

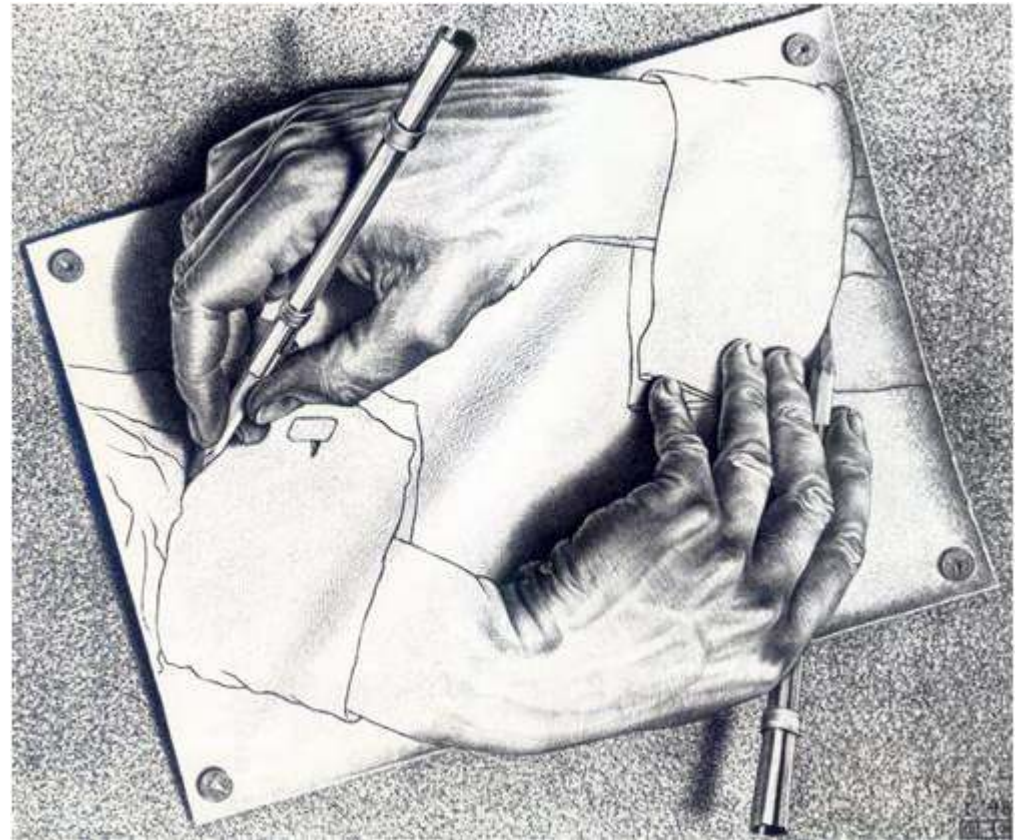
# 3D genome sequencing

## Hi-C data analysis



# Outline

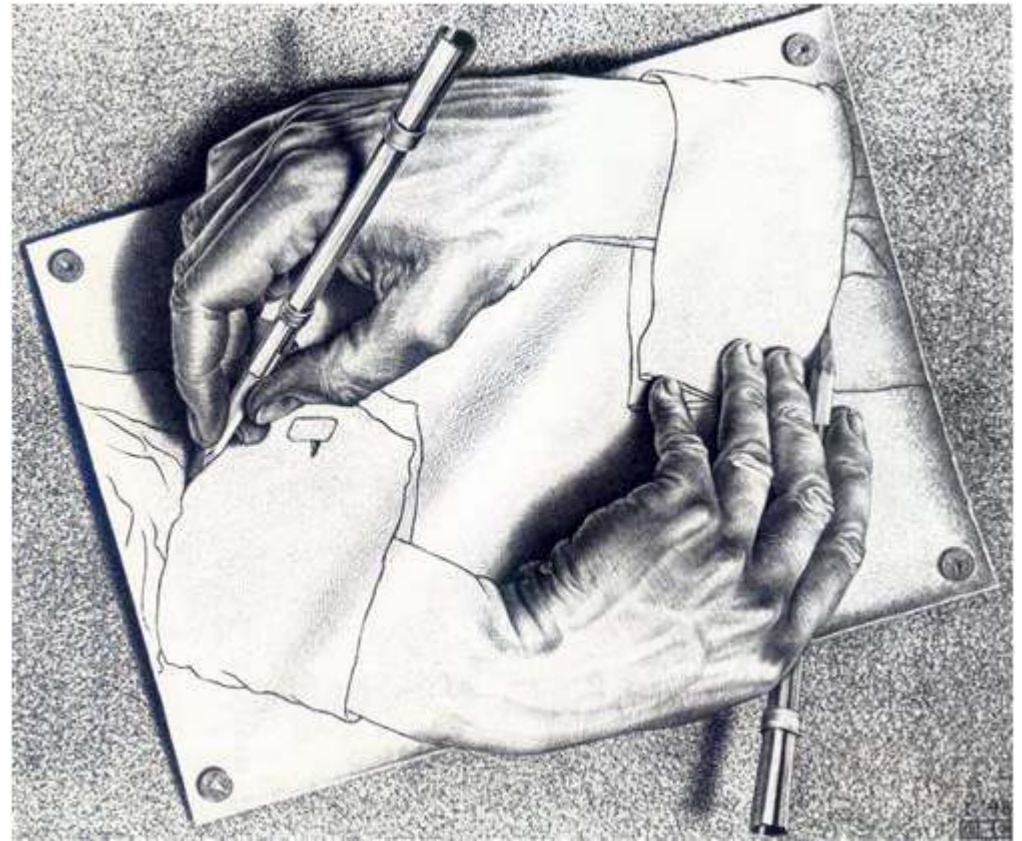
- ◆ Why
- ◆ What
- ◆ How
- ◆ Where
- ◆ Who



*M.C. Escher, 1948*

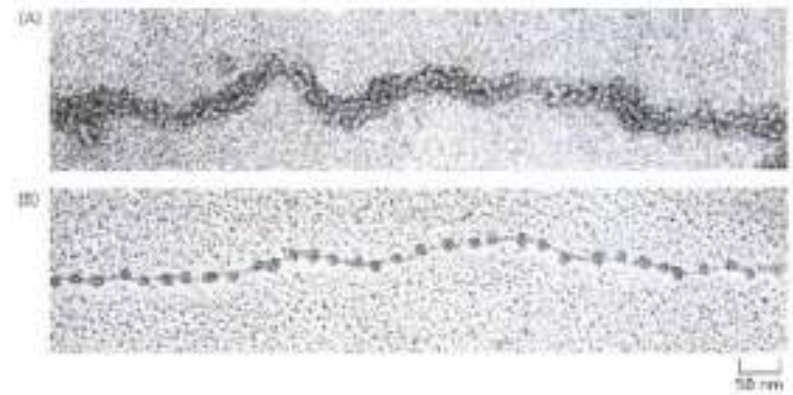
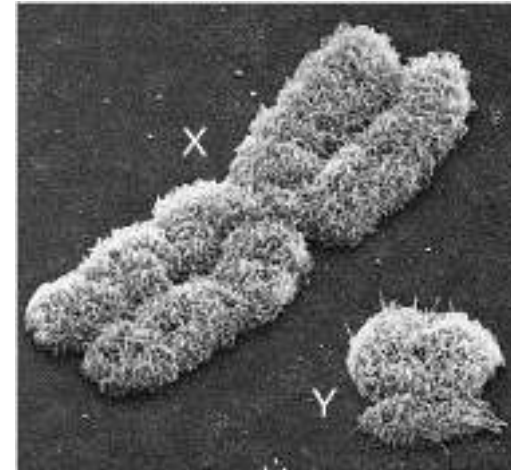
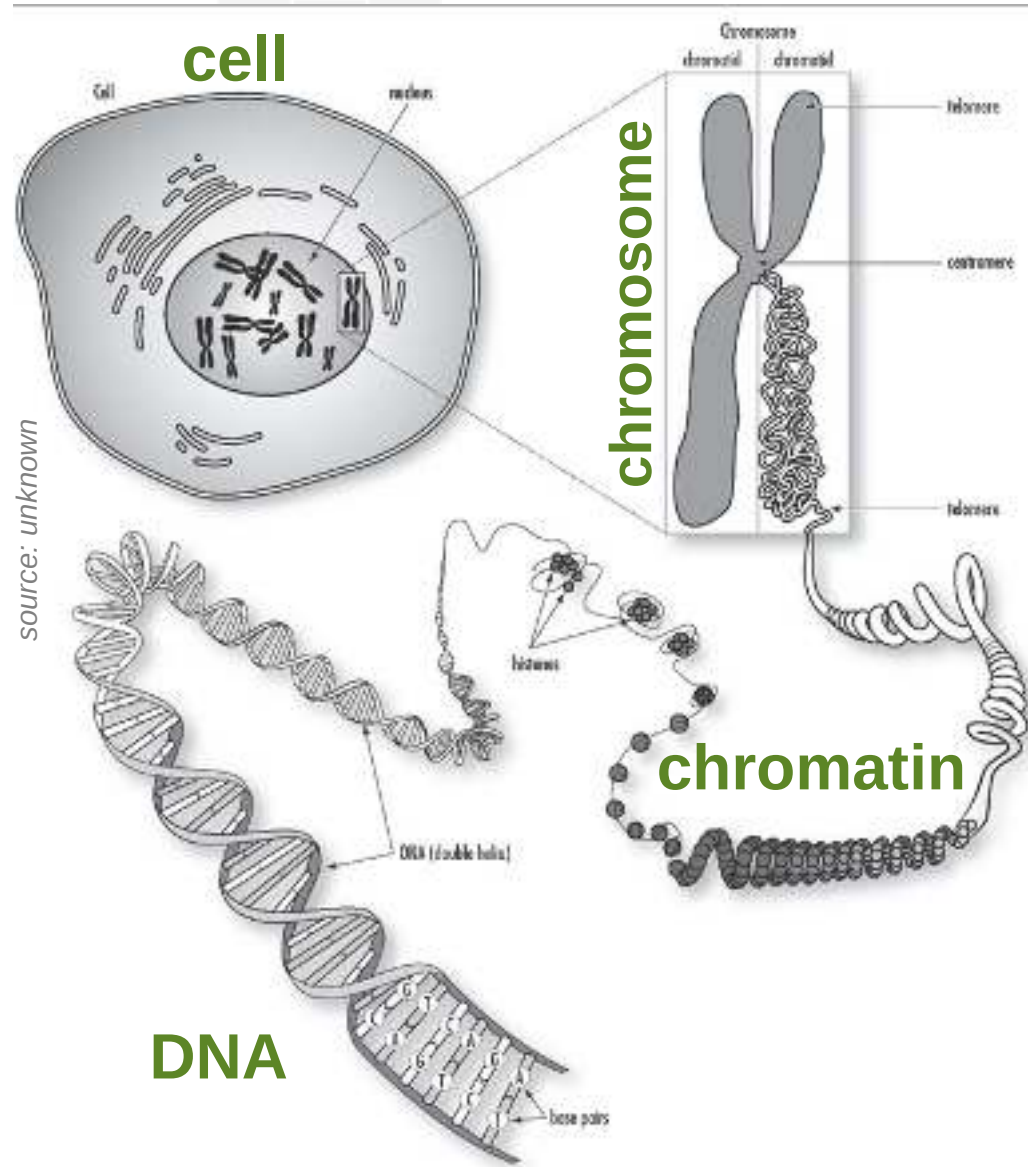
# Outline

- ◆ Why ←
- ◆ What
- ◆ How
- ◆ Where
- ◆ Who



*M.C. Escher, 1948*

# Life, cell, chromosome & DNA





# Life, cell, chromosome & DNA

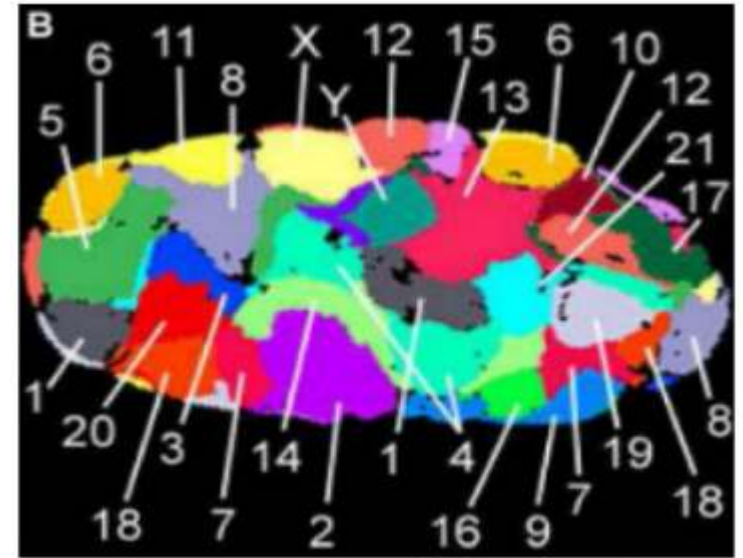


?

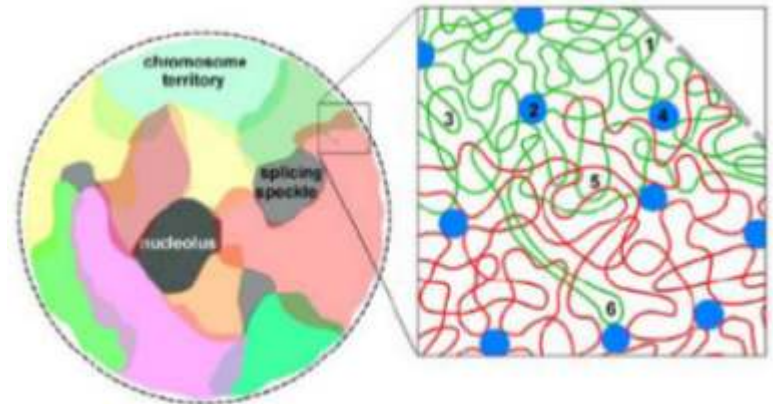
# Life, cell, chromosome & DNA



≠



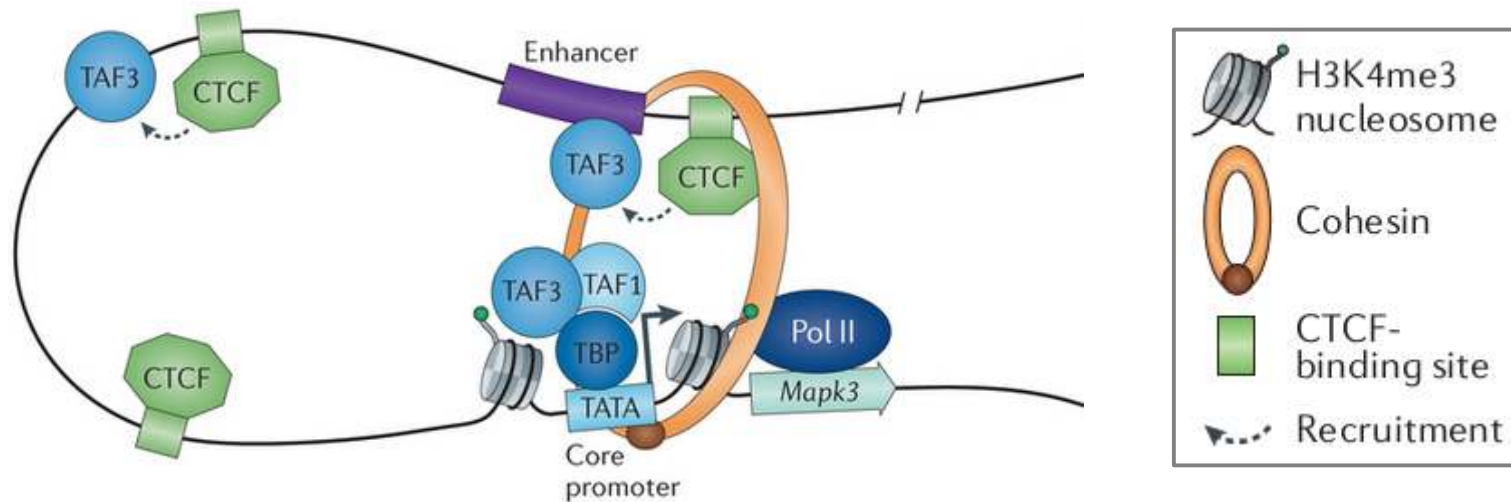
*Bolzer A et al. 2005*



*Branco MR & Pombo A, 2006*

**Genome 3D structure is organized**

# From structure to function



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

**3D structure regulates 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>2,4</sup> and Musa M. Mhlanga<sup>1,2\*</sup>

<sup>1</sup>Gene Expression and Biophysics Group, Synthetic Biology Emerging Research Area, Biociences 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: ymde@mhlabs.org  
<http://dx.doi.org/10.1016/j.cell.2015.06.051>

### SUMMARY

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

2015). These highly sensitive assays can nascent mRNA and have revealed the FISH foci in a fraction of the population

Nucleic Acids Research Advance Access published February 4, 2015

*Nucleic Acids Research* 2015, 43, 1093-1103  
doi:10.1093/nar/gkv046



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

Dario 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 Borschliwer,<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>1,17</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

## 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>1,4,\*</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

22. G. A. Gray, *Nat. Rev. Genet.* 5, 338-345 (2004).  
22. S. B. Carroll, *Dev. Biol.* 25, 25-38 (2008).

and K. P. Ziegler and the Exeter Sequencing Service for genome sequencing services. This work was also supported

## 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

Learning from Previews

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

Regard et al. *Nat. Genet.* 47, 1155-1163 (2015)  
\*Only available in the Cell Research Network. For more information, visit [www.cell.com](http://www.cell.com)

Structural variants are common in the human genome, but their contributions to human disease are unclear. We used CRISPR-Cas9 to disrupt some structural variants and found that they can alter gene expression in olfactory enhancer-proximal transducers, altered spatial-temporal gene expression, and developmental disorders.

### 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> Diyi An,<sup>4</sup> Ramya Ravindran,<sup>1,2</sup> Esteban G. Mazzoni,<sup>1,2</sup> Dmitry 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 CTCF-binding factor (CTCF) functions

<sup>1</sup>Physical Biosciences Division, Lawrence Berkeley National Lab

3D structure regulates gene expression

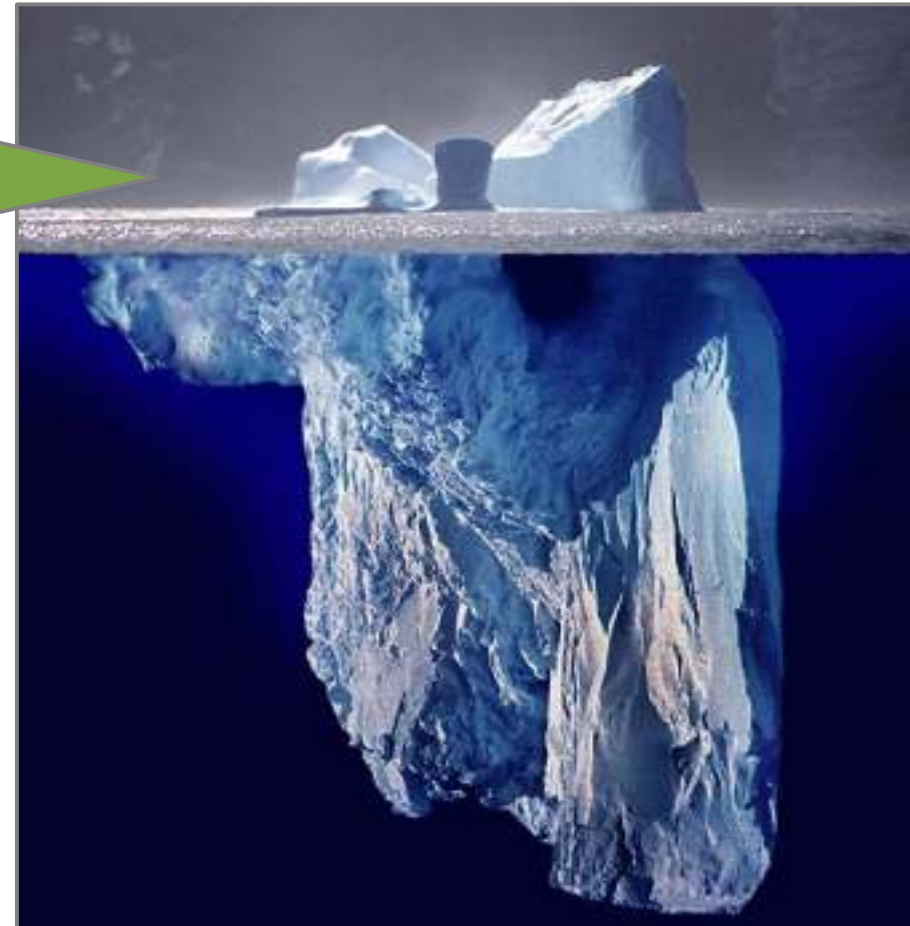




# From structure to function



?



**3D structure regulates gene expression**

# From structure to function



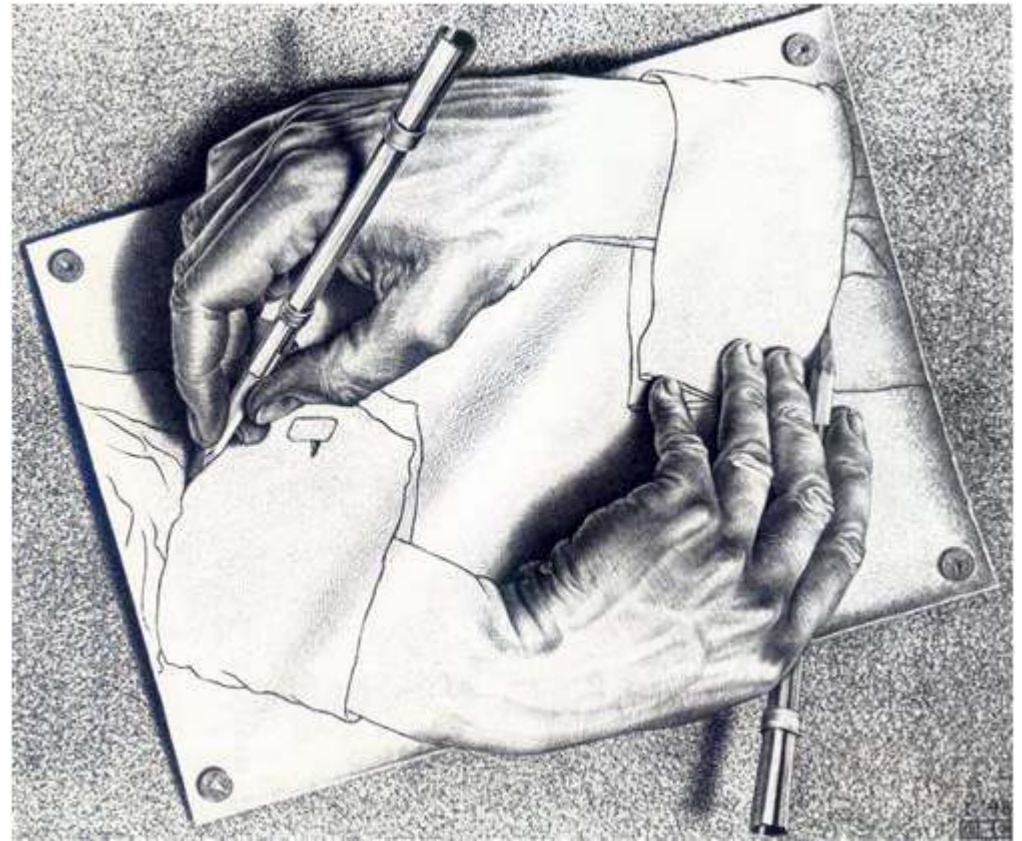
?



**3D structure regulates gene expression**

# Outline

- ◆ Why ←
- ◆ What
- ◆ How
- ◆ Where
- ◆ Who

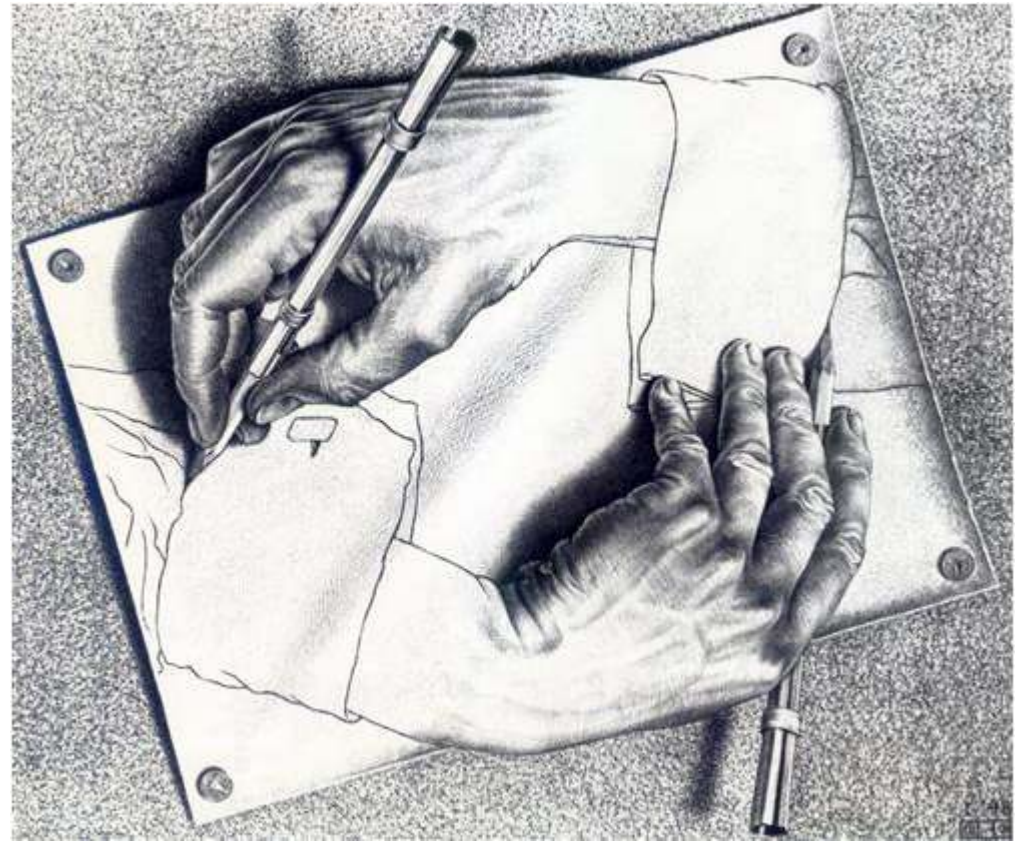


*M.C. Escher, 1948*



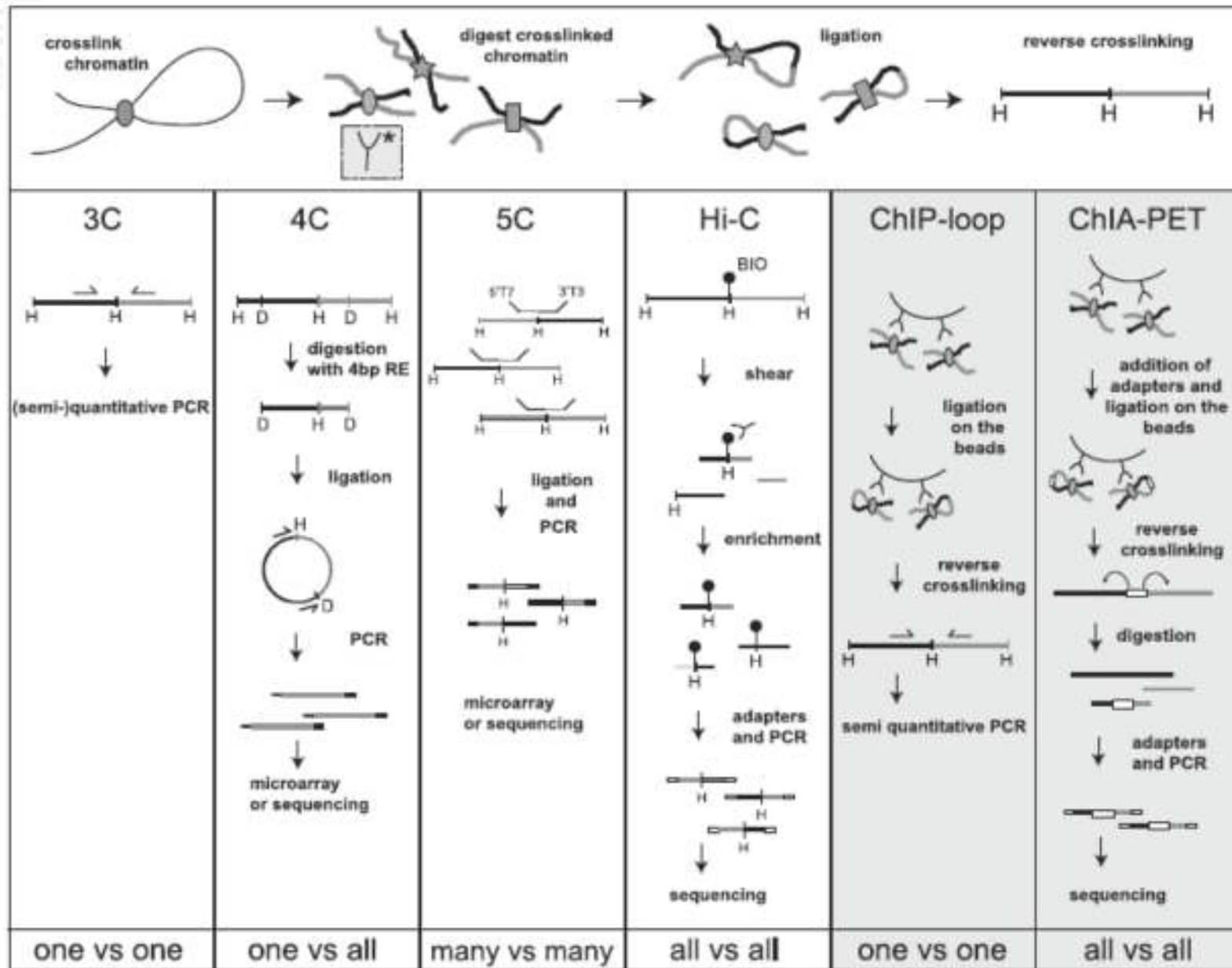
# Outline

- ◆ Why
- ◆ What ←
- ◆ How
- ◆ Where
- ◆ Who



*M.C. Escher, 1948*

# Chromosome Conformation Capture assays

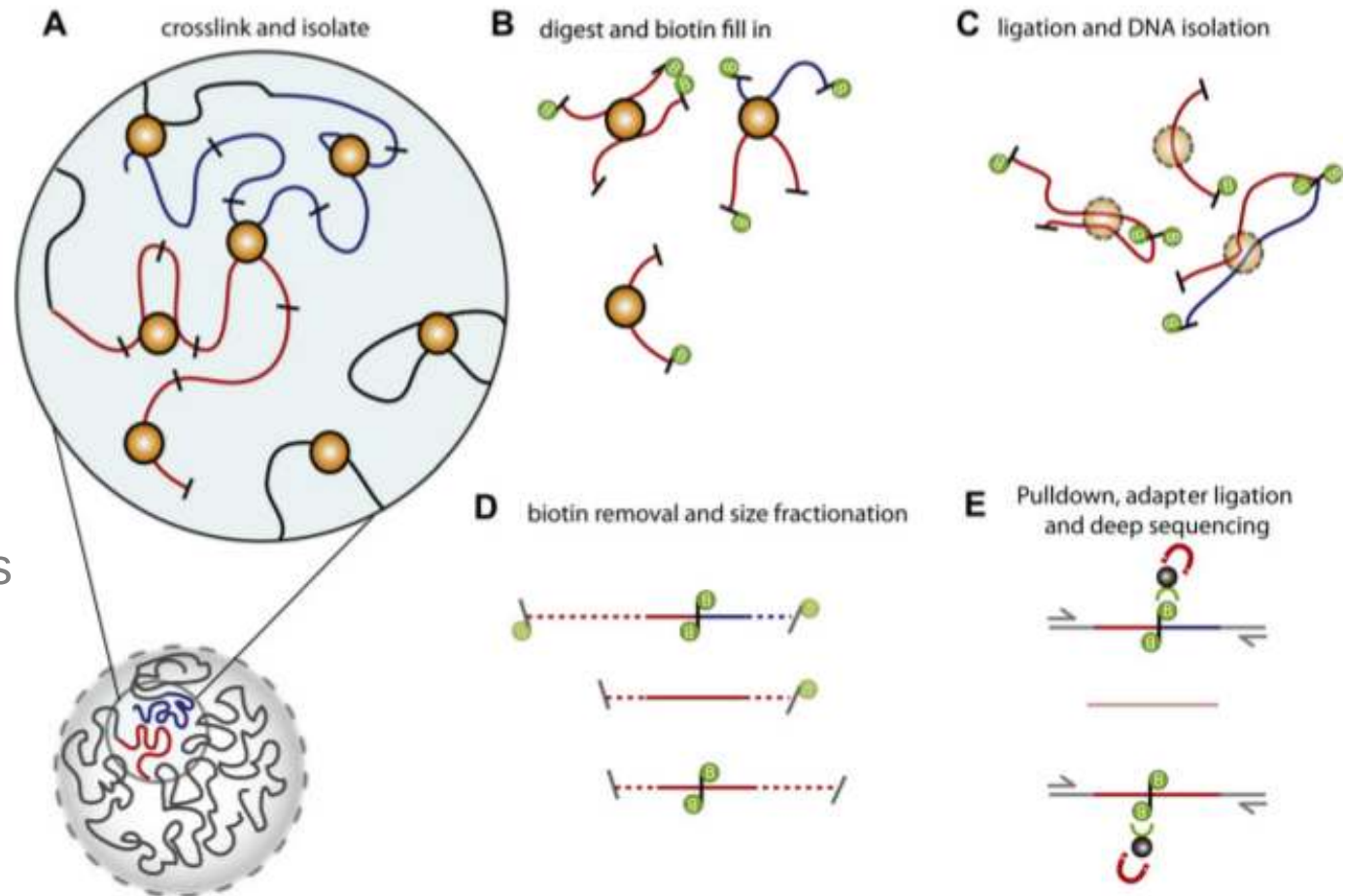


# Hi-C: the experiment

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*

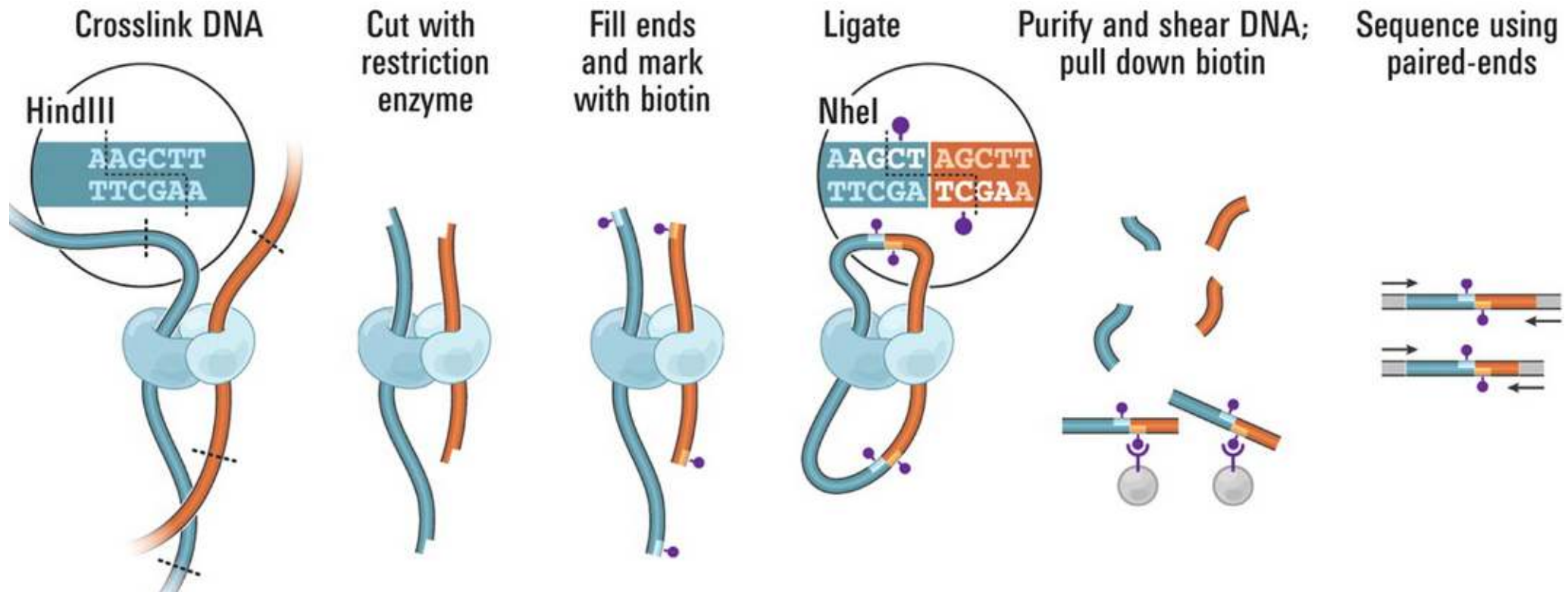
- ◆ crosslink DNA (“fixation”)
- ◆ cleave genome with restriction enzyme
- ◆ biotin-mark and ligate extremities
- ◆ fragment, select biotin-marked junctions
- ◆ sequence fragments (paired-ends)





# Hi-C: the experiment

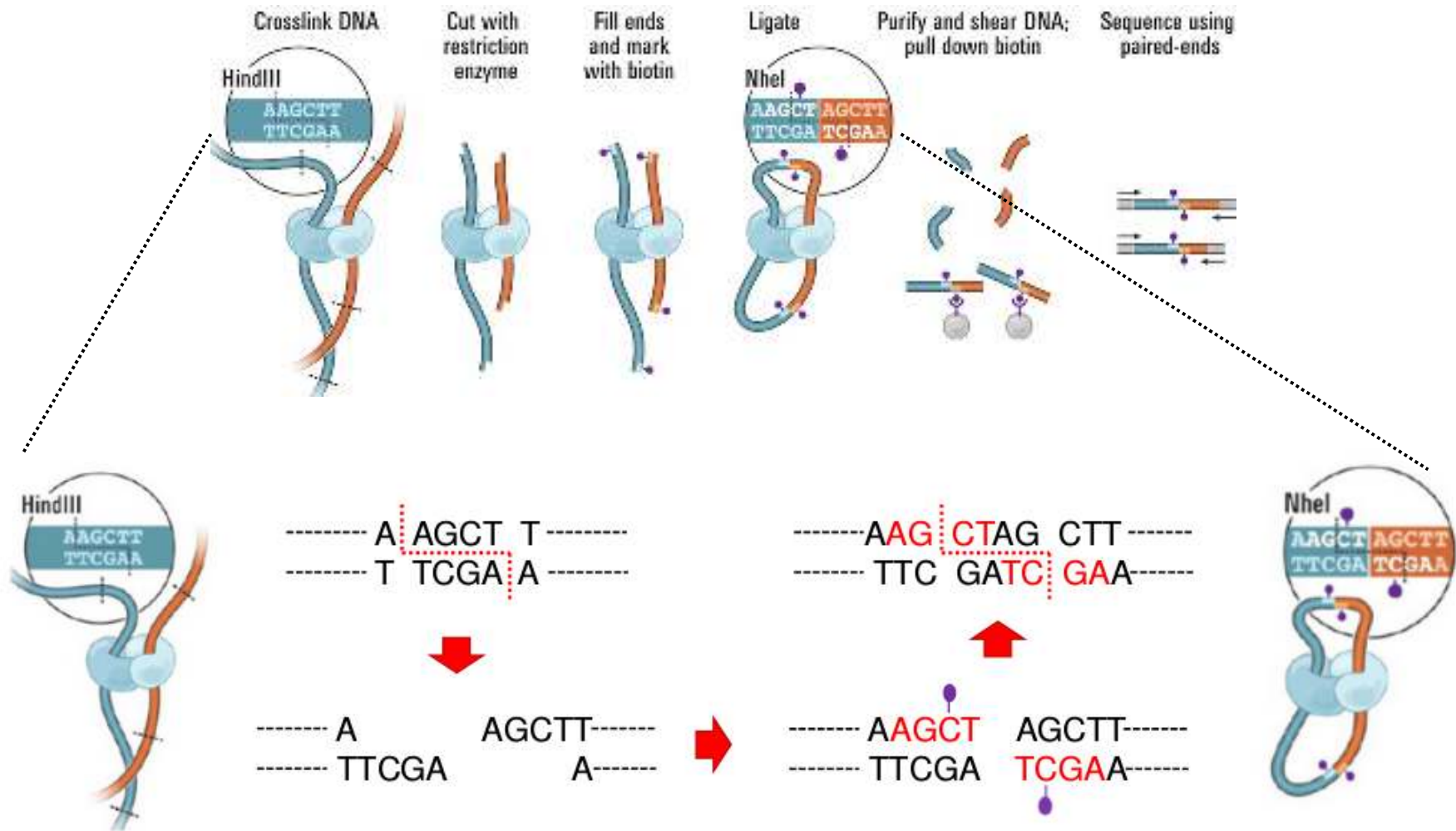
Hi-C: high-throughput chromatin conformation capture  
(Lieberman-Aiden et al, Science, 2009, Rao et al, Cell, 2014)



*Rao et al, Cell, 2014*

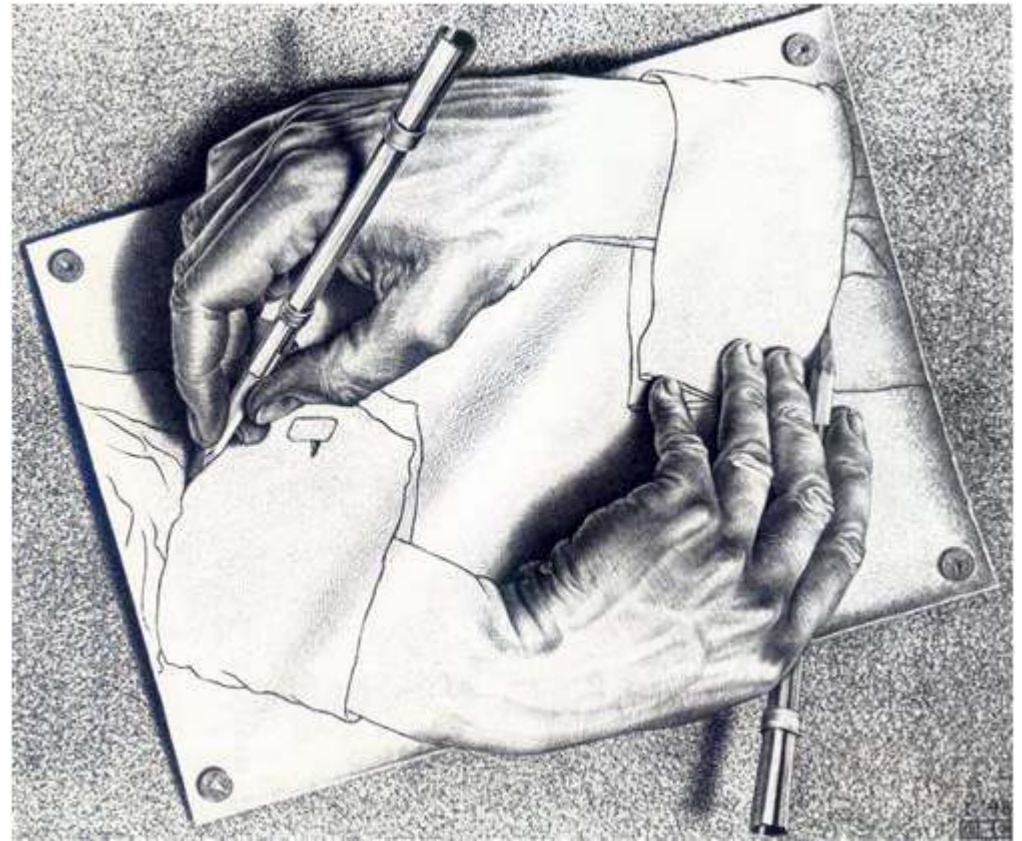
# Hi-C: the experiment

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



# Outline

- ◆ Why
- ◆ What ←
- ◆ How
- ◆ Where
- ◆ Who

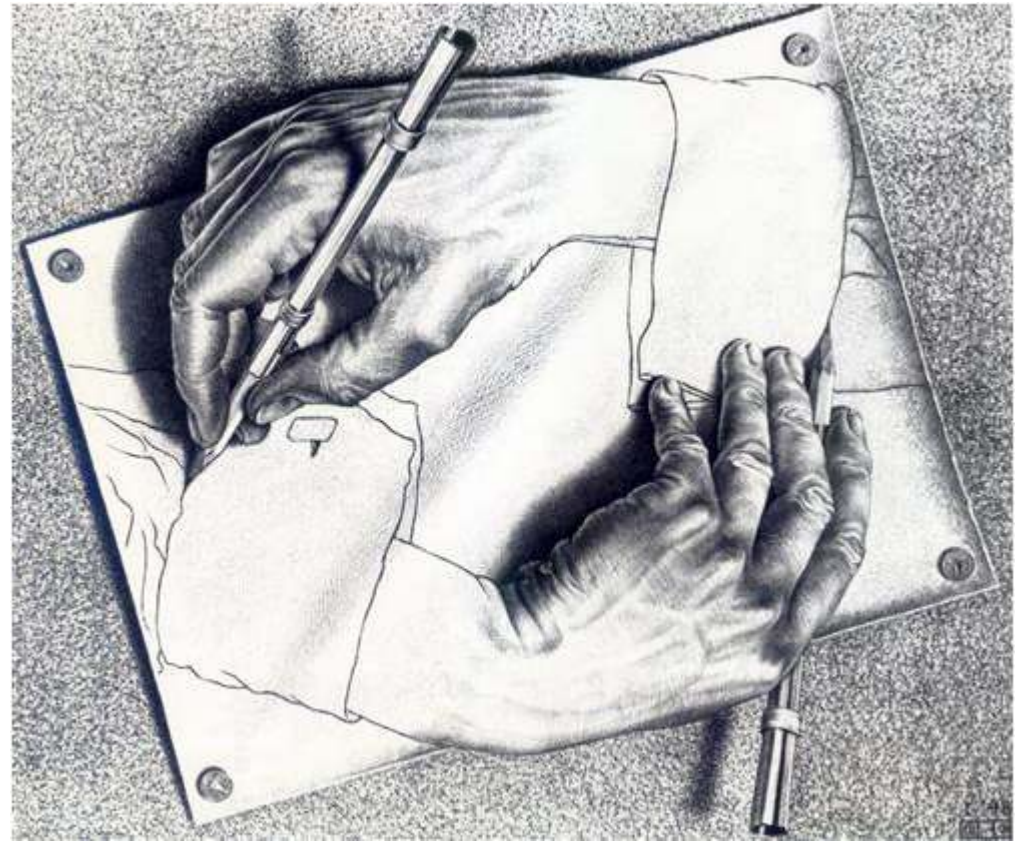


*M.C. Escher, 1948*



# Outline

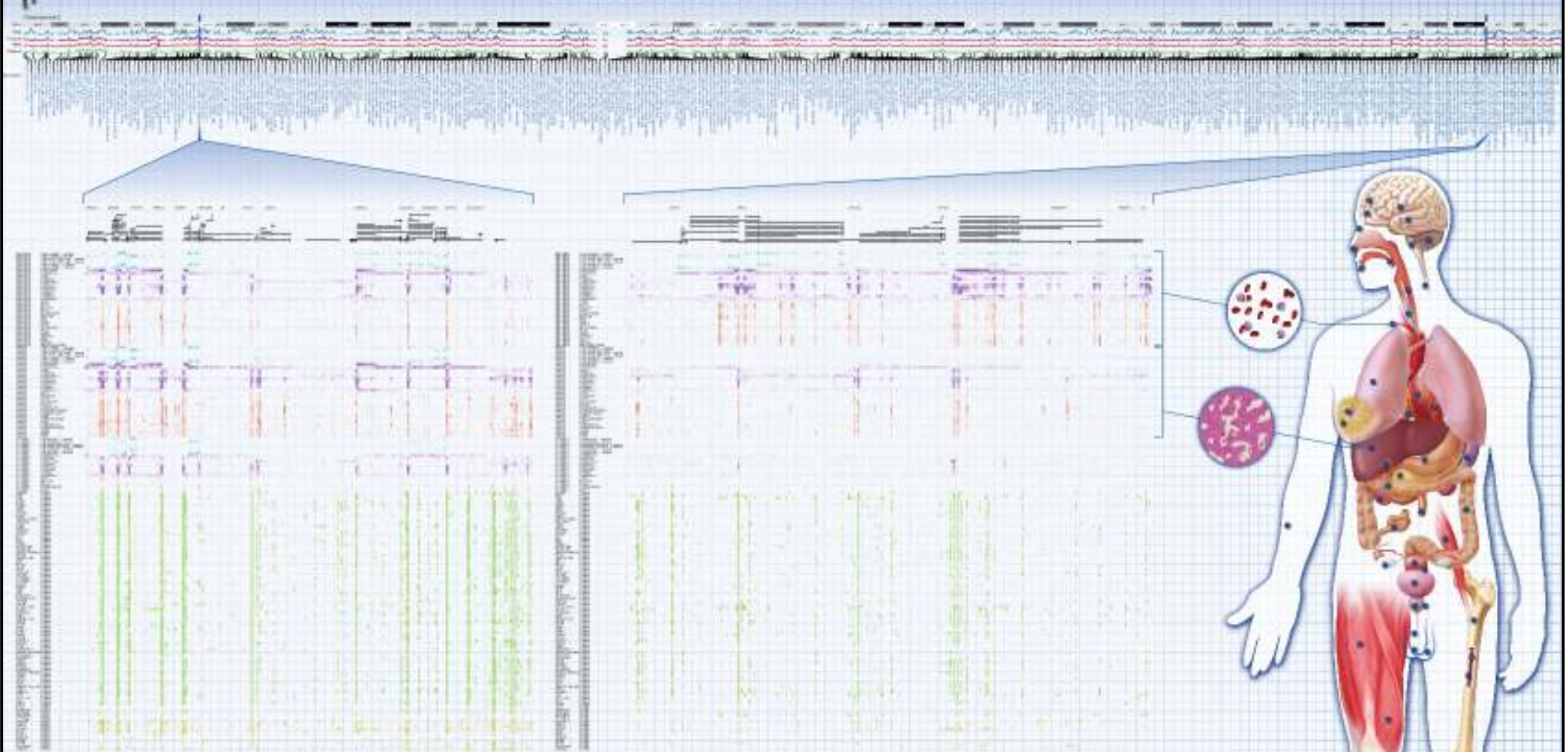
- ◆ Why
- ◆ What
- ◆ How ←
- ◆ Where
- ◆ Who



*M.C. Escher, 1948*

# the ENCODE project

## ENCYCLOPEDIA OF DNA ELEMENTS



Produced in association with  
**nature**

Produced with support from  
**illumina**

The genome browser interface of the ENCODE data can be accessed at <http://encode.igb.uci.edu>. The tracks shown in this figure are a subset of the data available in the ENCODE browser. The tracks shown in this figure are a subset of the data available in the ENCODE browser. The tracks shown in this figure are a subset of the data available in the ENCODE browser.

[www.nature.com](http://www.nature.com)



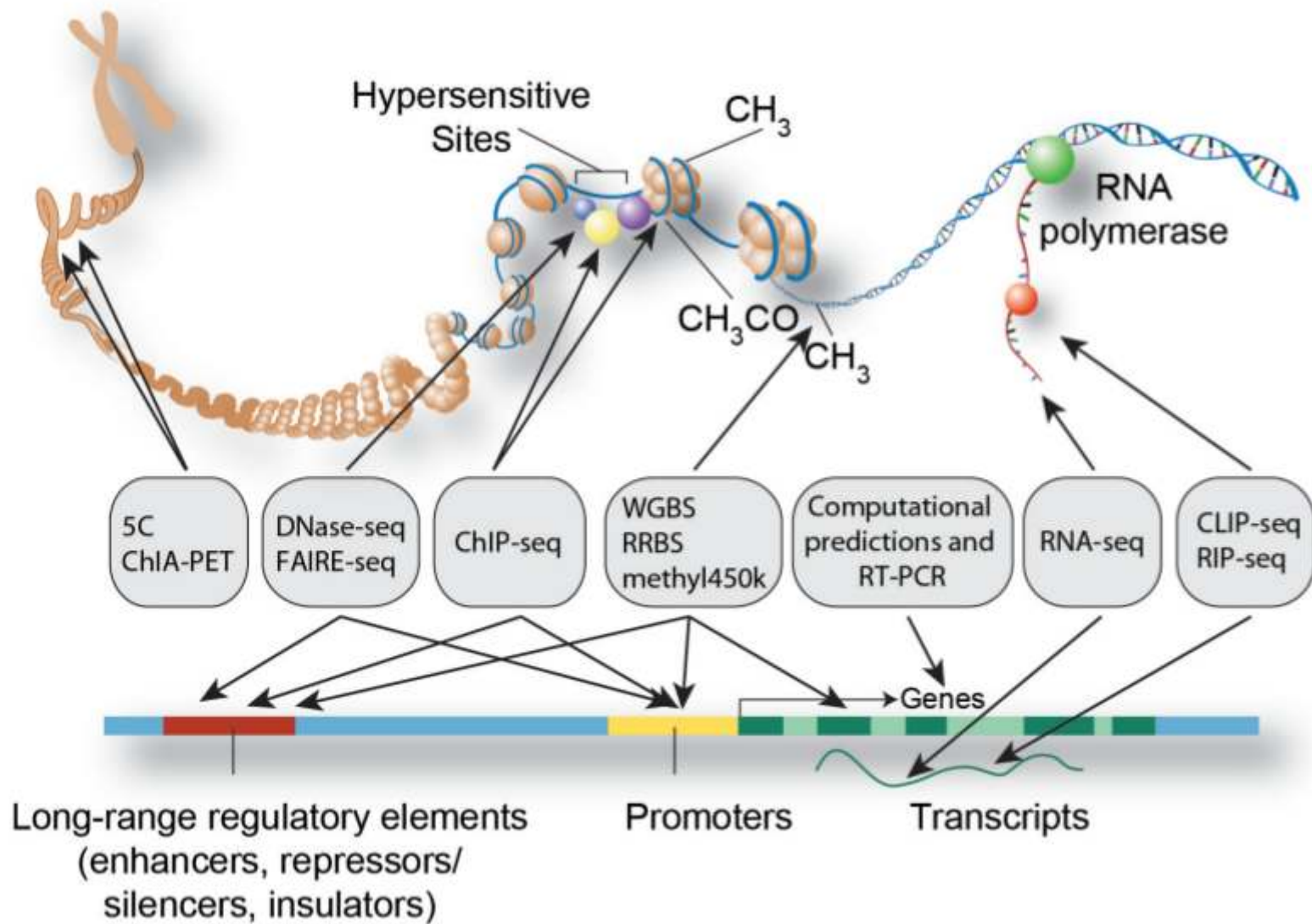
# The ENCODE project

- ◆ Encyclopedia of DNA elements
- ◆ goal: characterize (annotate) all functional elements of the human genome
- ◆ mainly USA, UK, Spain, Singapore, Japan
- ◆ 32 institutes
- ◆ 440 scientists
- ◆ \$300M budget from NHGRI (pilot phase + production, 2003-2012)



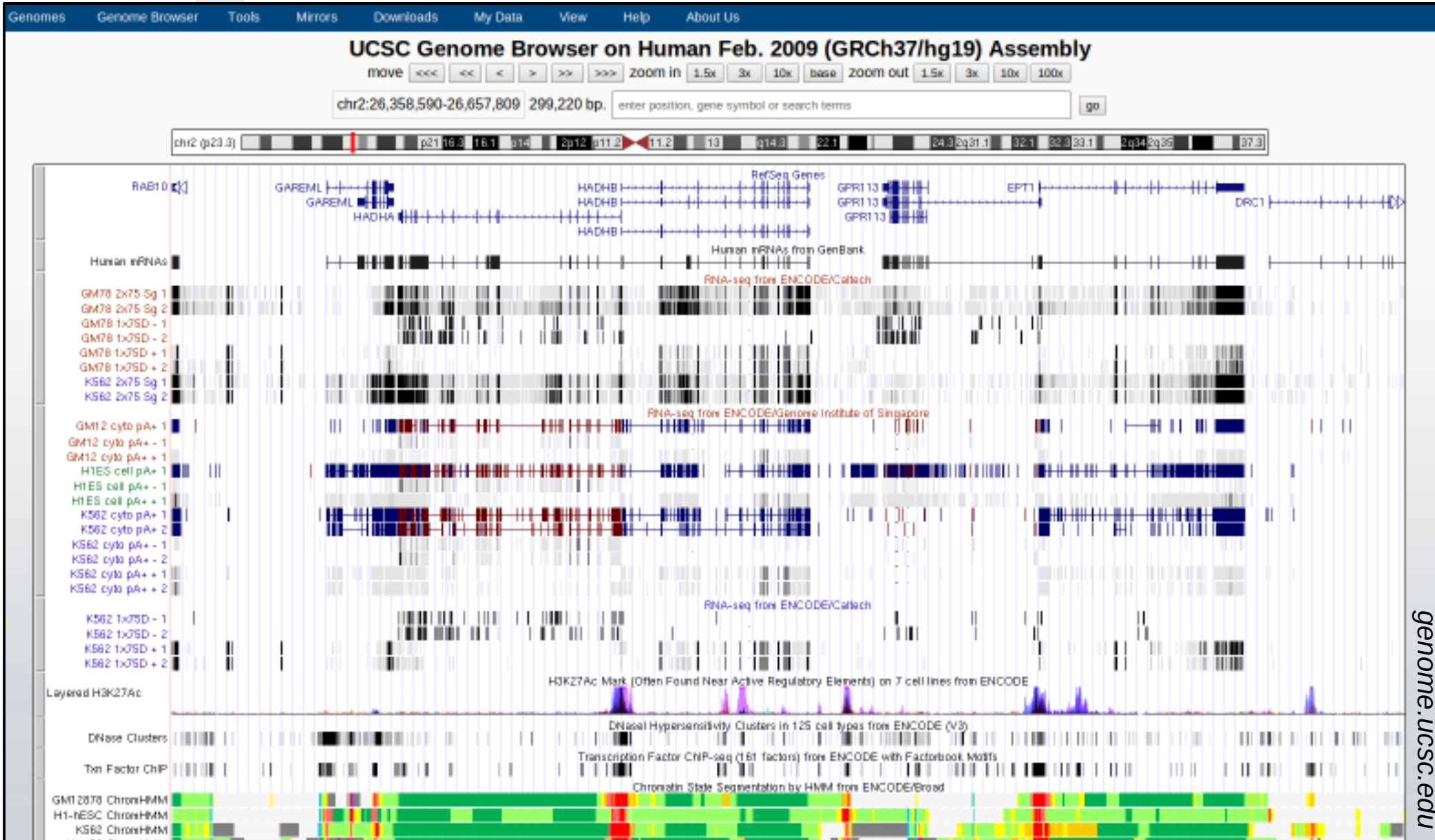


# ENCODE experiments

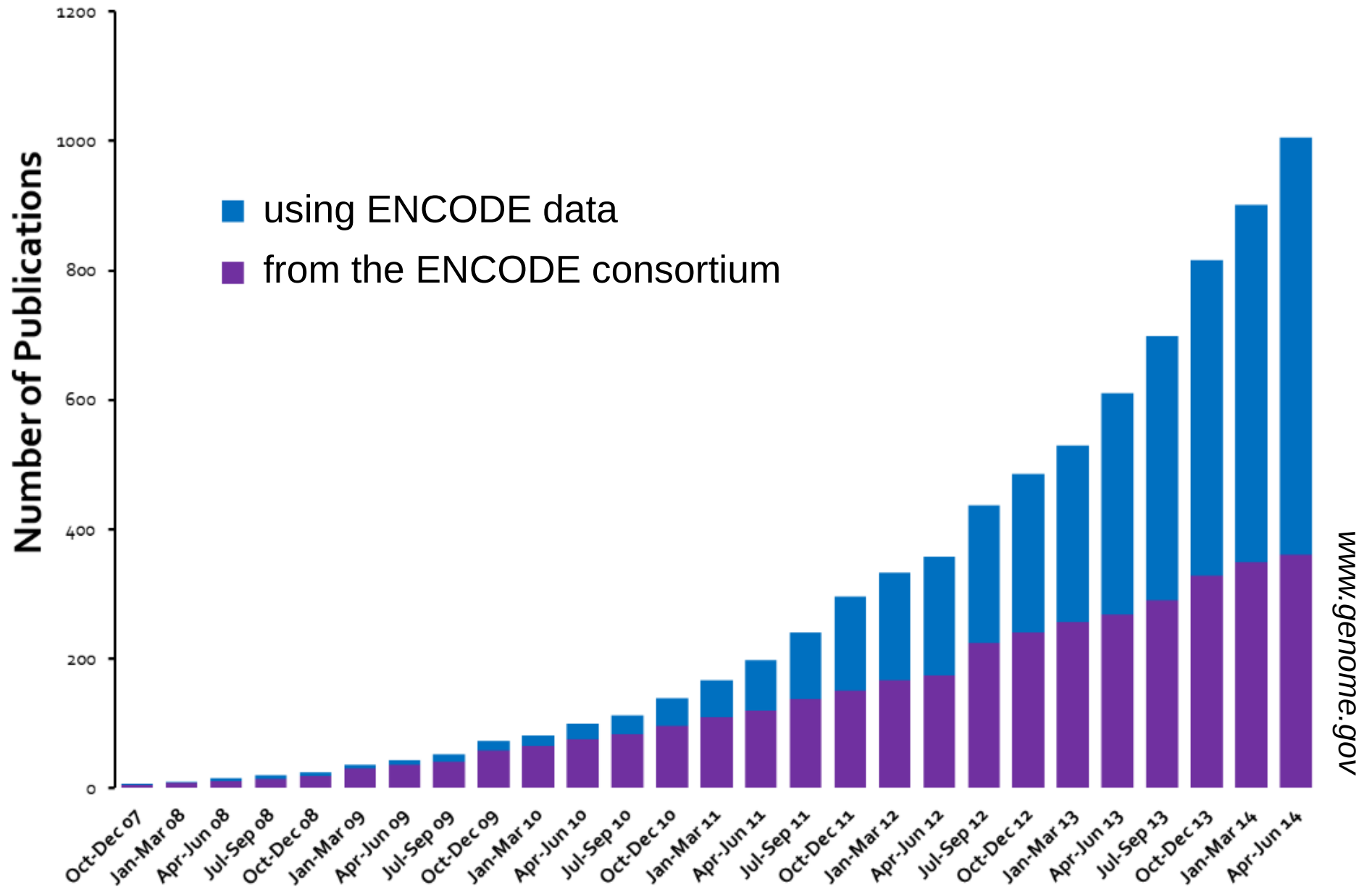


Modified from PLoS Biol 9:e1001046, 2011 & Science 306:636, 2004  
Image credits: Darryl Leja (NHGRI), Ian Dunham (EBI), Michael Pazin (NHGRI)

# ENCODE data



# ENCODE publications

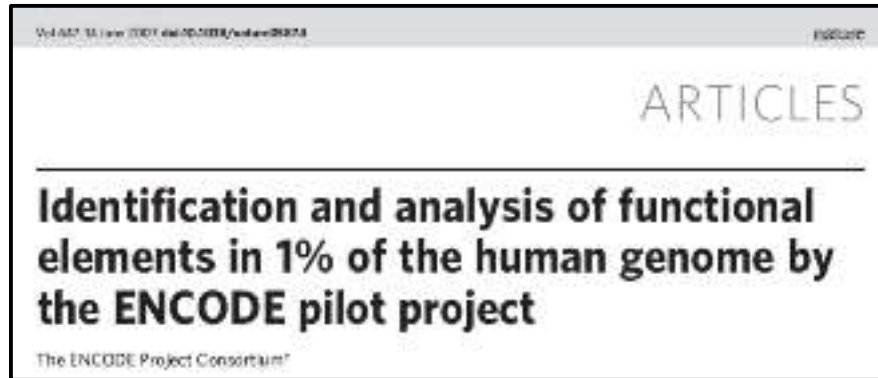


www.genome.gov

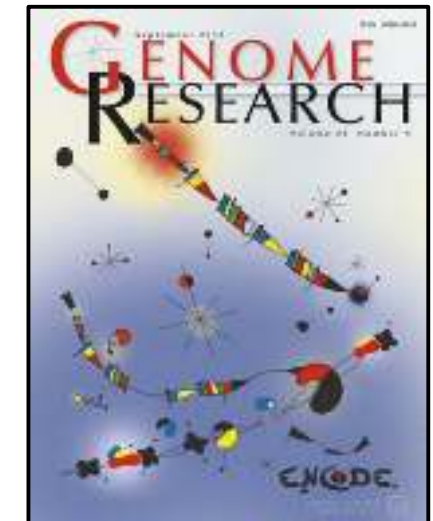
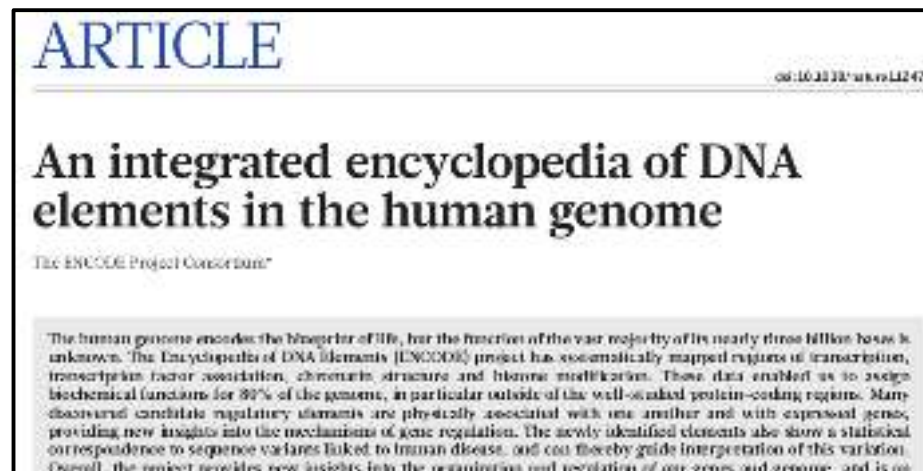


# ENCODE publications

## 2007: pilot project



## 2012: 30 articles



HEALTH

# 'Junk DNA' Debunked

Studies Find Human Genomic Makeup Is Vastly Messier; New

By GAUTAM NAIK and ROBERT LEE HOTZ

Updated Sept. 5, 2012 2:01 p.m. ET

# NEWS HEALTH

Home UK Africa Asia Australia Europe Latin America Mid-East US & Canada Business Health

## Human genome 'more active than

Human genomics

# The new world of DNA

A long-term effort to catalogue all the bits of the human genome that do

— Back to Original Article

## ENCODE project sheds light on

New findings from the ENCODE project provide insights into how DNA is used in personalized medicine.

September 05, 2012 | By Rosie Mestel and Eryn Brown

WORLD U.S. N.Y. / REGION BUSINESS TECHNOLOGY SCIENCE HEALTH

ENVIRONMENT

## Bits of Mystery DNA, Far From 'Junk,' Play

Crucial Role in Regulating Genes, Study Finds

Genetics

## Breakthrough study overturns theory of 'junk DNA' in genome

The international Encode project has found that about a fifth of the human genome regulates the 2% that makes proteins

Science News | Space | Night Sky | Roger Highfield | Dinosaurs | Evolution | Steve Jones | Science

Home Video News World

Politics Investigations

Science News | Space | Night Sky | Roger Highfield | Dinosaurs | Evolution | Steve Jones | Science

HOME » SCIENCE » SCIENCE NEWS

## Worldwide army of scientists cracks the 'junk DNA' code

One of the biggest mysteries in genetics has been solved after an international team of hundreds of scientists uncovered the secrets of "junk DNA".



## DNA project interprets 'book of life'

By Elizabeth Landau . CNN September 5, 2012 -- Updated 1841 GMT (0241 HKT)



(CNN) -- Our genes play a major role in determining who we are, but a lot of information about them has been mysterious.

That's why an international team of scientists has worked out what the working parts of the genome

they mean for the human body as we know it.

The project is called the Encyclopedia of DNA Elements (ENCODE)



# The FAANG project

FAANG: Functional Analysis of ANimal Genomes

Andersson et al, Genome Biology, 2015

*Coordinated international action to accelerate Genome to Phenome - the Functional Annotation of Animal Genome (FAANG) project*



wikipedia.org



# The FAANG project

FAANG: Functional Analysis of ANimal Genomes

[www.faang.org](http://www.faang.org)

=> 3 pilot projects

Home | Project description | Publications | CO-FAANG | more | Wiki | Login

**FAANG**  
Functional Annotation of Animal Genomes

A coordinated international action to accelerate genome to phenome

FAANG aims to:

- Standardize core assays and experimental protocols
- Coordinate and facilitate data sharing
- Establish an infrastructure for analysis of these data
- Provide high quality functional annotation of animal genomes

★ Sign up here to take part in its activities  
★ Contact respective working committees to get involved

**Working groups**

- Steering Committee
- Animals, Samples and Assays (ASA)
- Bioinformatics and Data Analysis (BKDA)
- Communication (COM)
- Metadata and Data Sharing (M&DS)

White paper authors | FAANG Contributors | FAANG Signatories

# FR-AgENCODE: overview

- ◆ a French pilot project of FAANG
- ◆ goal: improve functional annotation of livestock genomes
- ◆ founding: INRA, France (300KE)
- ◆ 4 INRA sites, 9 labs, 58 scientists
- ◆ 4 species: pig, chicken, cattle, goat
- ◆ primary targets: liver & blood cells (CD4+ & CD8+)
- ◆ molecular assays: RNA-seq, Hi-C & ATAC-seq
- ◆ duration: 2015-2017



*E. Giuffra,  
INRA GABI*

# FR-AgENCODE: overview

## Sampling: 40+ tissues

(liver, CD4+, CD8+, sperm, plasma, heart, lung, skin, fat, duodenum, ileum, jejunum, cerebellum, frontal lobe, olfactory bulb, trigeminal ganglia, hypothalamus, pancreas, adrenals, kidney, muscle, bone, joints, spleen, lymphatic nodes, peyer's patches, ovary, oocytes, oviduct, uterus, mammary gland, acini, testis, seminal vesicle, etc)

2x ♂  
2x ♀



*Sus scrofa*  
(Large White)



*Gallus gallus*  
(White Leghorn)



*Bos Taurus*  
(Holstein)

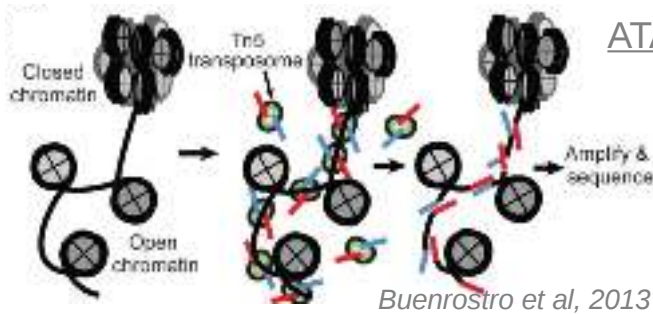


*Capra hircus*  
(Alpine)

=> **INRA CRB-Anim biorepository**

ATAC-seq: chromatin accessibility

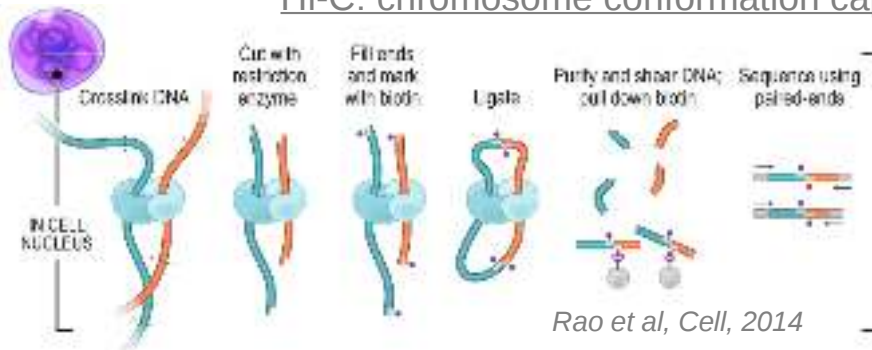
RNA-seq: long & short RNAs



## Molecular assays: 3 target tissues

transcriptome & chromatin structure profiling  
polyA+ RNA-seq (mRNAs & lncRNAs, 130M RP/lib)  
small RNA-seq (miRNAs & <200nt RNAs, 40MR/lib)  
Hi-C (130M RP/lib) & ATAC-seq (40M RP/lib)

Hi-C: chromosome conformation capture



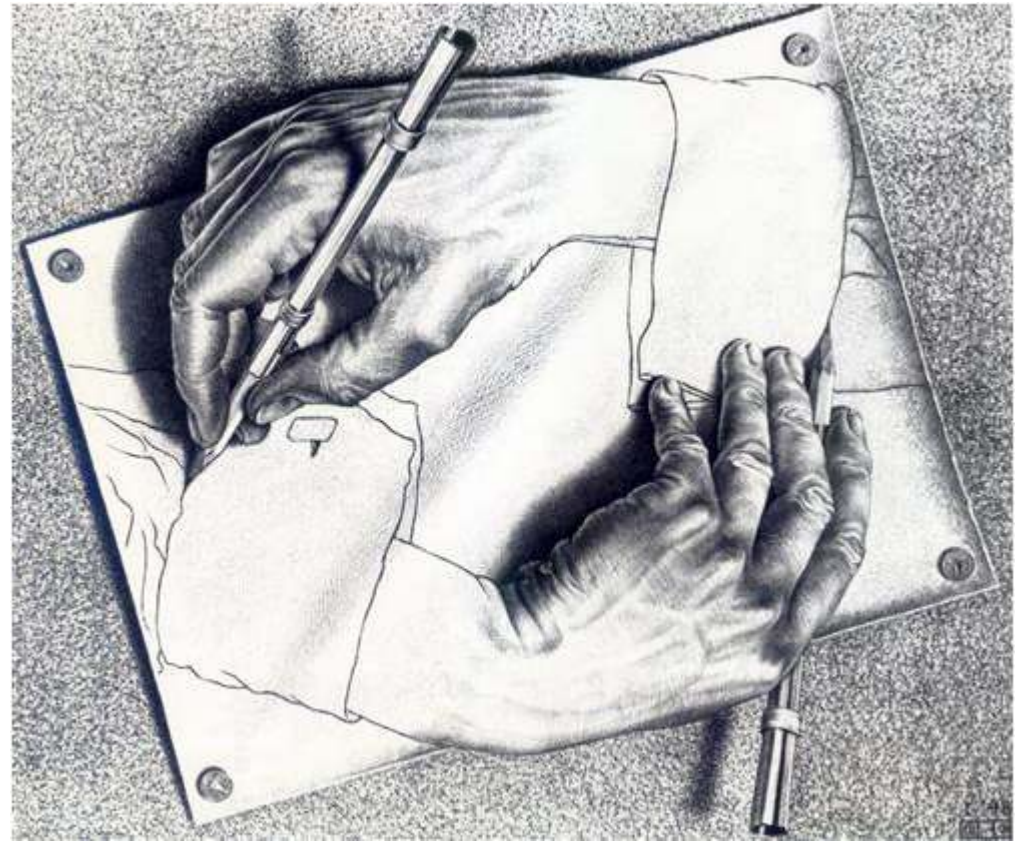
## Bioinformatics data analysis

genome annotation, gene expression, lncRNAs & sRNAs annotation/prediction, chromosome interaction matrices & contact heatmaps, allele-specific expression, chimeric transcripts detection, comparative genomics, etc



# Outline

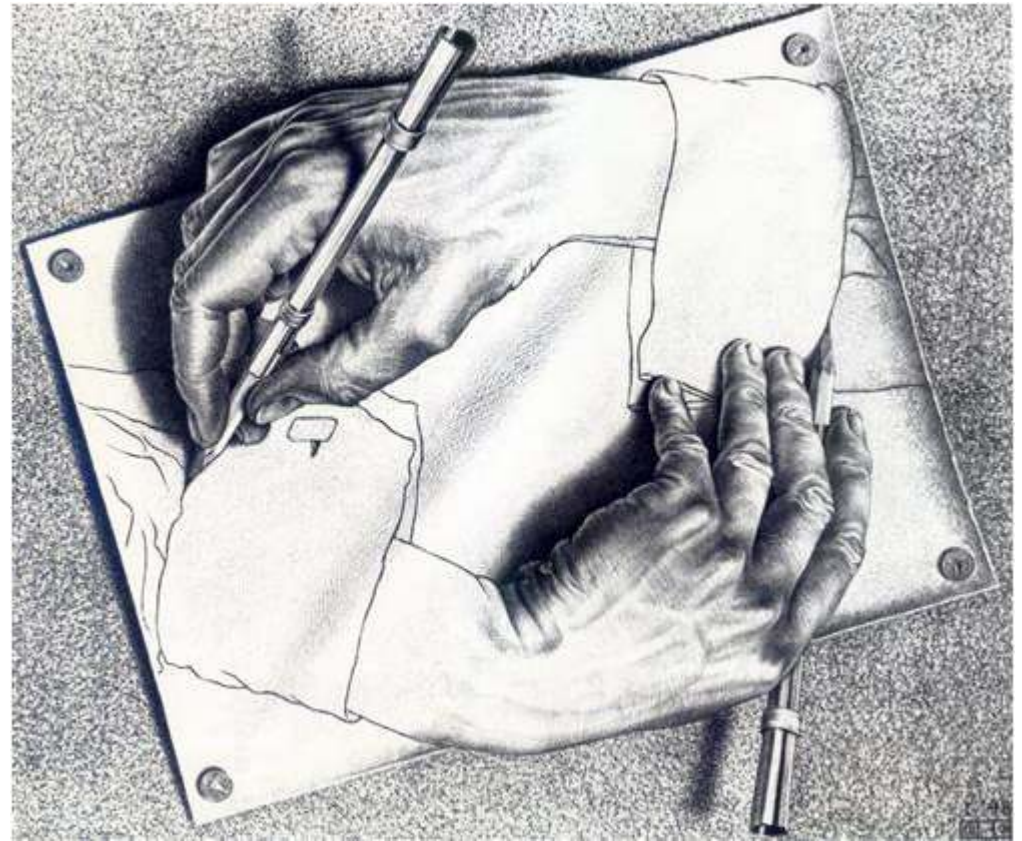
- ◆ Why
- ◆ What
- ◆ How ←
- ◆ Where
- ◆ Who



*M.C. Escher, 1948*

# Outline

- ◆ Why
- ◆ What
- ◆ How
- ◆ Where ←
- ◆ Who



*M.C. Escher, 1948*

# FR-AgENCODE: sampling

2x ♂

2x ♀



*Sus scrofa*  
(Large White)



*Gallus gallus*  
(White Leghorn)



*Bos Taurus*  
(Holstein)



*Capra hircus*  
(Alpine)

- ◆ 34 somatic tissues + 13 reproductive tissues (8 female + 5 male):

liver, sperm, CD4+, CD8+, plasma, heart, lung, skin, fat, duodenum, ileum, jejunum, cerebellum, frontal lobe, olfactory bulb, trigeminal ganglia, hypothalamus, pancreas, adrenals, kidney, muscle, bone, joints, spleen, lymphatic nodes, peyer's patches, ovary, oocytes, oviduct, uterus, mammary gland, acini, testis, seminal vesicle, etc...

- ◆ total: 2,000 to 6,000 samples



M. Tixier-  
Boichard,  
INRA GABI

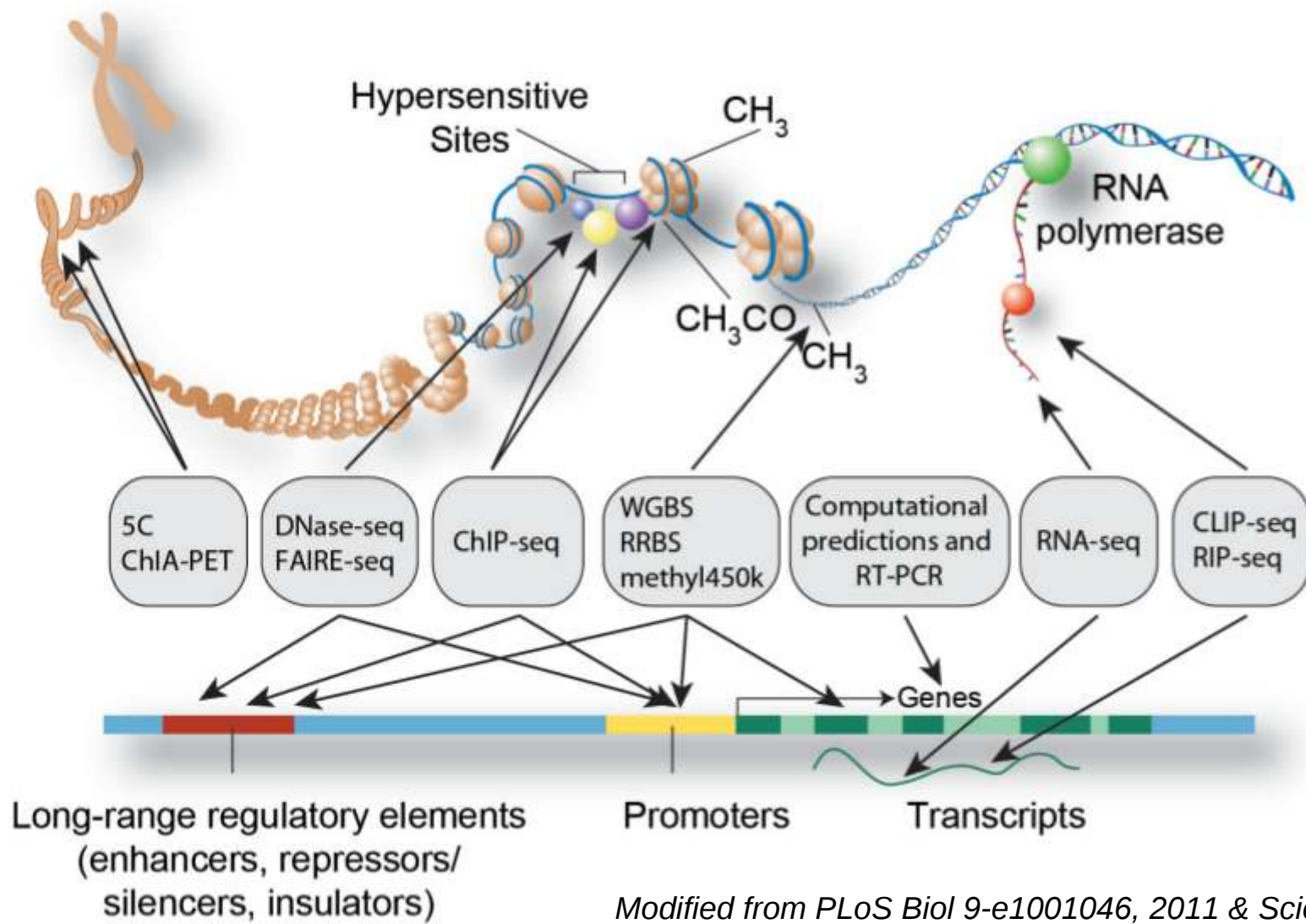
S. Fabre  
INRA  
GenPhySE



**INRA CRB-Anim biorepository**

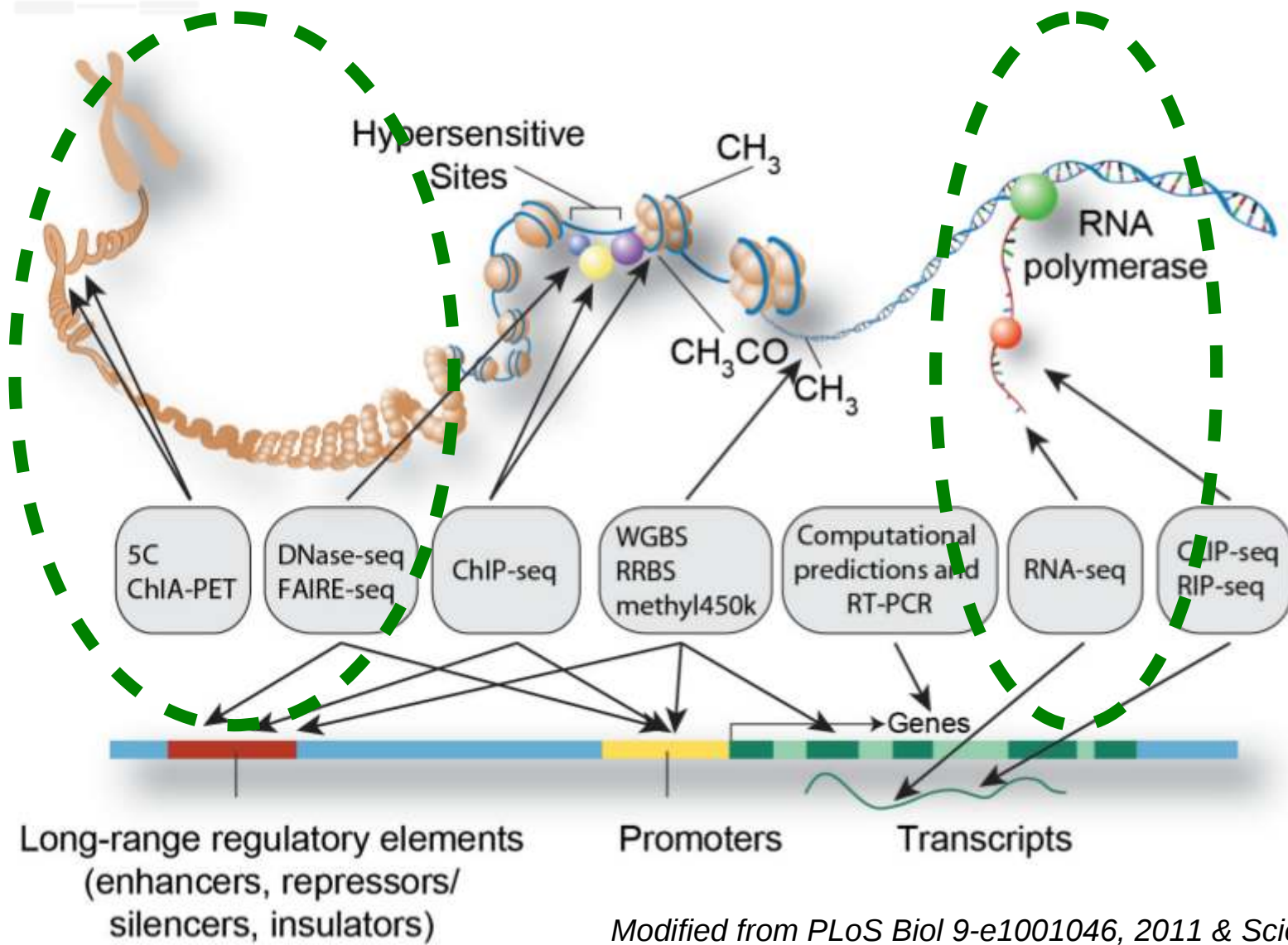


# FR-AgENCODE: molecular assays



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

# FR-AgENCODE: molecular assays



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

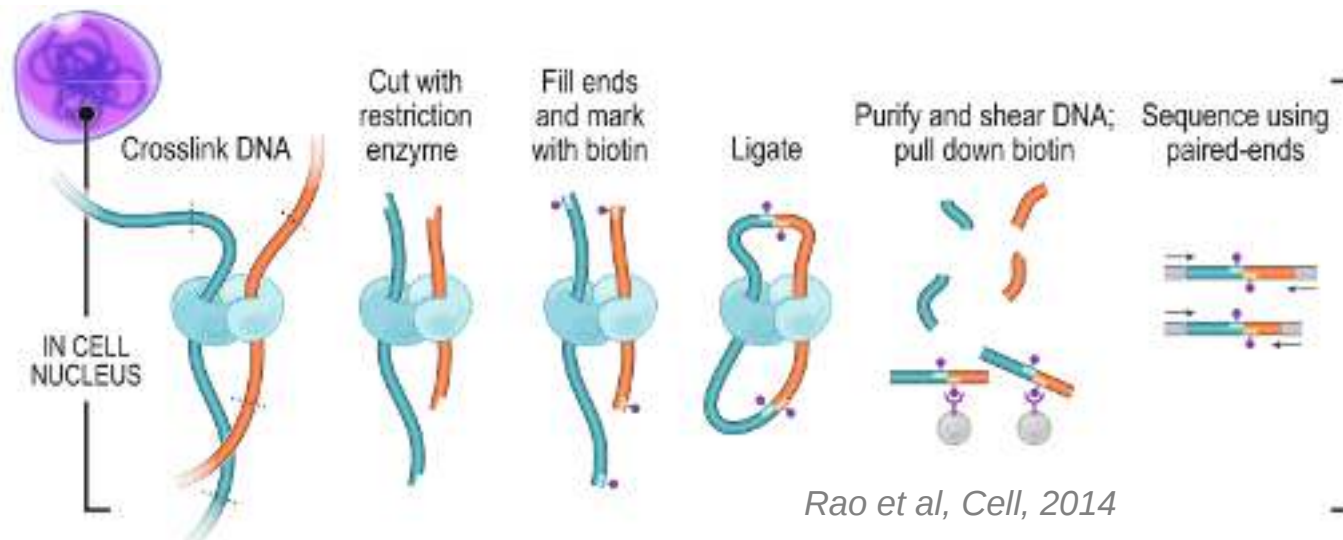
# FR-AgENCODÉ: molecular assays

- ◆ RNA-seq
  - directional protocol, Illumina Hi-Seq3000
  - ◆ polyA+ RNAs (mRNAs + lncRNAs)  
100M+ pairs (2x150bp) per sample (2/lane)
  - ◆ sRNAs (miRNAs and others)  
35M reads (1x50bp) per sample (3/lane)



# FR-AgENCODE: molecular assays

- ◆ RNA-seq  
directional protocol, Illumina Hi-Seq3000
  - ◆ polyA+ RNAs (mRNAs + lncRNAs)  
100M+ pairs (2x150bp) per sample (2/lane)
  - ◆ sRNAs (miRNAs and others)  
35M reads (1x50bp) per sample (3/lane)
- ◆ Hi-C: chromosome conformation capture  
(Lieberman-Aiden et al, Science, 2009, Rao et al, Cell, 2014)

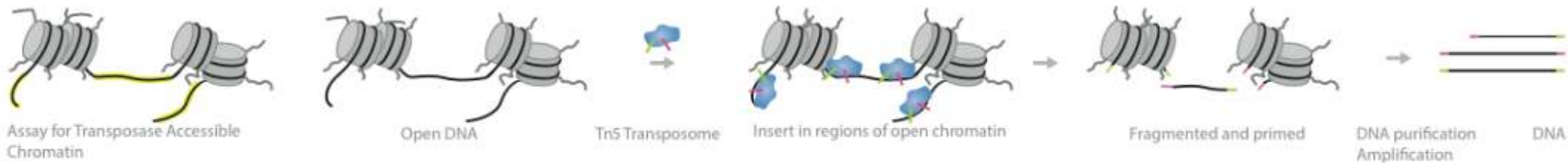


# FR-AgENCODE: molecular assays

- ◆ RNA-seq
  - directional protocol, Illumina Hi-Seq3000
    - ◆ polyA+ RNAs (mRNAs + lncRNAs)  
100M+ pairs (2x150bp) per sample (2/lane)
    - ◆ sRNAs (miRNAs and others)  
35M reads (1x50bp) per sample (3/lane)
- ◆ Hi-C: chromosome conformation capture
  - ◆ 16 samples (4 replicates, 4 species)
  - ◆ 130M read pairs per sample (2 lib/lane)

# FR-AgENCODE: molecular assays

- ◆ RNA-seq
  - directional protocol, Illumina Hi-Seq3000
    - ◆ polyA+ RNAs (mRNAs + lncRNAs)  
100M+ pairs (2x150bp) per sample (2/lane)
    - ◆ sRNAs (miRNAs and others)  
35M reads (1x50bp) per sample (3/lane)
- ◆ Hi-C: chromosome conformation capture
  - ◆ 16 samples (4 replicates, 4 species)
  - ◆ 130M read pairs per sample (2 lib/lane)
- ◆ ATAC-seq: chromatin accessibility



*Adapted from: [www.illumina.com/techniques](http://www.illumina.com/techniques)*



# FR-AgENCODÉ: molecular assays

- ◆ RNA-seq  
directional protocol, Illumina Hi-Seq3000
  - ◆ polyA+ RNAs (mRNAs + lncRNAs)  
100M+ pairs (2x150bp) per sample (2/lane)
  - ◆ sRNAs (miRNAs and others)  
35M reads (1x50bp) per sample (3/lane)
- ◆ Hi-C: chromosome conformation capture
  - ◆ 16 samples (4 replicates, 4 species)
  - ◆ 130M read pairs per sample (2 lib/lane)
- ◆ ATAC-seq: chromatin accessibility



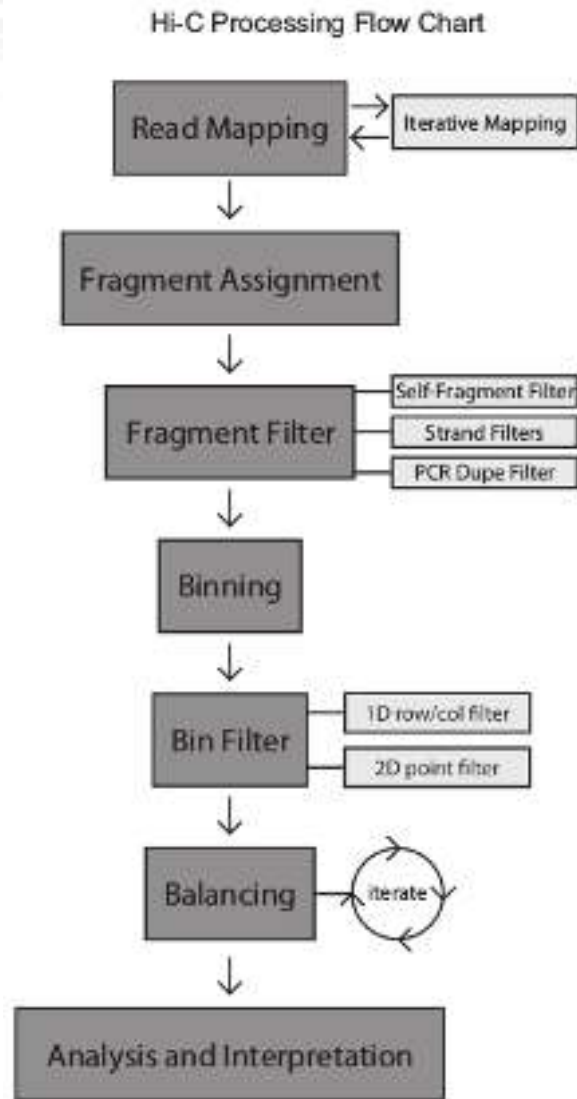
*D. Esquerré  
INRA GetPlage*



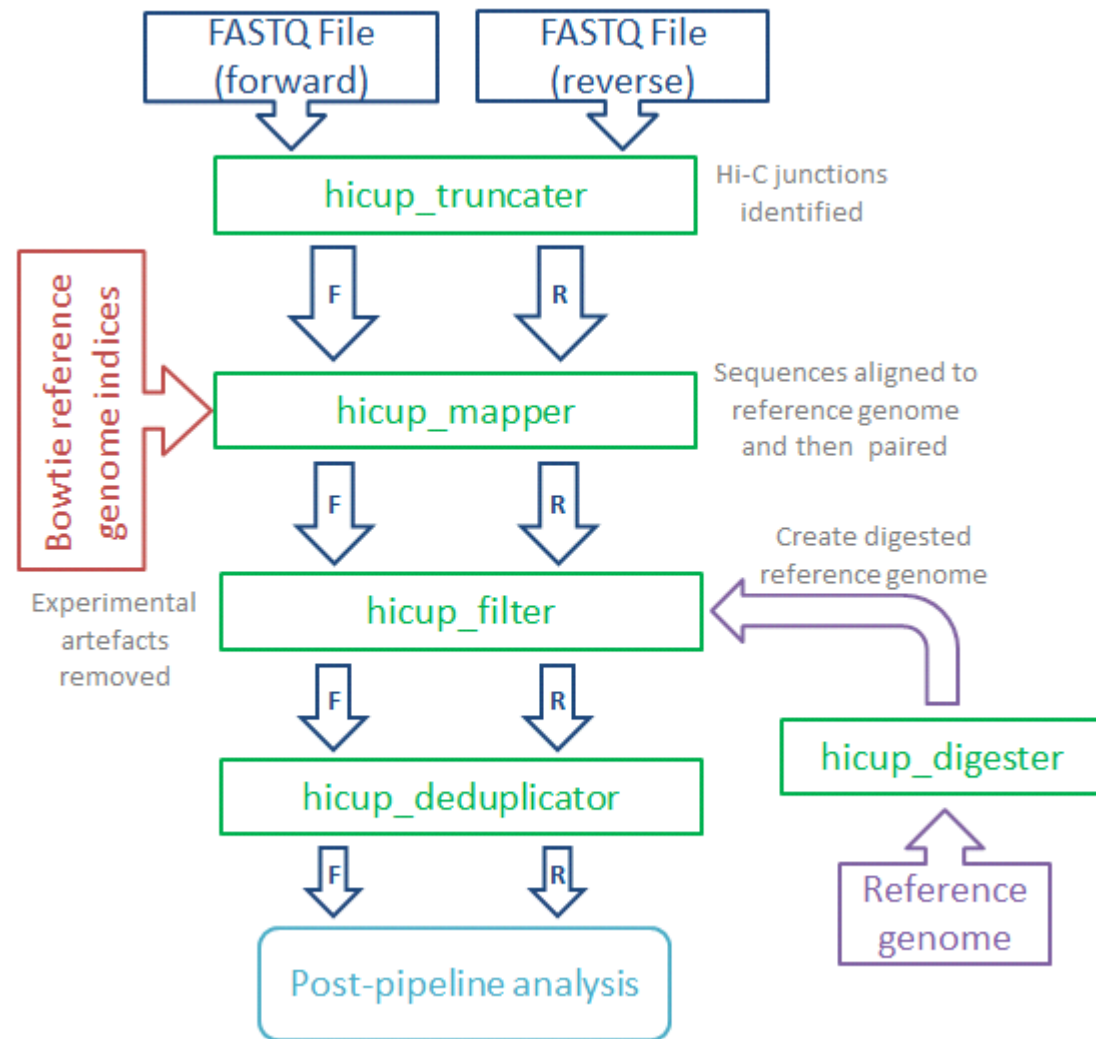
*H. Acloque,  
INRA GenPhySE*

➔ **Comparative analysis of genome topology and expression**

# Hi-C data analysis: overview



Lajoie et al, 2015



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

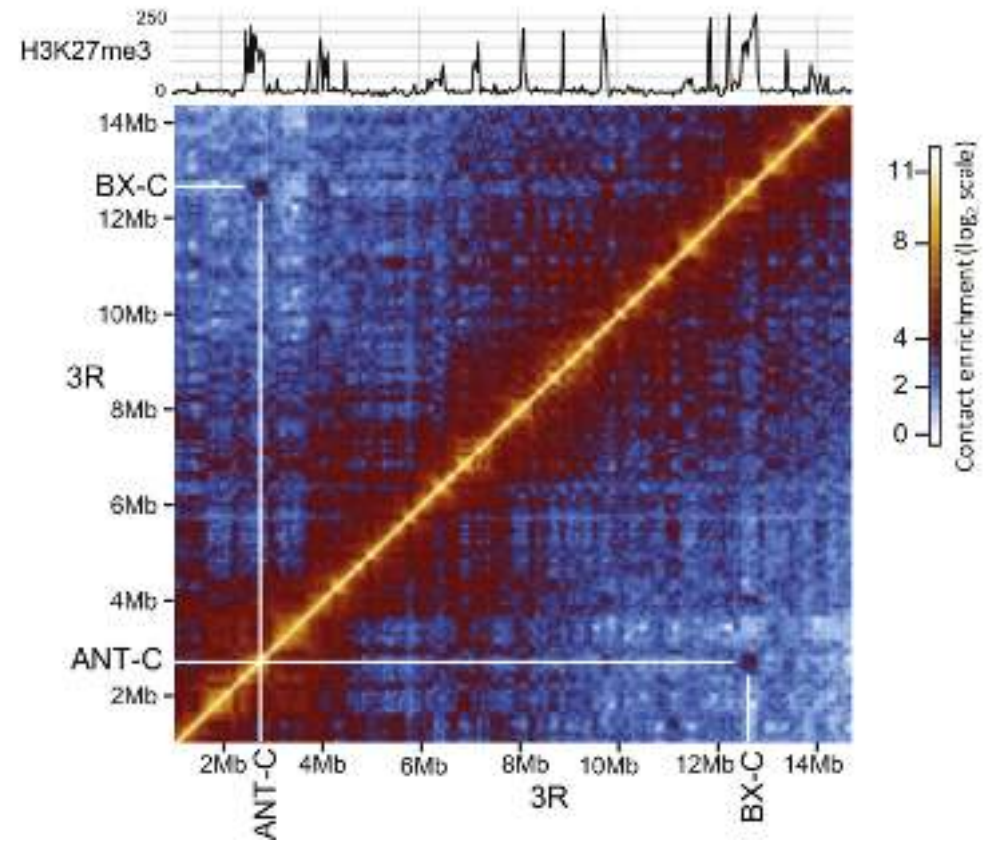
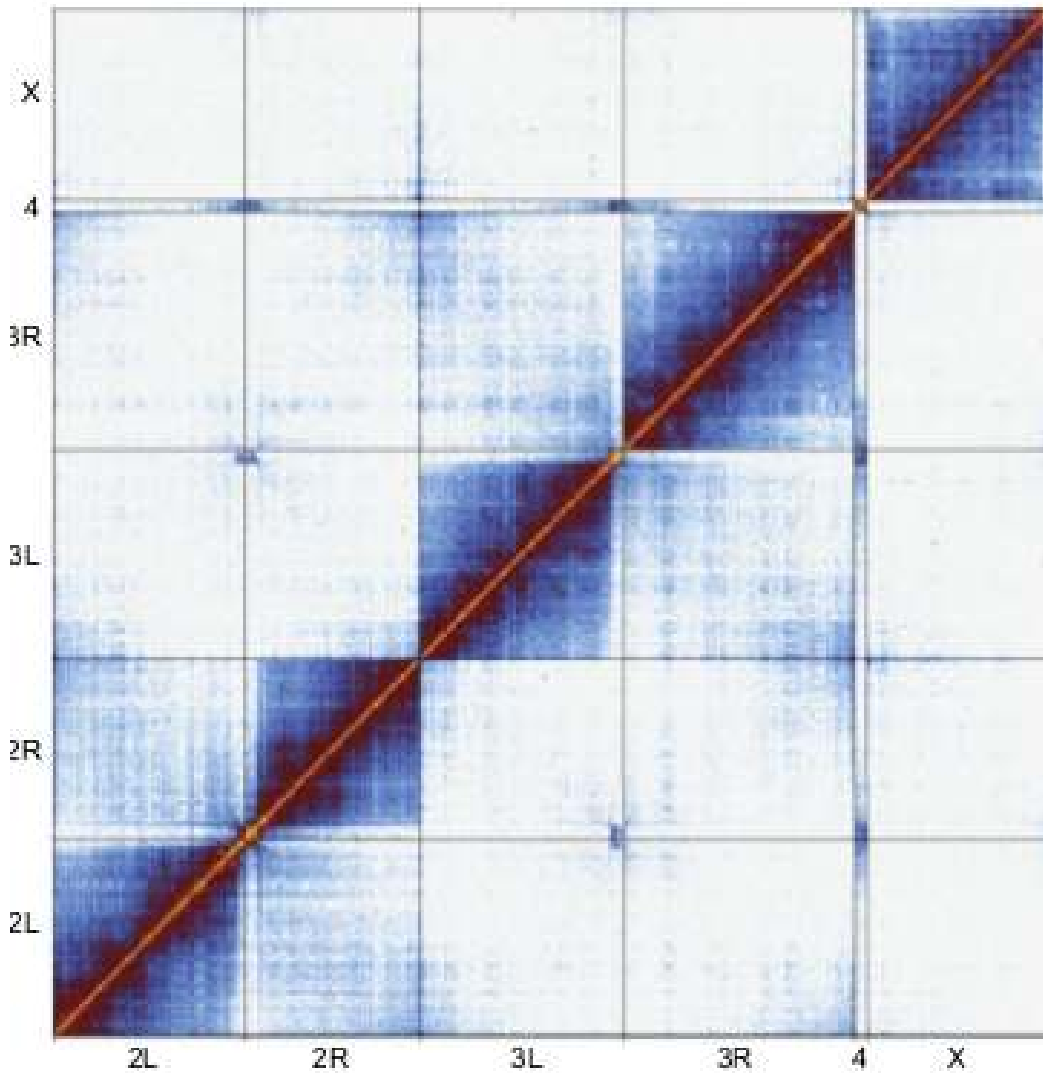


# Hi-C data analysis: overview

- ◆ 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: the contact matrix



*Sexton et al 2012*





# 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

## A FAST ALGORITHM FOR MATRIX BALANCING

PHILIP A. KNIGHT\* AND DANIEL RUIZ†

**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.

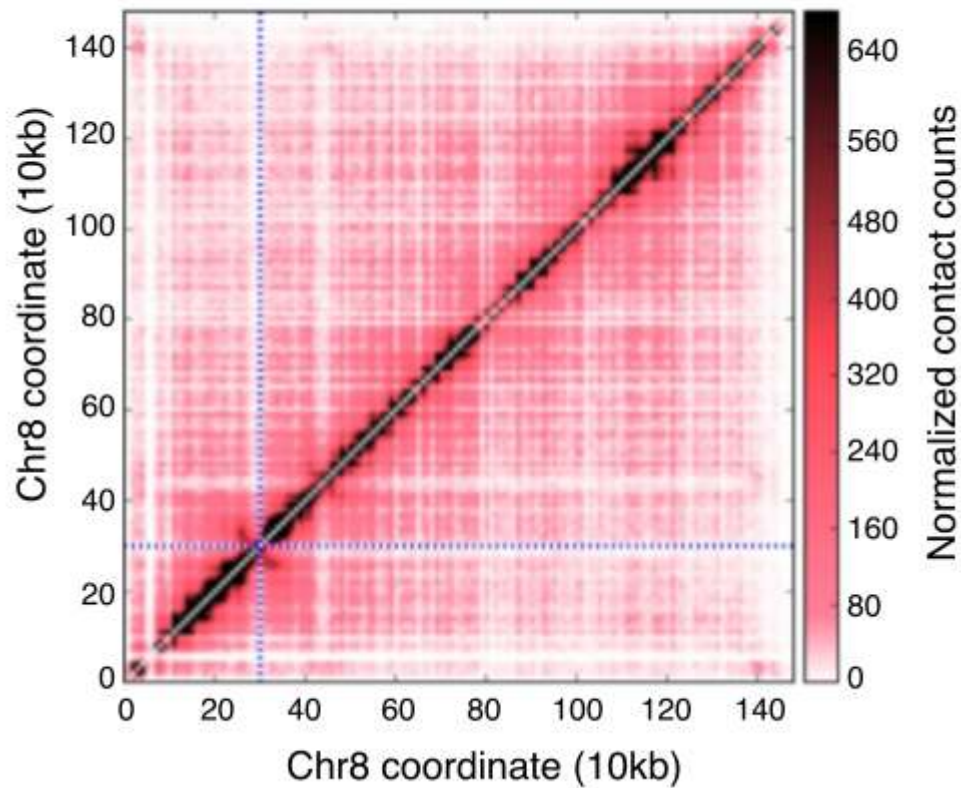
**1. Introduction.** For at least 70 years, scientists in a wide variety of disciplines have attempted to transform square nonnegative matrices into doubly stochastic form by applying diagonal scalings. That is, given  $A \in \mathbb{R}^{n \times n}$ ,  $A \geq 0$ , find diagonal matrices  $D_1$  and  $D_2$  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*

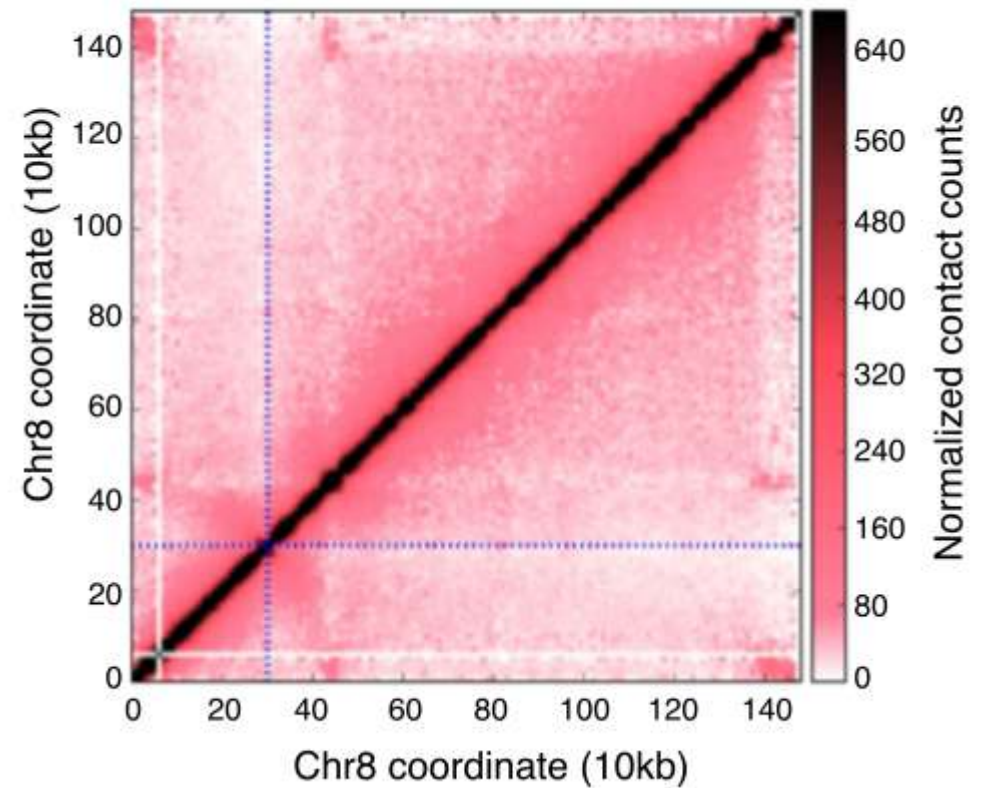


# Hi-C data analysis: matrix normalization

**(a)** Raw contact map

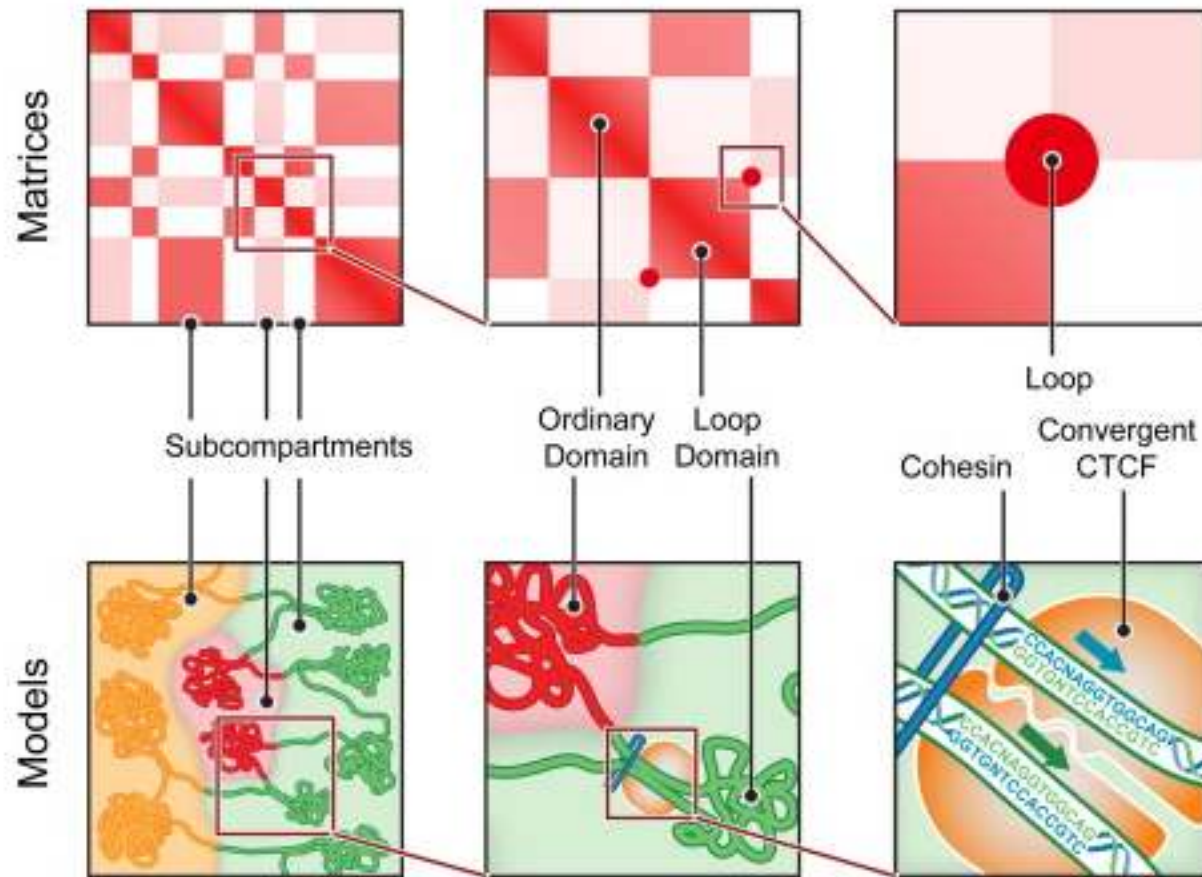


**(b)** Normalized contact map



*Ay & Noble, Genome Biology, 2015*

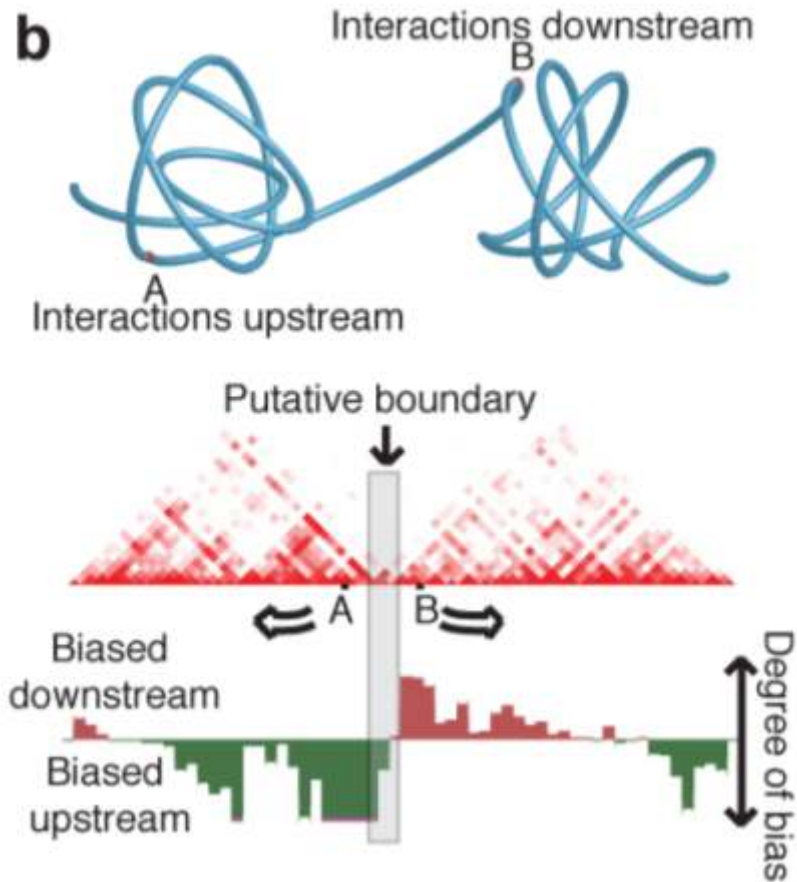
# Hi-C data analysis: finding topologically associated domains (TADs)



*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

Published in final edited form as:  
Nature. 2012;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>1,2,7</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 nuclei of eukaryotic cells, interphase chromosomes occupy distinct chromosome territories (CT).

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 Rätsch

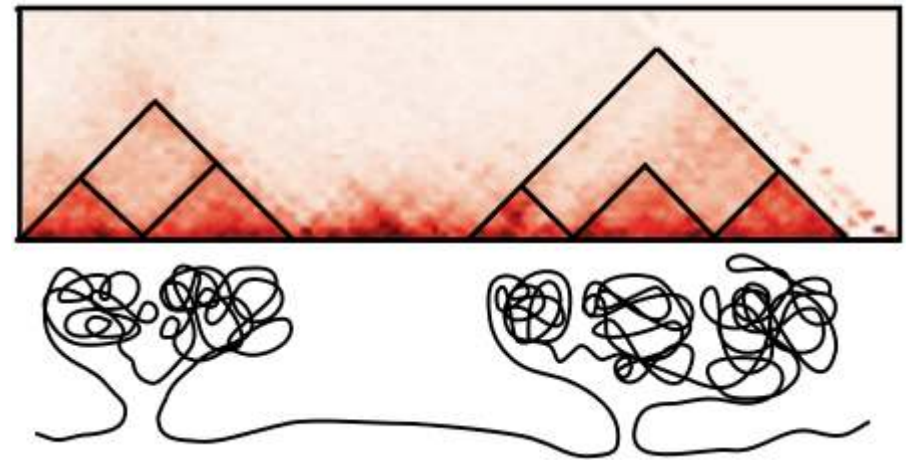
### ABSTRACT

**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 hierarchy of nested TADs. This model is based on empirical distributions of contact frequencies within TADs, where positions that are located near a greater number of contacts than

resulting in a contact matrix  $A$ , where  $A_{ij}$  is the number of contacts between bias  $i$  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 TAD statistic from the contact matrix  $A_{ij}$ ; (2) methods that exploit the 2D structure of the contact matrix.

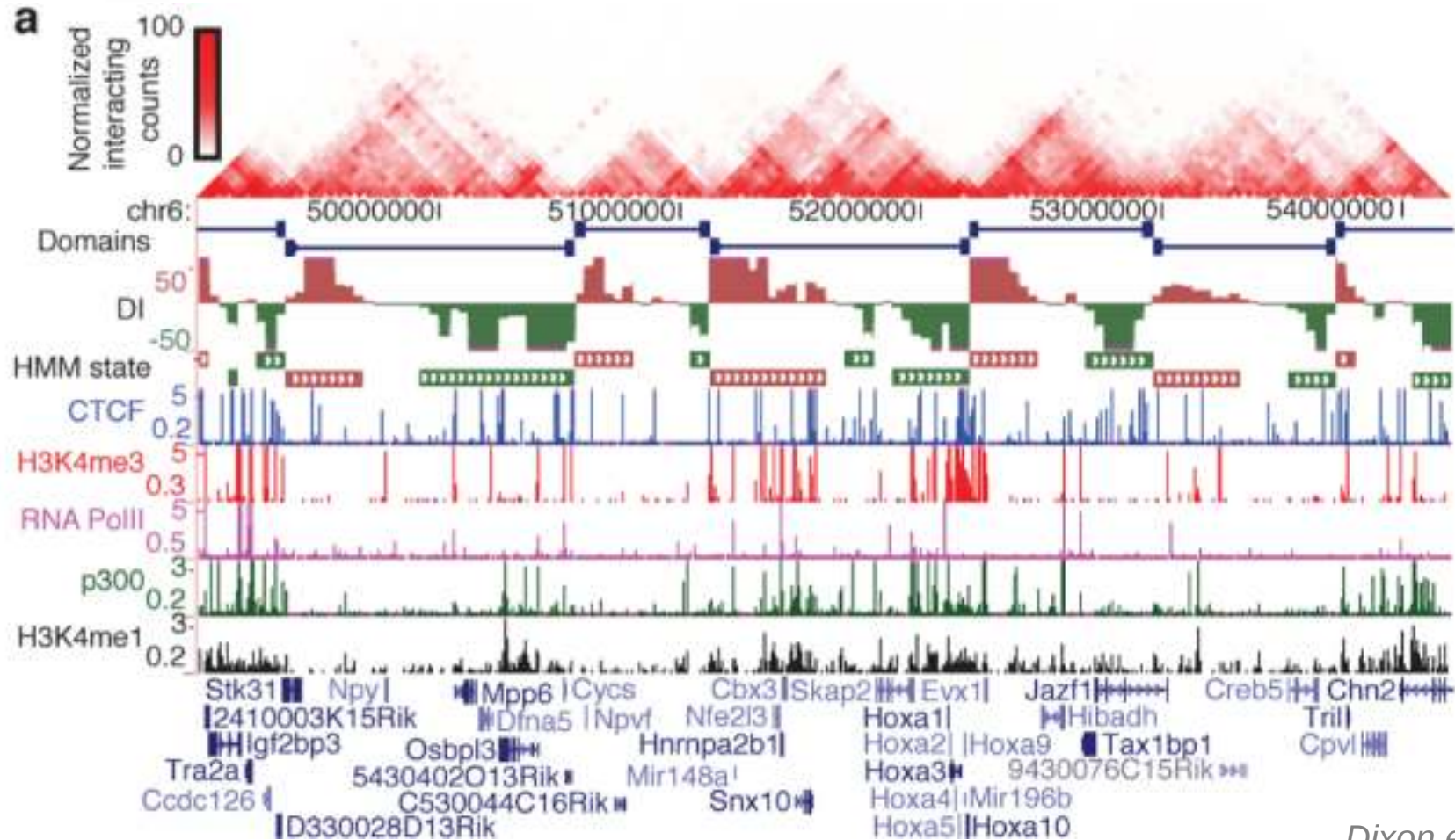
Dixon et al. (2012) compute a 1D "directionality index" (DI) from the contact matrix. This index defines whether contacts have an upstream bias, downstream bias or no bias. Next, they use a hidden Markov model (HMM) to partition the genome into regions defined by changes in the directionality index.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



Weinreb & Raphael, *Bioinformatics*, 2015



# Hi-C data analysis: integrative analysis



*Dixon et al 2012*

# Hi-C data analysis: FR-AgENCODER pipeline

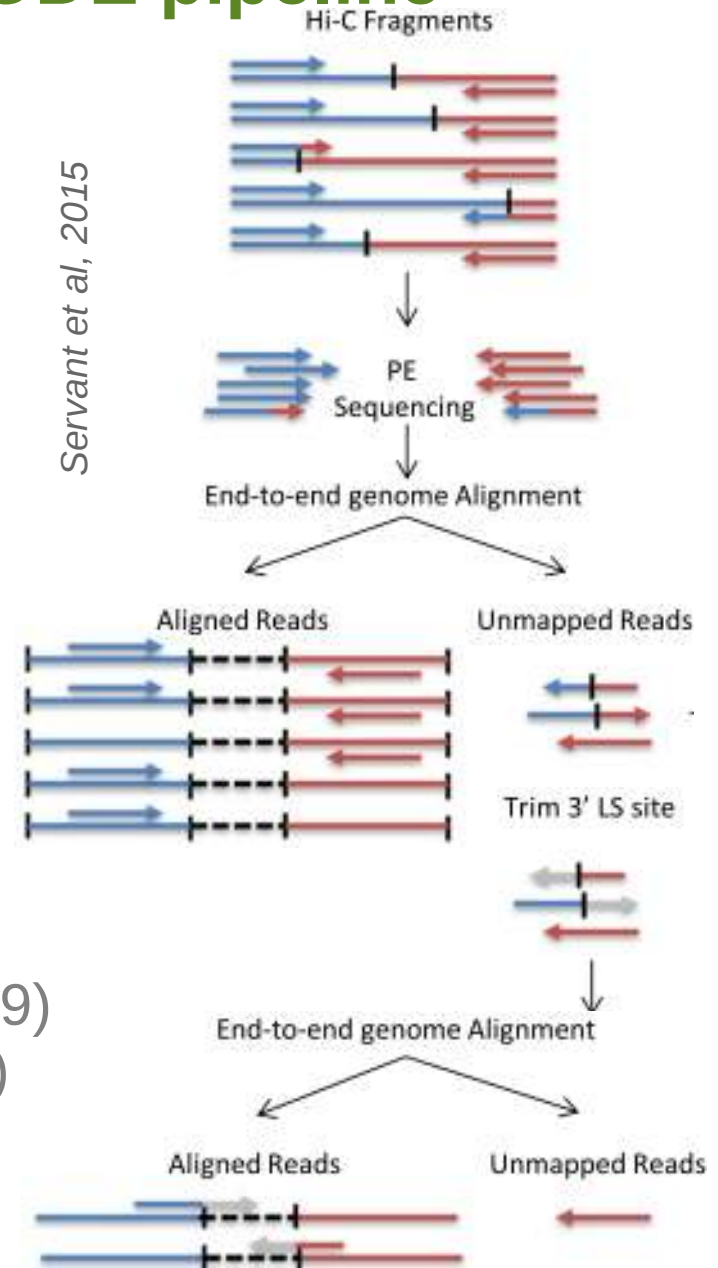
## Data analysis

### Pipeline

- ▶ Trim reads (ligation site)
- ▶ Map on reference genome
- ▶ Discard inconsistent pairs
- ▶ Count reads in pairs of genomic bins & generate contact matrix
- ▶ Normalize contact matrix (non parametric, matrix balancing)
- ▶ Identify Topologically Associated Domains, *cis* and *trans* interactions

### Software

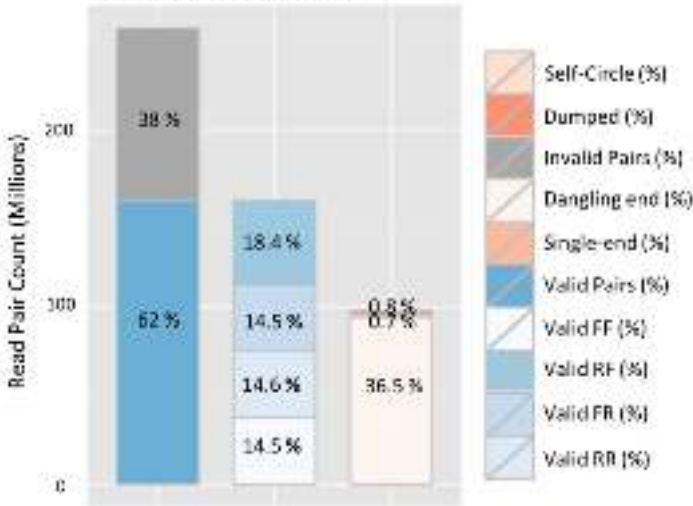
- ▶ HiC-Pro pipeline (Servant et al 2015)
- ▶ Bowtie2 mapping (Langmead et al, 2009)
- ▶ ICE normalization (Imakaev et al, 2012)
- ▶ HiTC display (Servant et al, 2012)
- ▶ HiFive pipeline (Sauria et al, 2015)



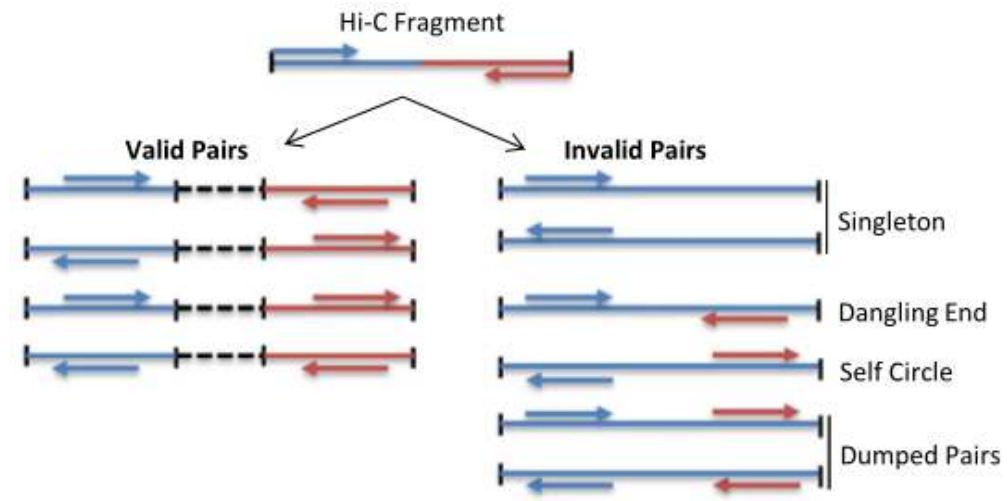
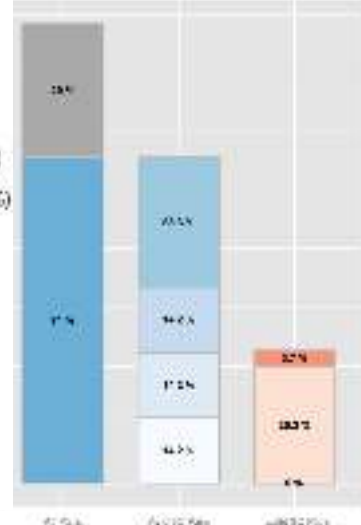
# FR-AgENCODE Hi-C preliminary results

## Read pairs status after mapping

Read Pair Filtering (IMR90)



Read Pair Filtering (CH12)



Servant et al 2015

**Dixon et al data** (human, from Servant et al 2015) **Rao et al data** (mouse, CH12 cells)

**62%**  
**valid pairs**

**71%**  
**valid pairs**

# FR-AgENCODE Hi-C preliminary results

## Read pairs status after mapping

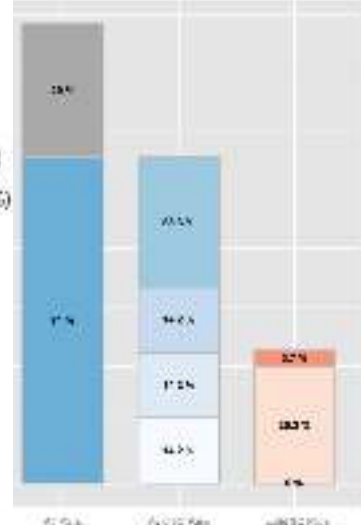
Read Pair Filtering (IMR90)



**Dixon et al data**  
(human, from Servant et al 2015)

**62%**  
**valid pairs**

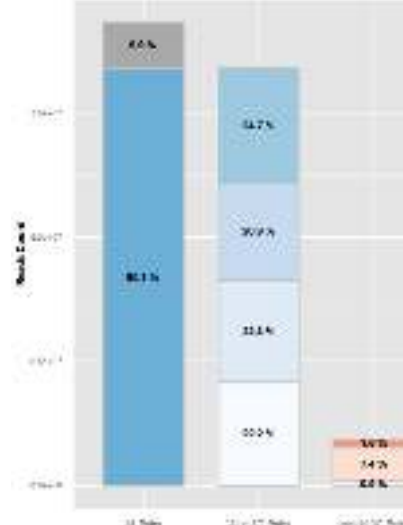
Read Pair Filtering (IMR90)



**Rao et al data**  
(mouse, CH12 cells)

**71%**  
**valid pairs**

Read Pair Filtering (IMR90)



**FR-AgENCODE data**  
(mouse, STO cells)

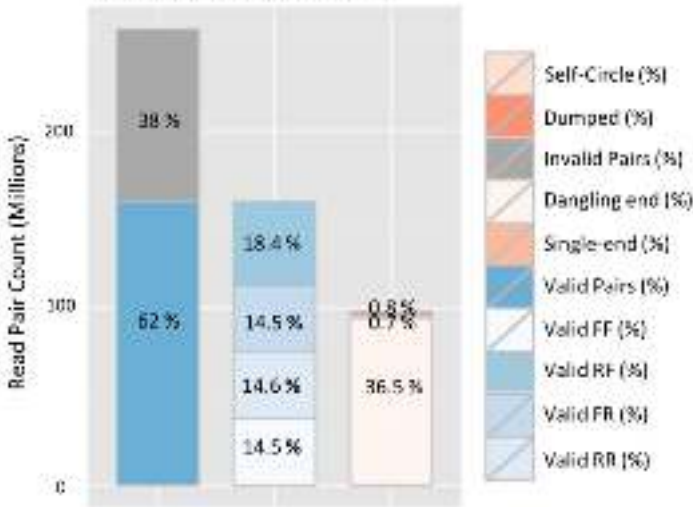
**90%**  
**valid pairs**



# FR-AgENCODE Hi-C preliminary results

## Read pairs status after mapping

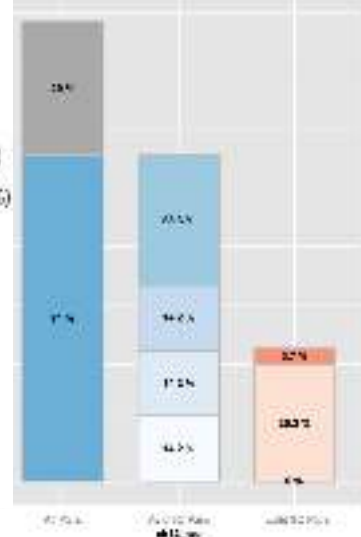
Read Pair Filtering (IMR90)



**Dixon et al data**  
(human, from Servant et al 2015)

**62%**  
**valid pairs**

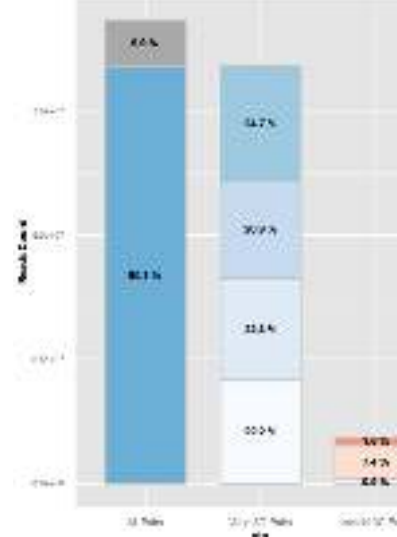
Read Pair Filtering (IMR90)



**Rao et al data**  
(mouse, CH12 cells)

**71%**  
**valid pairs**

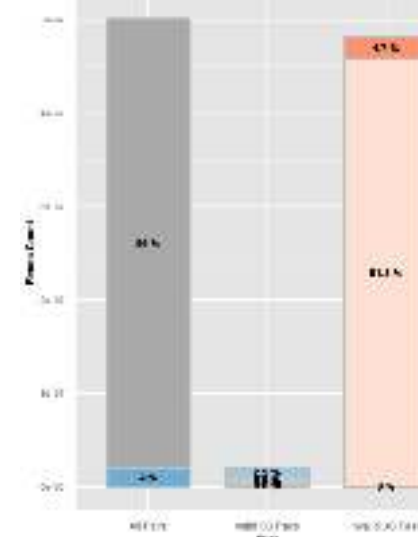
Read Pair Filtering (IMR90)



**FR-AgENCODE data**  
(mouse, STO cells)

**90%**  
**valid pairs**

Read Pair Filtering (IMR90)

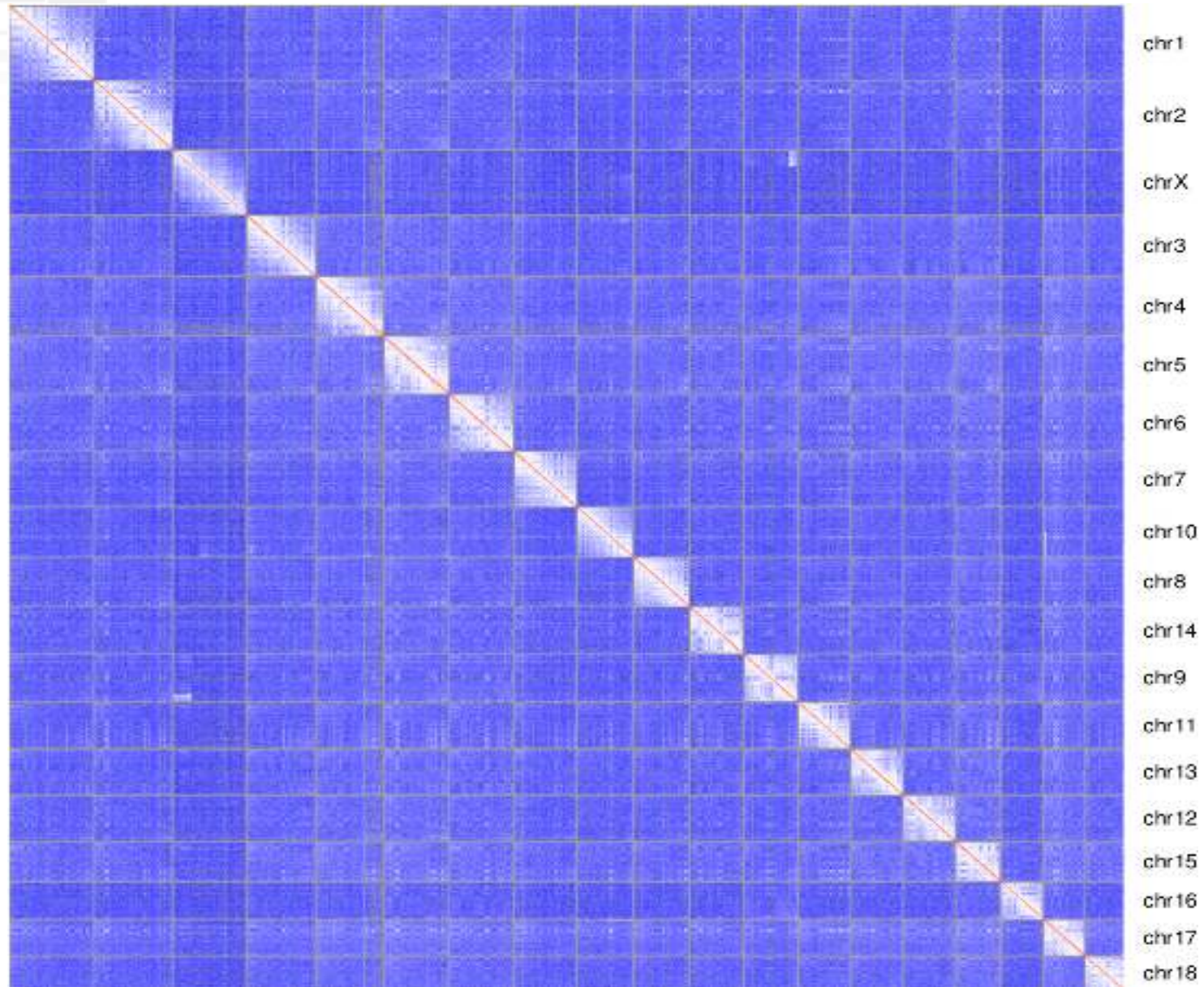


**FR-AgENCODE data**  
(pig, hepatocytes)

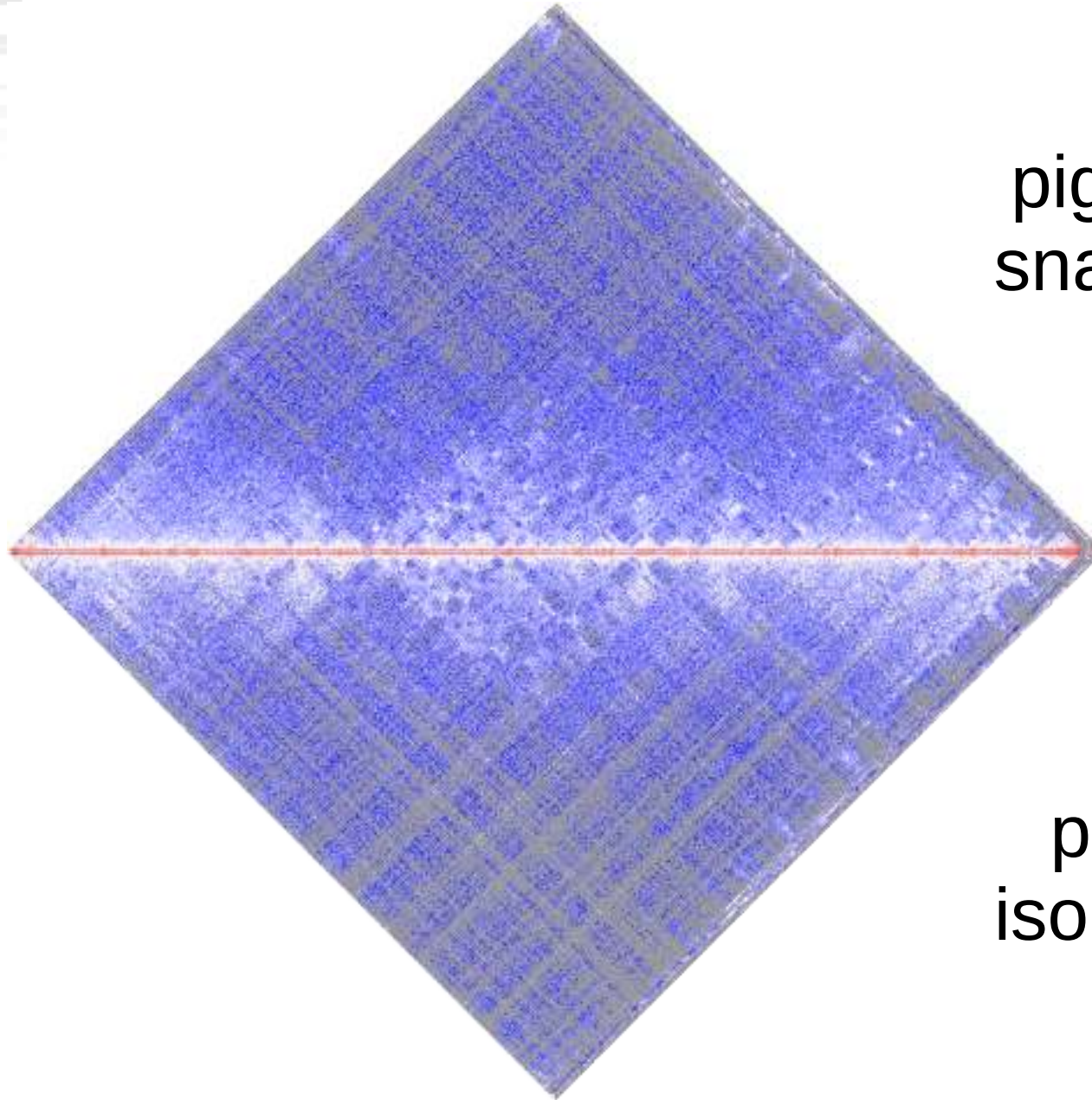
**4%**  
**valid pairs**

# FR-AgENCODER: Hi-C preliminary set-up

Mouse STO cells,  
whole genome



# FR-AgENCODE: Hi-C preliminary set-up



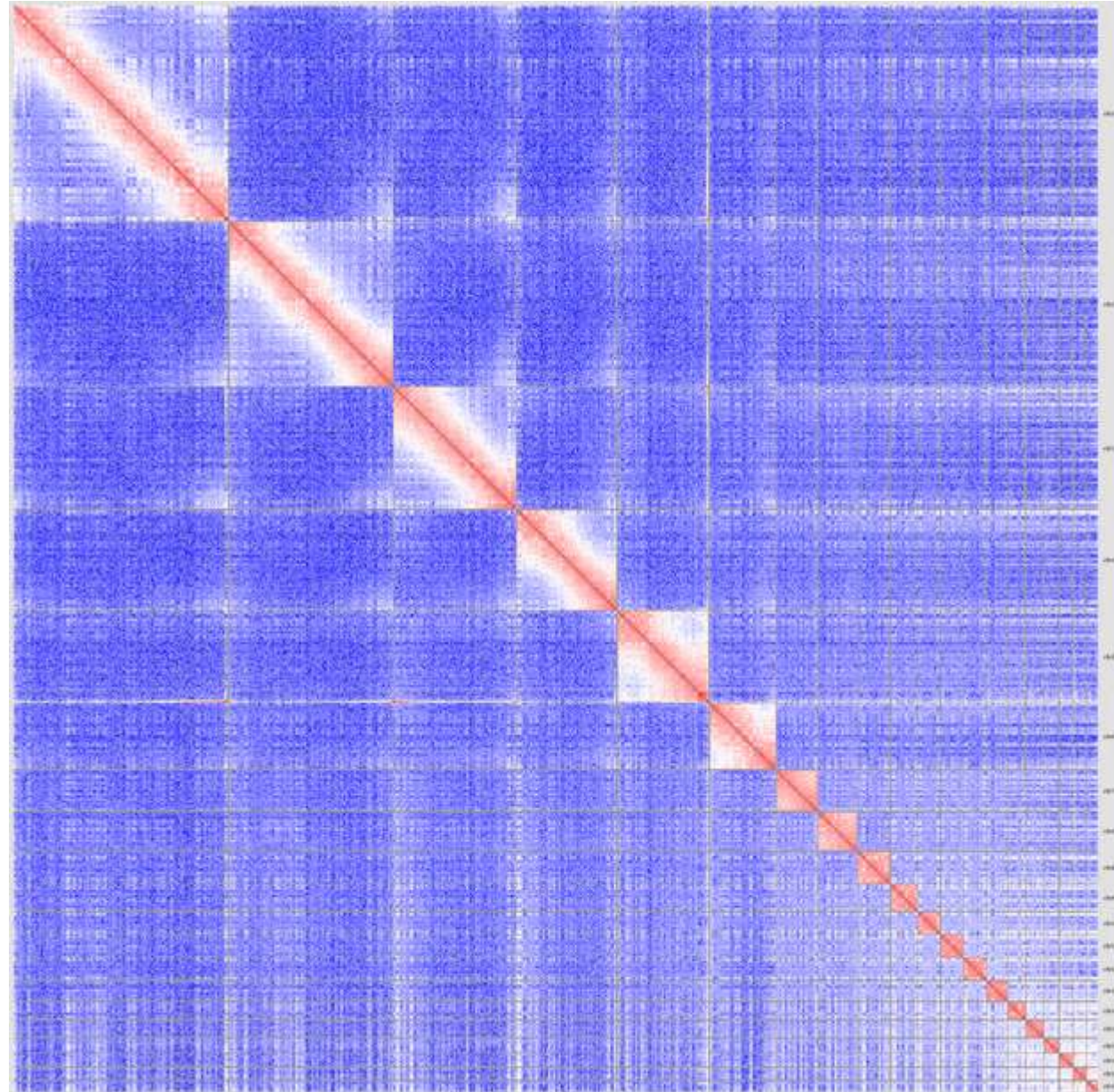
pig muscle,  
snap frozen,  
chr1

pig liver,  
isopentane,  
chr1



# FR-AgENCODE: Hi-C preliminary results

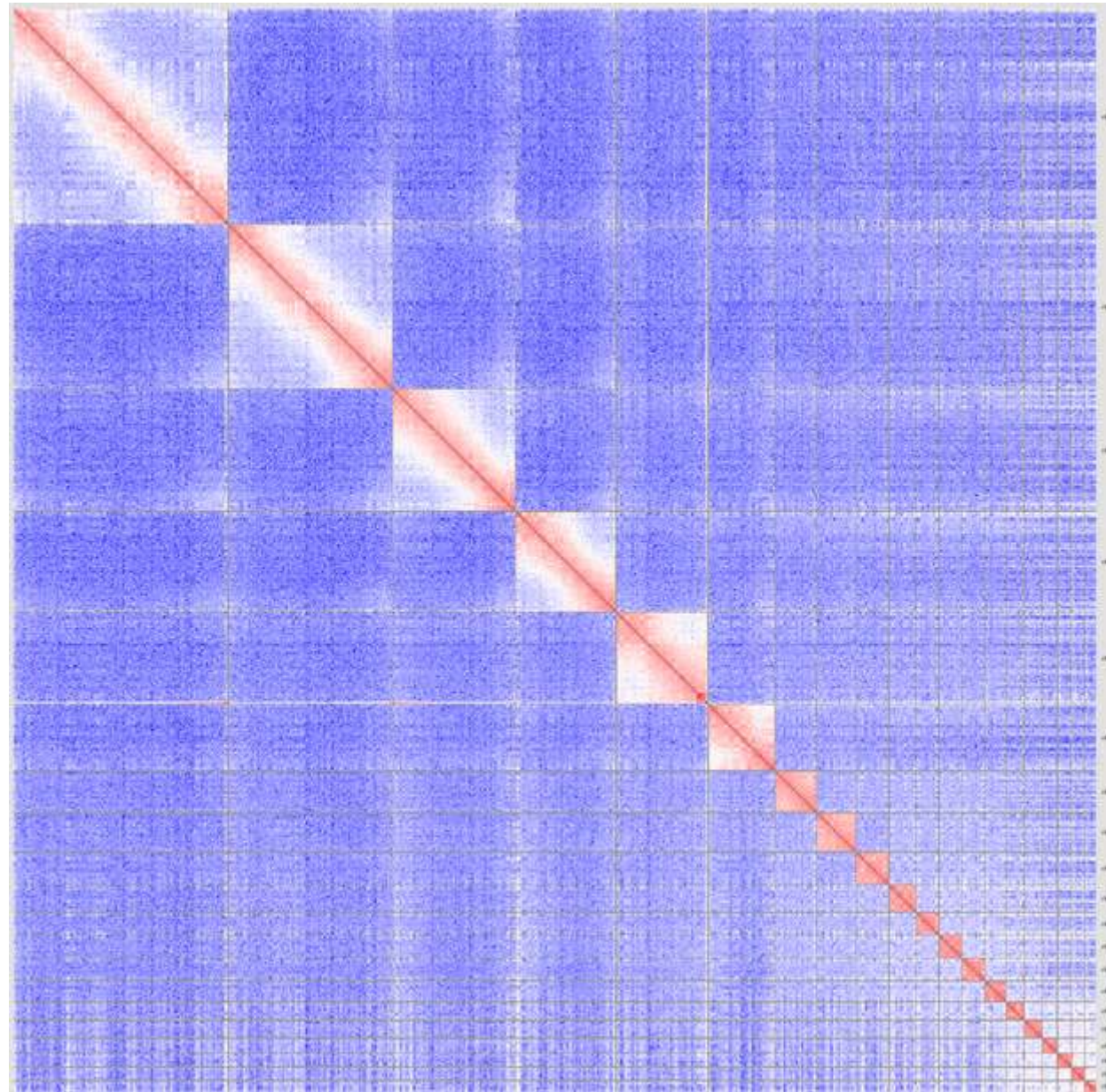
chicken  
liver,  
animal #2





# FR-AgENCODE: Hi-C preliminary results

chicken  
liver,  
animal #4





# FR-AgENCODE: data management and analysis

- ◆ Data production
  - ◆ RNA-seq: 2.2T
  - ◆ Hi-C: 0.8T
  - ◆ ATAC-seq: 1.1T

=> Expected total: > 3T of raw sequence data
- ◆ Data storage
  - ◆ GenoToul ng6
  - ◆ EMBL-EBI (FAANG rapid data release policy)
- ◆ Data analysis
  - ◆ INRA units
  - ◆ EMBL-EBI (FAANG analysis pipelines)

# FR-AgENCODE in FAANG

Bioinformatics and Data Analysis Committee



- ◆ aim: define standard pipelines
- ◆ Working Groups
  - ◆ transcriptome: RNA-seq, lncRNA-seq, sRNA-seq
  - ◆ regulation: ChIP-seq
  - ◆ methylation: WGBS, RRBS
  - ◆ chromatin structure: Hi-C, DNase-seq, ATAC-seq
- ◆ activities: teleconferences, seminars, hackatons
  - ◆ identification of reference datasets
  - ◆ list and benchmark tools
  - ◆ publicly report



# Summary / Conclusion

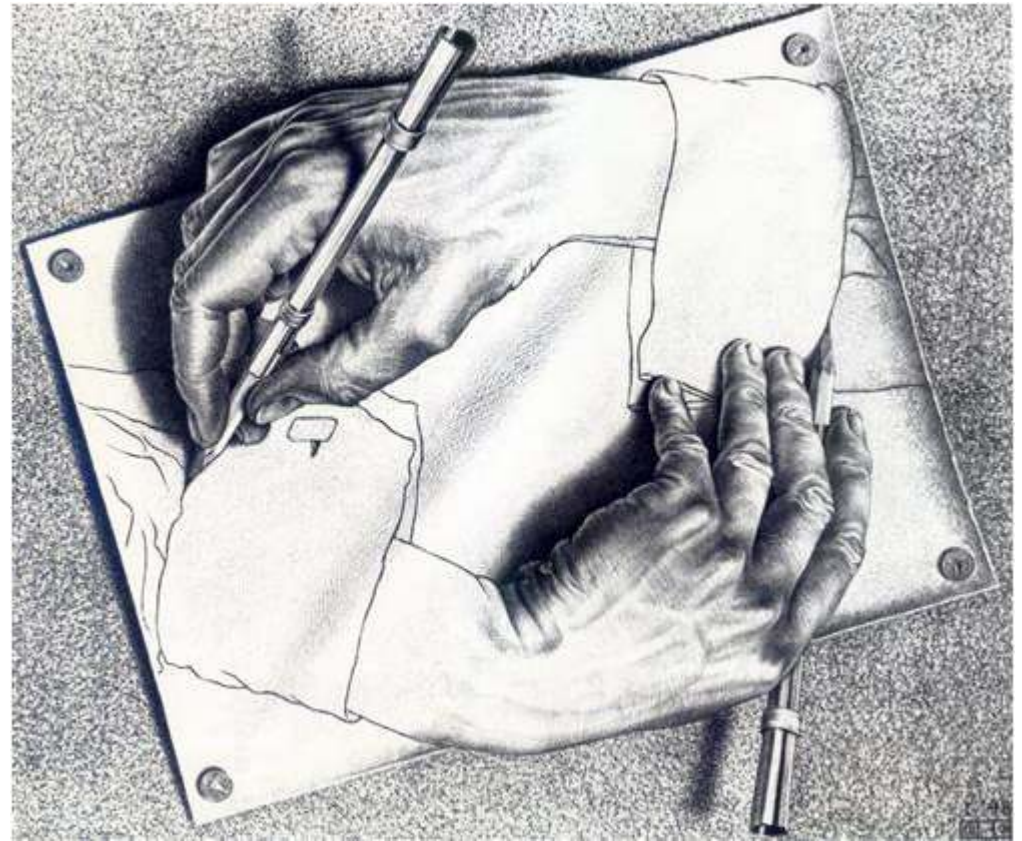
- ◆ Hi-C allows to capture 3D conformation of the chromatin
- ◆ INRA protocol for tissue samples from livestock species
- ◆ INRA contribution to the FAANG action and the Genome to Phenome challenge
  - ◆ FR-AgENCODE pilot project
  - ◆ Samples and Assays committee
  - ◆ Bioinformatics and Data Analysis committee
- ◆ Success story incoming?



*To Be Continued...*

# Outline

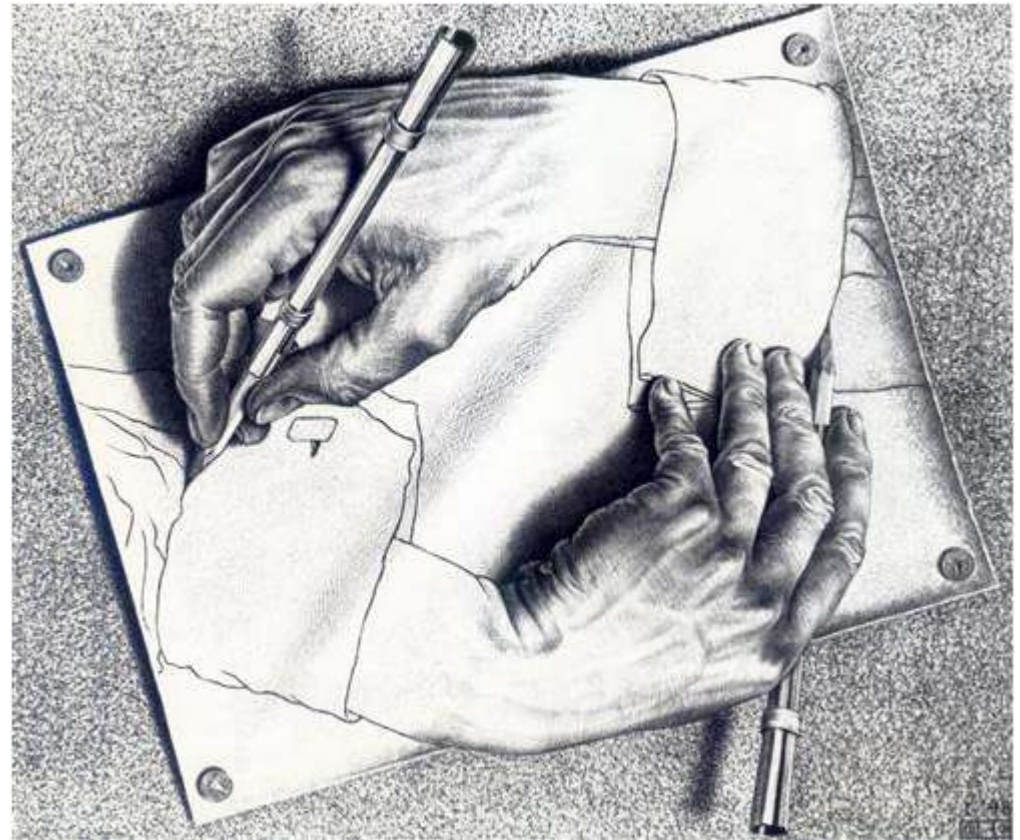
- ◆ Why
- ◆ What
- ◆ How
- ◆ Where ←
- ◆ Who



*M.C. Escher, 1948*

# Outline

- ◆ Why
- ◆ What
- ◆ How
- ◆ Where
- ◆ Who ←



*M.C. Escher, 1948*





# Acknowledgments

## FR-AgENCODE members

- ◆ Management: Elisabetta Giuffra (coordination), Sandrine Lagarrigue, Marie Hélène Pinard
- ◆ Sampling: Michèle Tixier-Boichard, Stéphane Fabre et al.
- ◆ Assays: **Diane Esquerré, Hervé Acloque** et al.
- ◆ Analysis: Christophe Klopp, Christine Gaspin, **David Robelin, Matthias Zytnicki, Sarah Djebali, Magali San Cristobal, Ignacio Gonzalez**, Kylie Munyard, Céline Noirot, Nathalie Villa Vialaneix, Gaelle Lefort, Marjorie Mersch, Frédérique Pitel et al.

## Hi-C team @ GenPhySE, INRA

- ◆ H. Acloque, M. Yerle, Y. Lahbib, F. Mompарт, M. Marti et al.

## FAANG B&DA committee

- ◆ L. Clarke, D. Zerbino, J. Reecy, P. Ross, L. Eory, M. Watson et al.