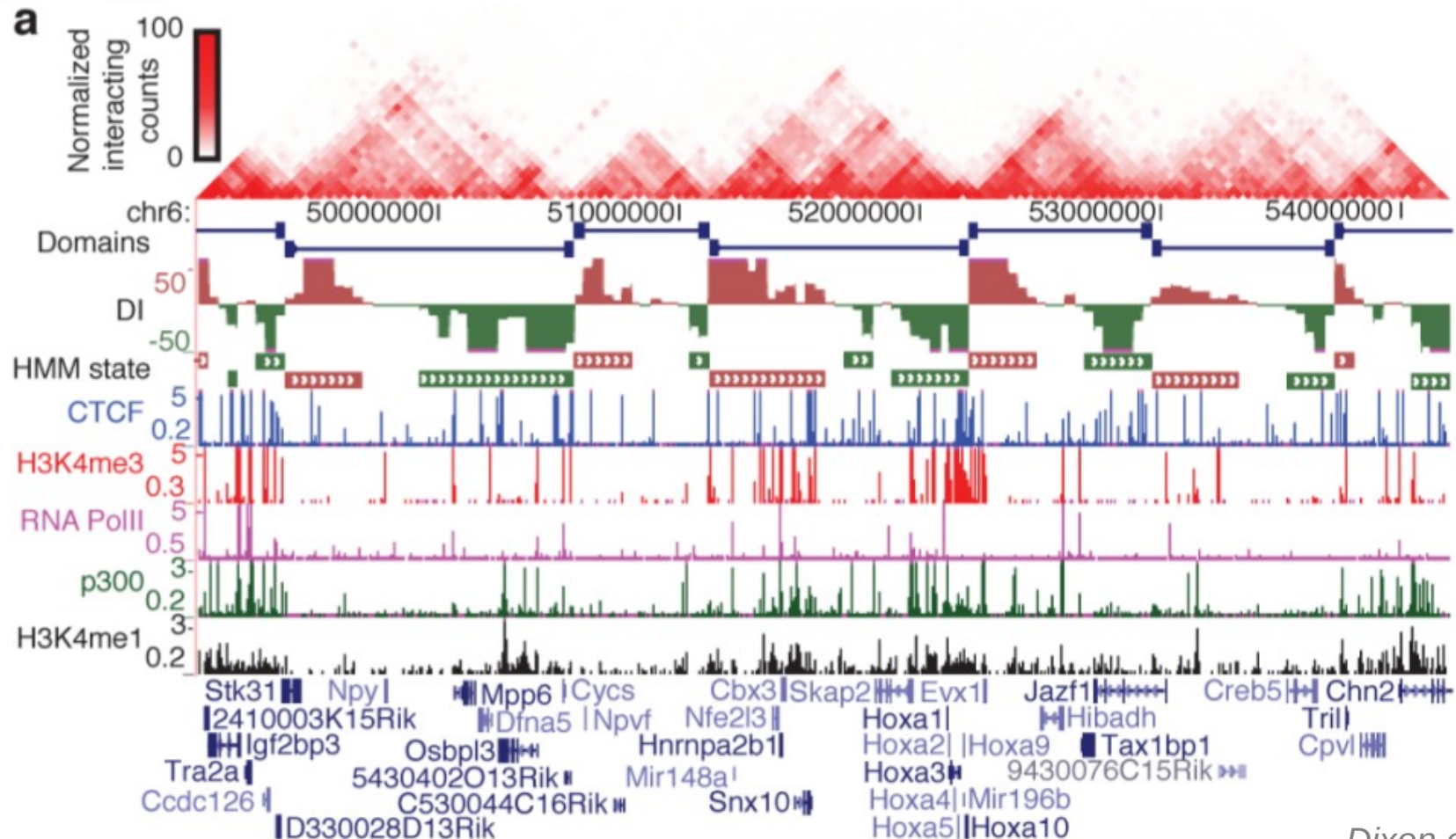
Dixon *et al.* (2012) compute a 1D “directionality index” (DI) from the contact matrix. This index defines whether contacts have an upstream bias, downstream bias or no bias. Next, they use a hidden Markov model (HMM) to partition the genome into regions defined by changes in the directionality index. Each transition into downstream bias marks the start of a domain and the next transition out of upstream bias marks its end. Sauria *et al.* (2014) introduce a 1D



Weinreb & Raphael, *Bioinformatics*, 2015

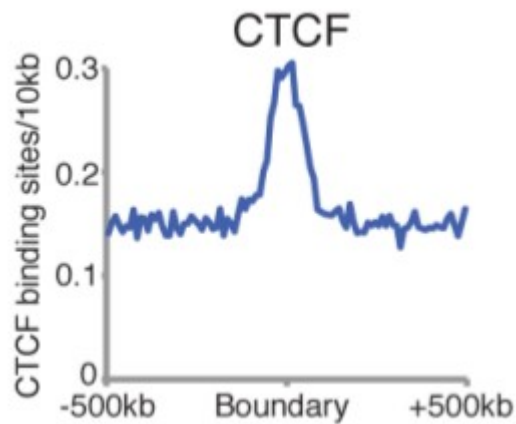


# Hi-C data analysis: integrative analysis

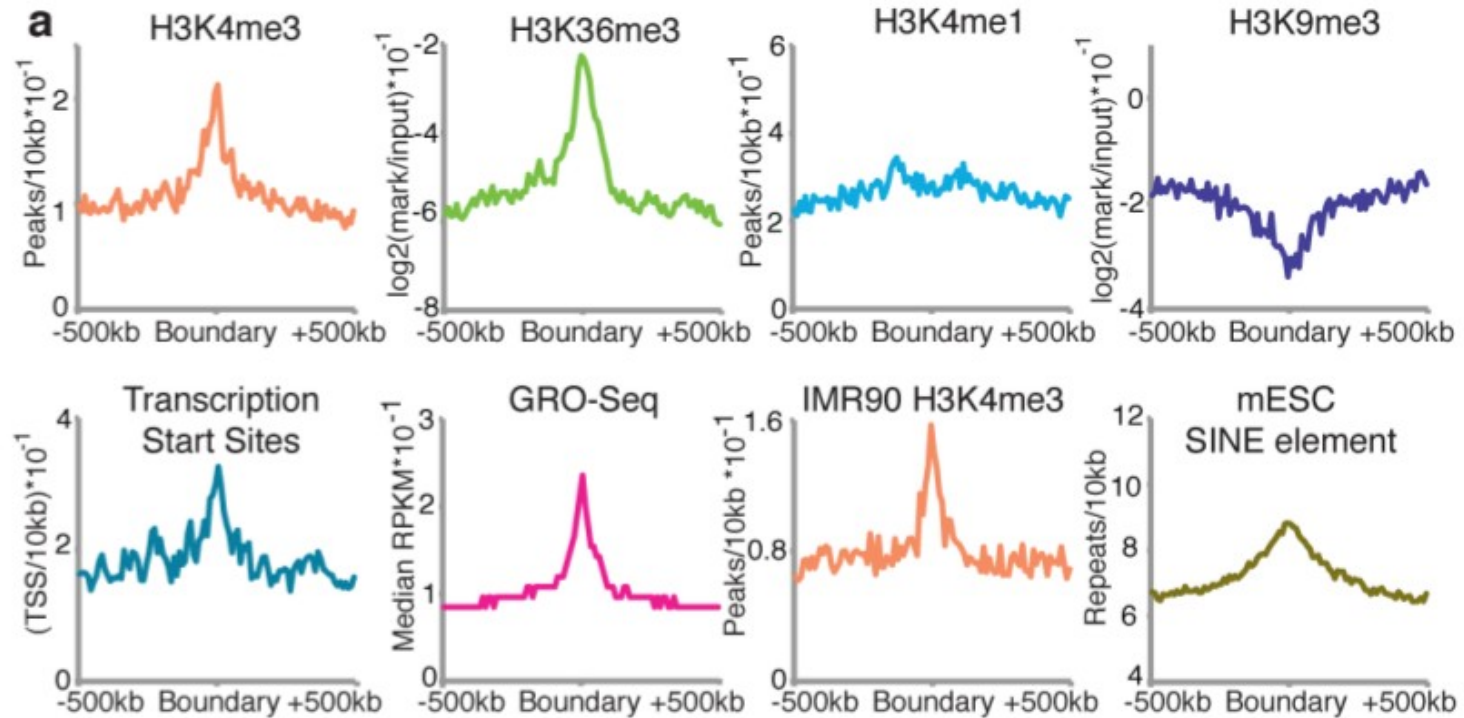


*Dixon et al 2012*

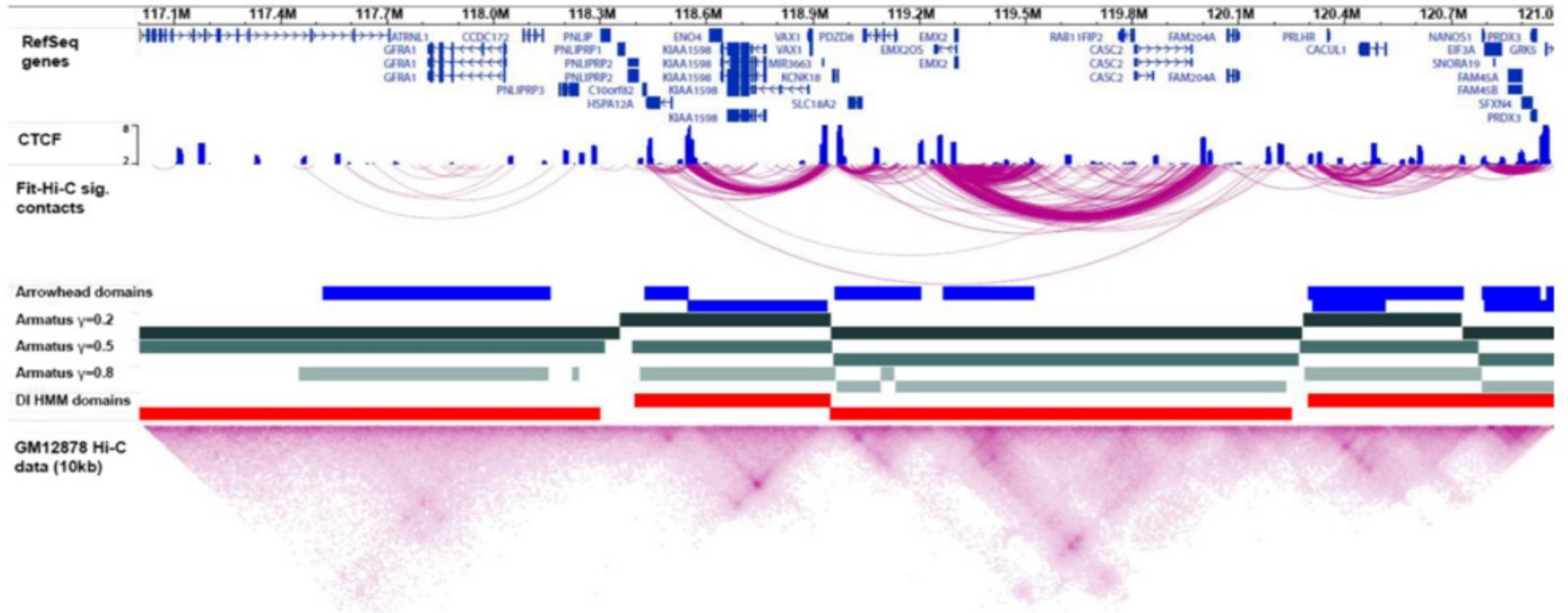
# Hi-C data analysis: comparisons



*Dixon et al 2012*



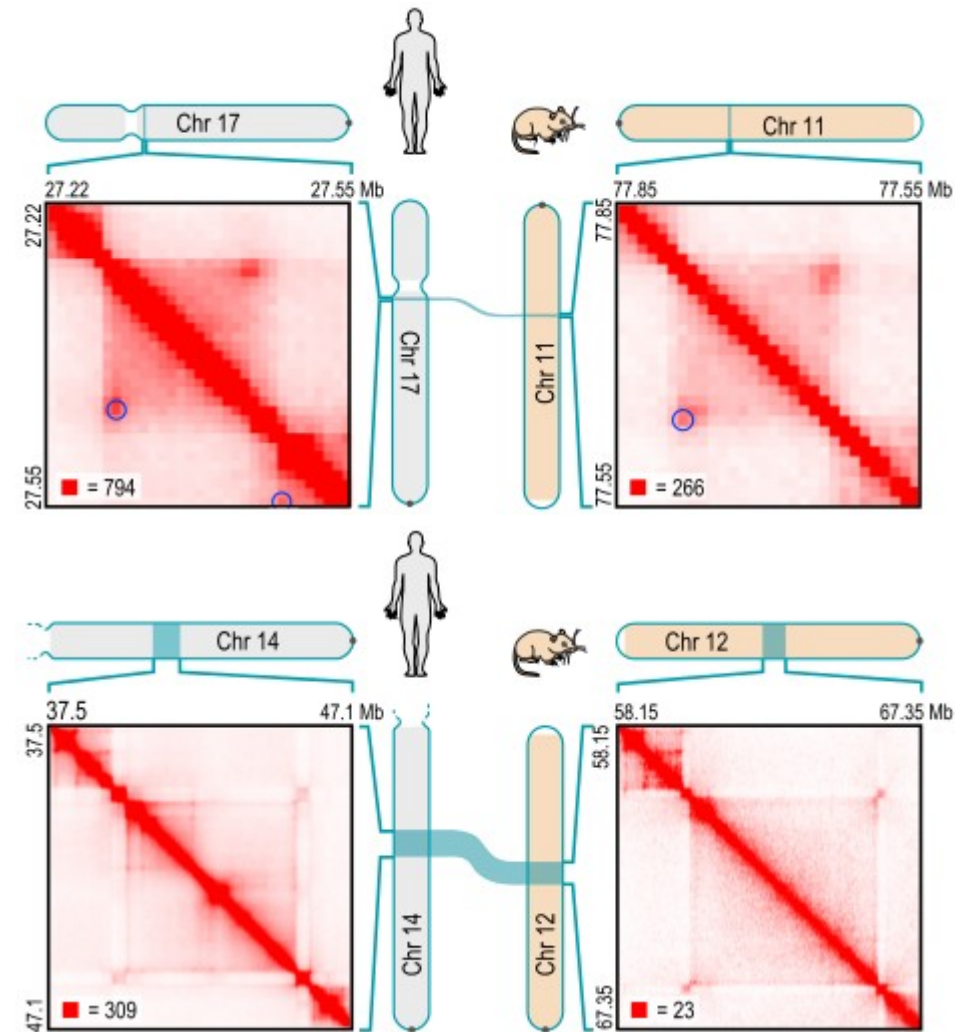
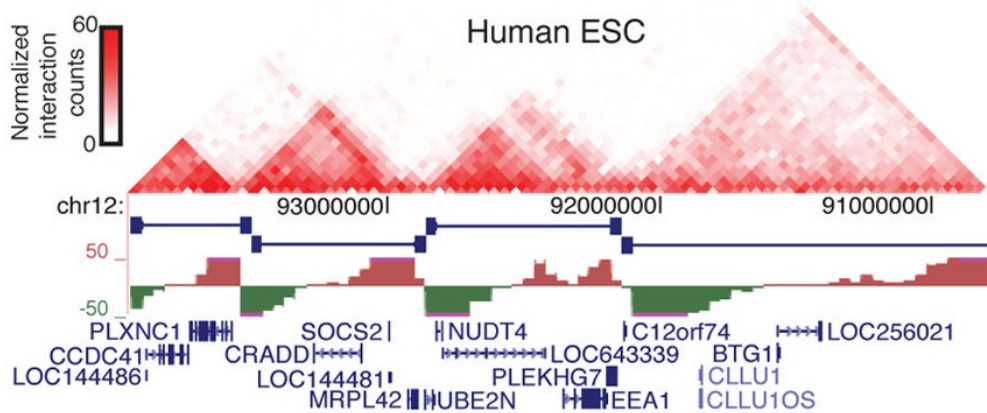
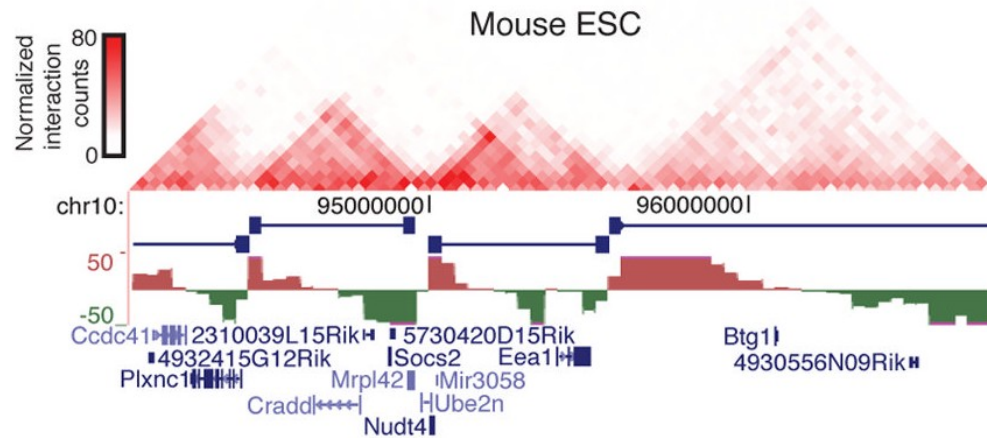
# Hi-C data analysis: integrative analysis



Ay & Noble, *Genome Biology*, 2015



# Hi-C data analysis: integrative analysis





# FR-AgENCODE: livestock genome annotation

2x♂  
2x♀



*Sus scrofa*  
(Large White)



*Gallus gallus*  
(White Leghorn)



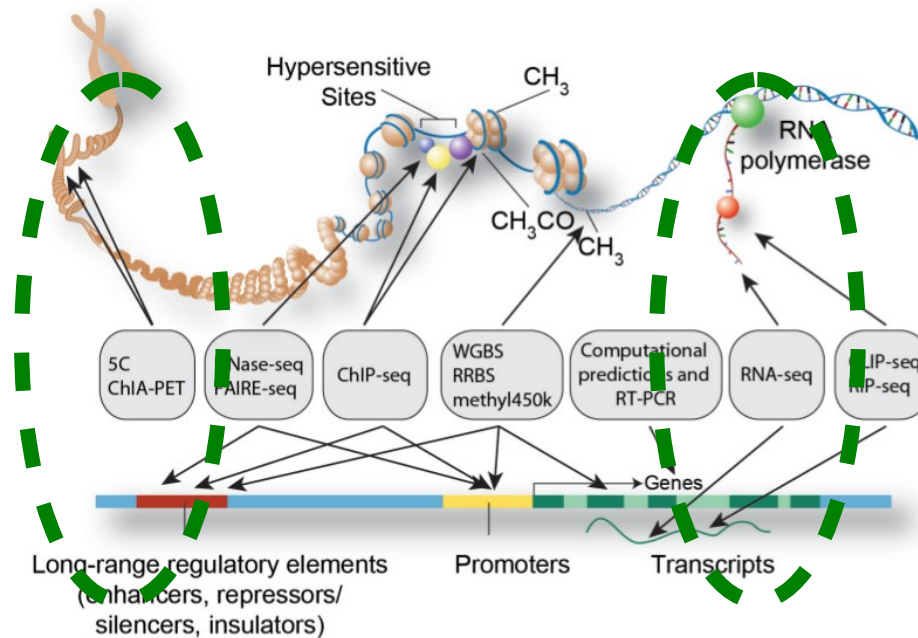
*Bos Taurus*  
(Holstein)



*Capra hircus*  
(Alpine)

- ◆ Sampling: 34 somatic tissues + 13 reproductive tissues  
=> INRA CRB-Anim biorepository

- ◆ Molecular assays
  - ◆ RNA-seq
  - ◆ Hi-C
- ◆ Data analysis

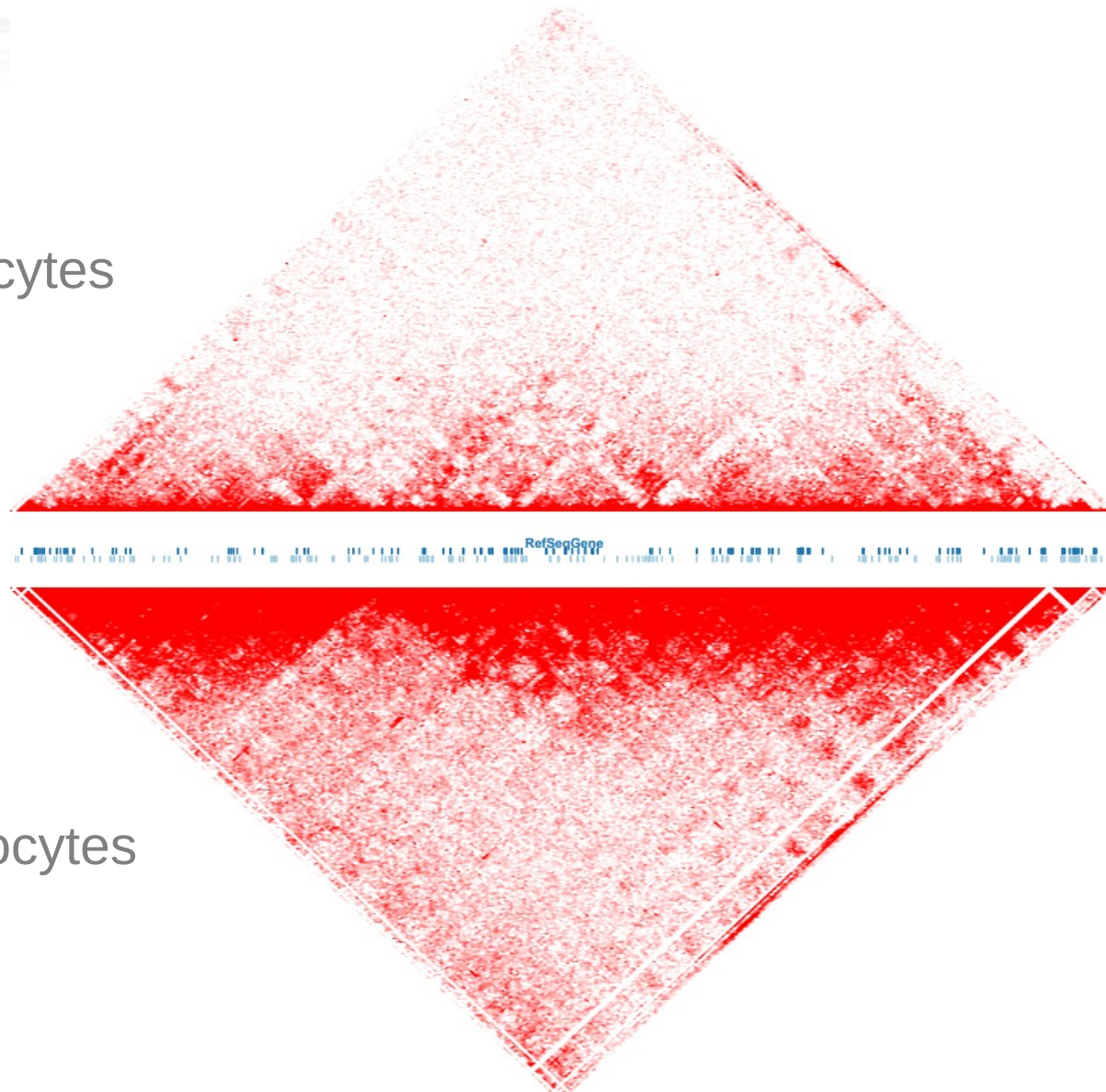


Modified from PLoS Biol 9-e1001046, 2011 & Science 306:636, 2004  
Image credits: Darryl Leja (NHGRI), Ian Dunham (EBI), Michael Pazin (NHGRI)

**Check out:**  
[www.faang.org](http://www.faang.org)

# FR-AgENCODÉ: livestock genome annotation

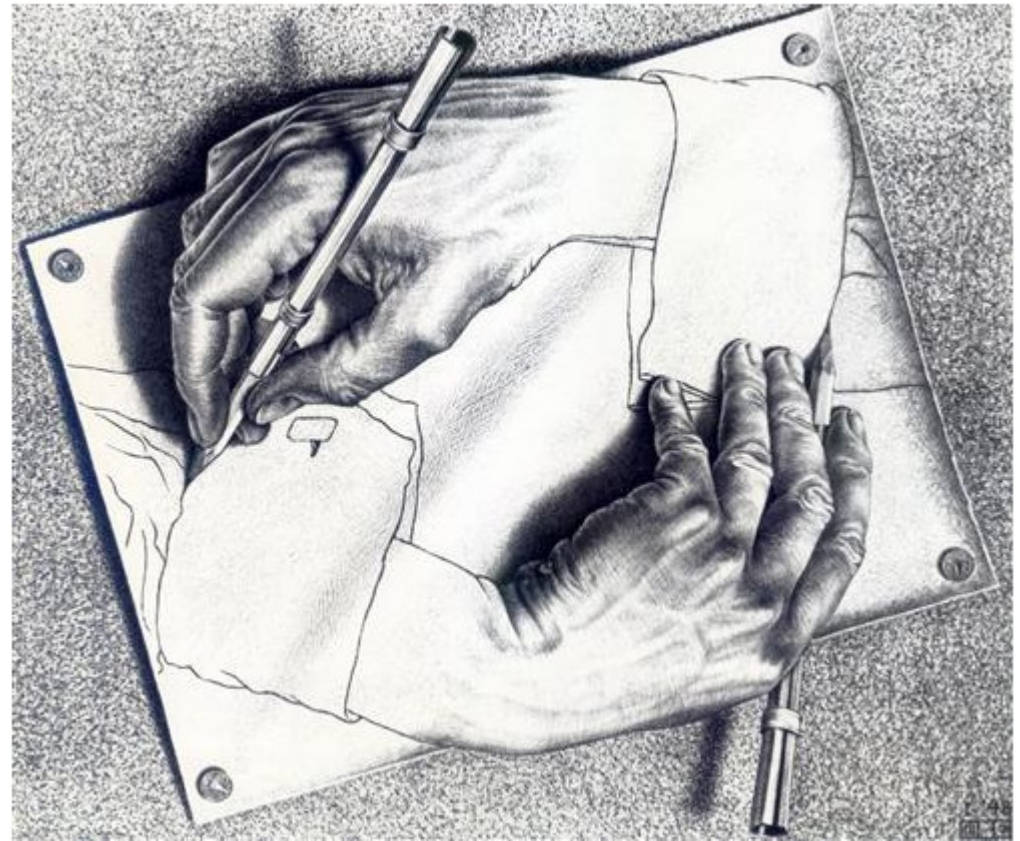
Pig hepatocytes



Pig lymphocytes

# Outline

- ◆ Biological context
- ◆ More biological context
- ➔ Hi-C data processing
  - ◆ map
  - ◆ filter
  - ◆ count
  - ◆ normalize
  - ◆ segment
  - ◆ compare
- ◆ Conclusion, discussion, NETBIO lunch

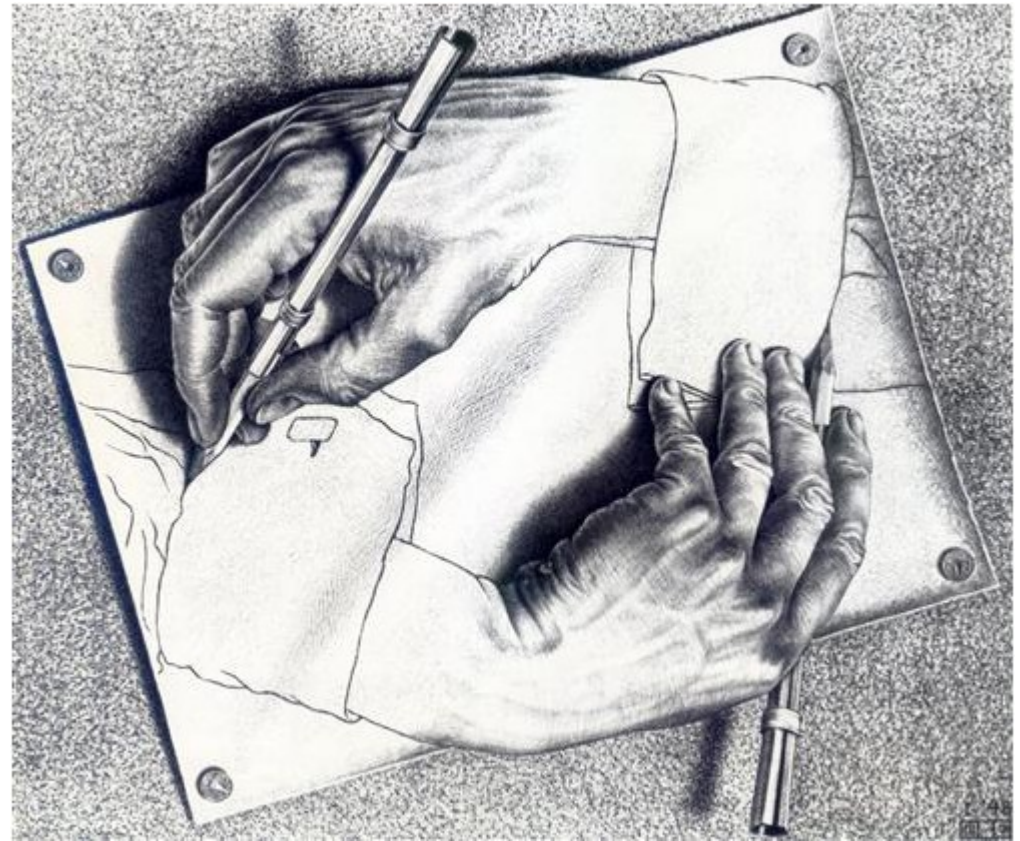


*M.C. Escher, 1948*



# Outline

- ◆ Biological context
- ◆ More biological context
- ◆ Hi-C data processing
  - ◆ map
  - ◆ filter
  - ◆ count
  - ◆ normalize
  - ◆ segment
  - ◆ compare



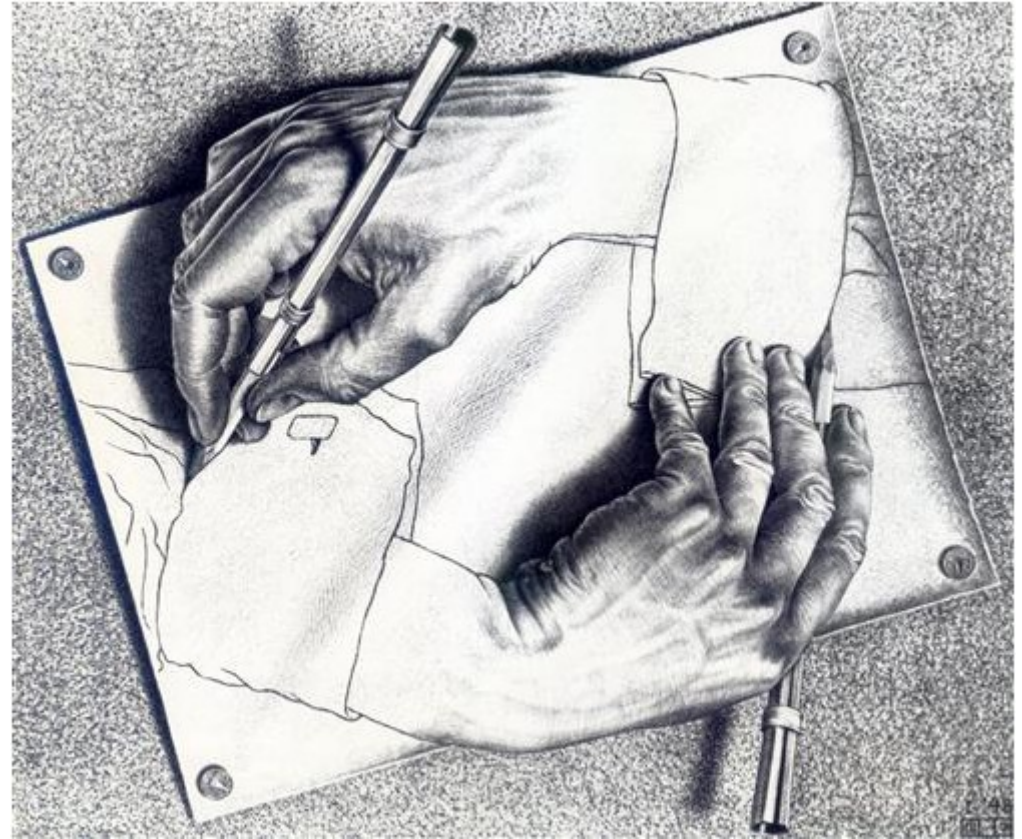
*M.C. Escher, 1948*

➔ Conclusion, discussion, NETBIO lunch



# Acknowledgements

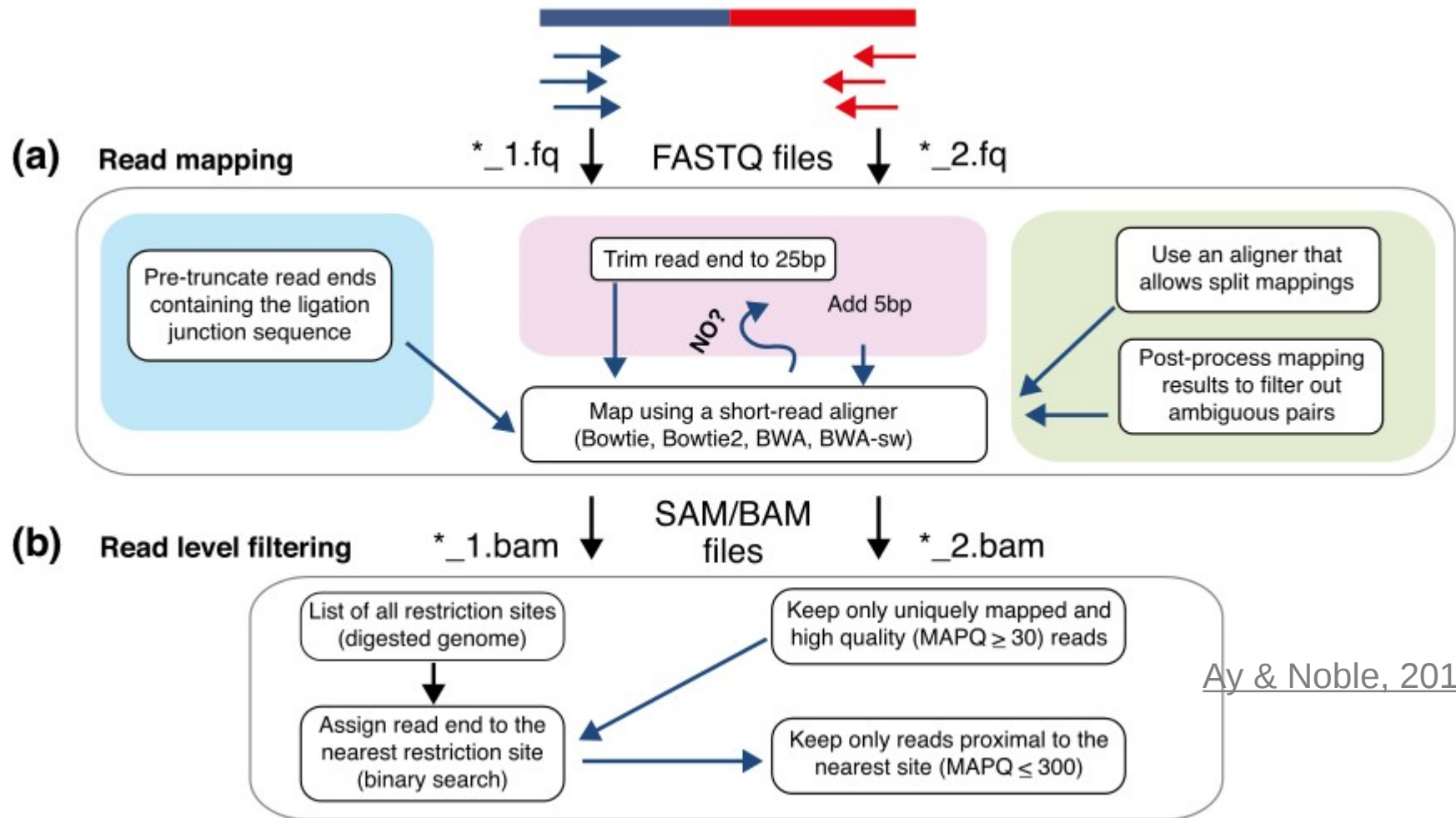
- ◆ NETBIO
- ◆ GenPhySE
  - ◆ Hervé Acloque
  - ◆ Diane Esquerré
  - ◆ Magali San Cristobal
  - ◆ Matthias Zytnicki
  - ◆ David Robelin
  - ◆ Martine Yerle
  - ◆ Yvette Lahbib
- ◆ Nicolas Servant
- ◆ FR-AgENCODE
  - ◆ Elisabetta Giuffra
- ◆ FAANG consortium
- ◆ Etc (sorry!)



*M.C. Escher, 1948*

# Analysis: example on Hi-C data

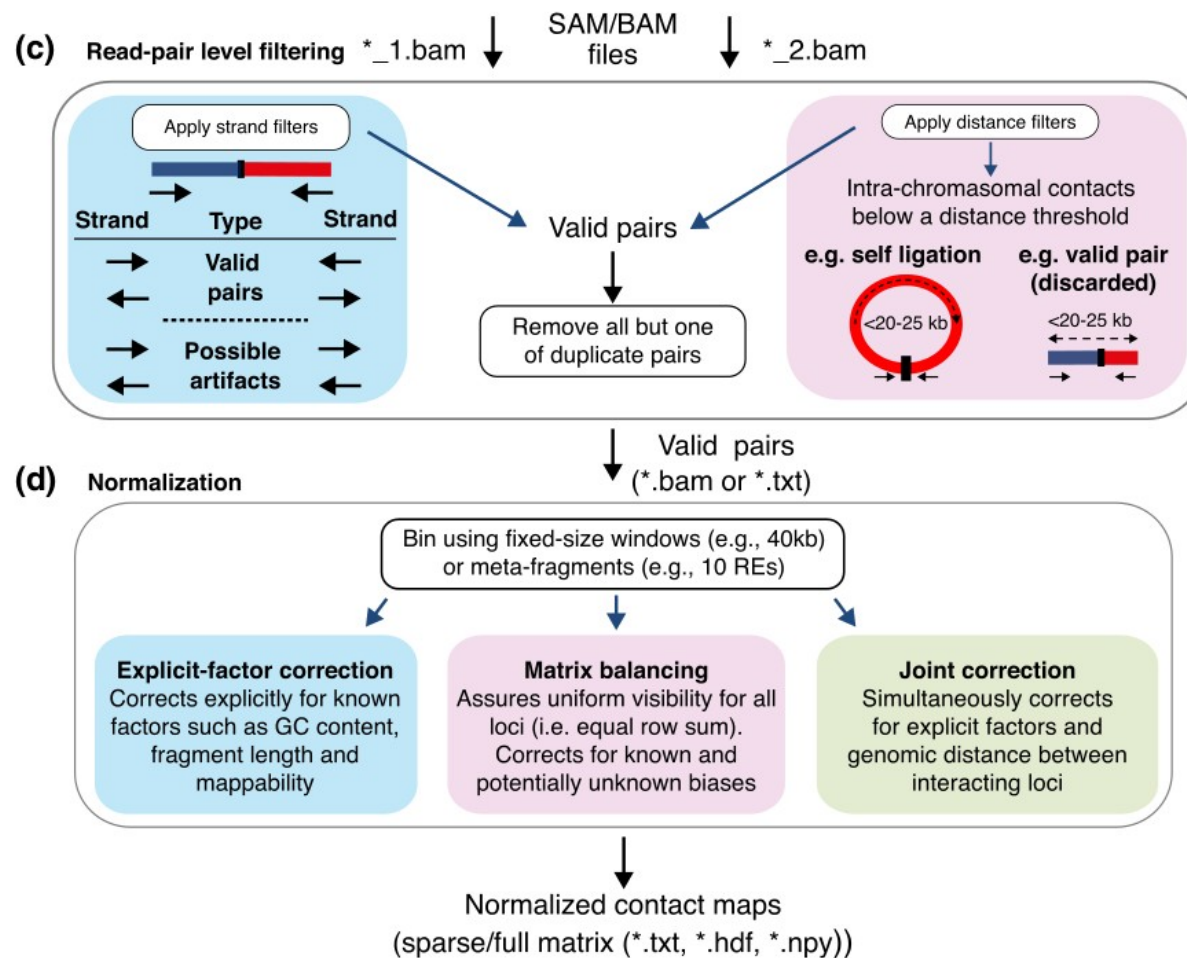
## Analysis pipeline (first part)



Ay & Noble, 2015

# Analysis: example on Hi-C data

## ◆ Analysis pipeline (first part)



Ay & Noble, 2015

◆ + TAD calling, differential analysis, integrative analysis, ...