

# Using statistical profiling to decipher hidden chromatin contacts resulting from repeated sequences

---

26/11/2024

Sébastien Gradit

Spatial Regulation of Genomes - Genomes & Genetics - UMR3525

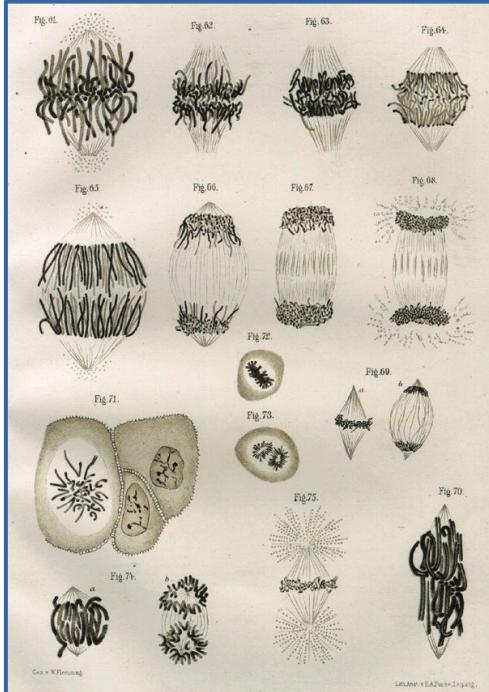
*Thesis supervisor: Axel Cournac*



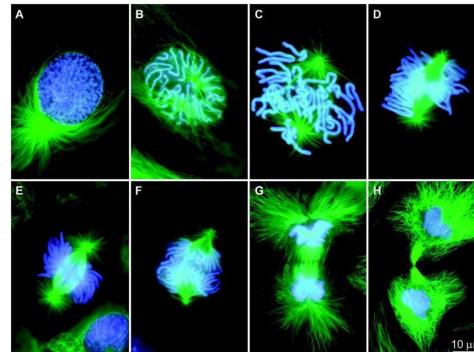
# Early observation of genomes organisation



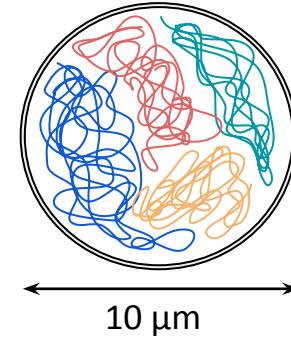
Lilium corceum  
drawings



DNA folding is a highly organized process necessary to fit DNA in nucleus



Human cell



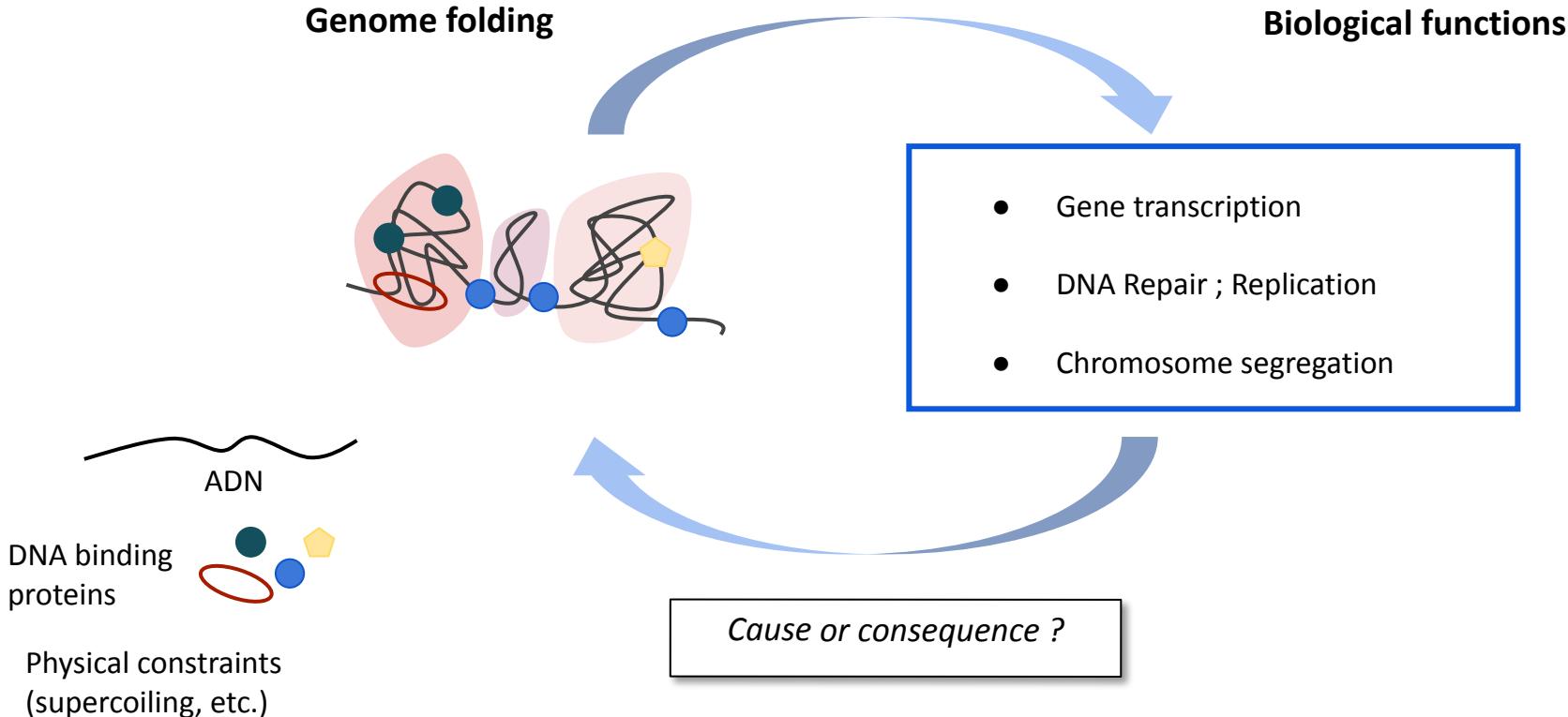
$$3.10^9 \text{ bp} = 2\text{m}$$

yeast



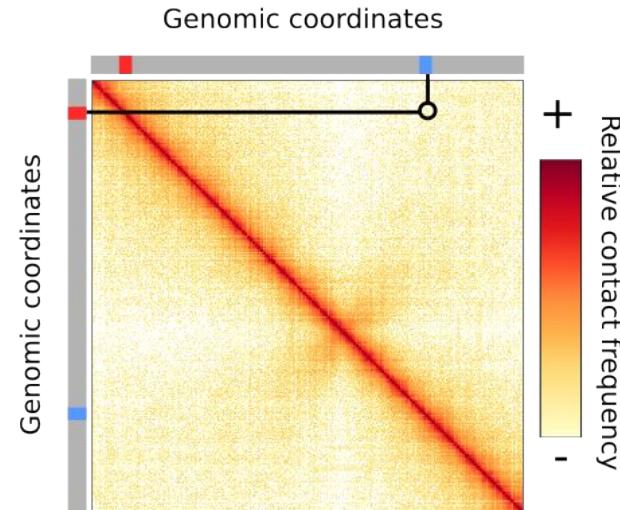
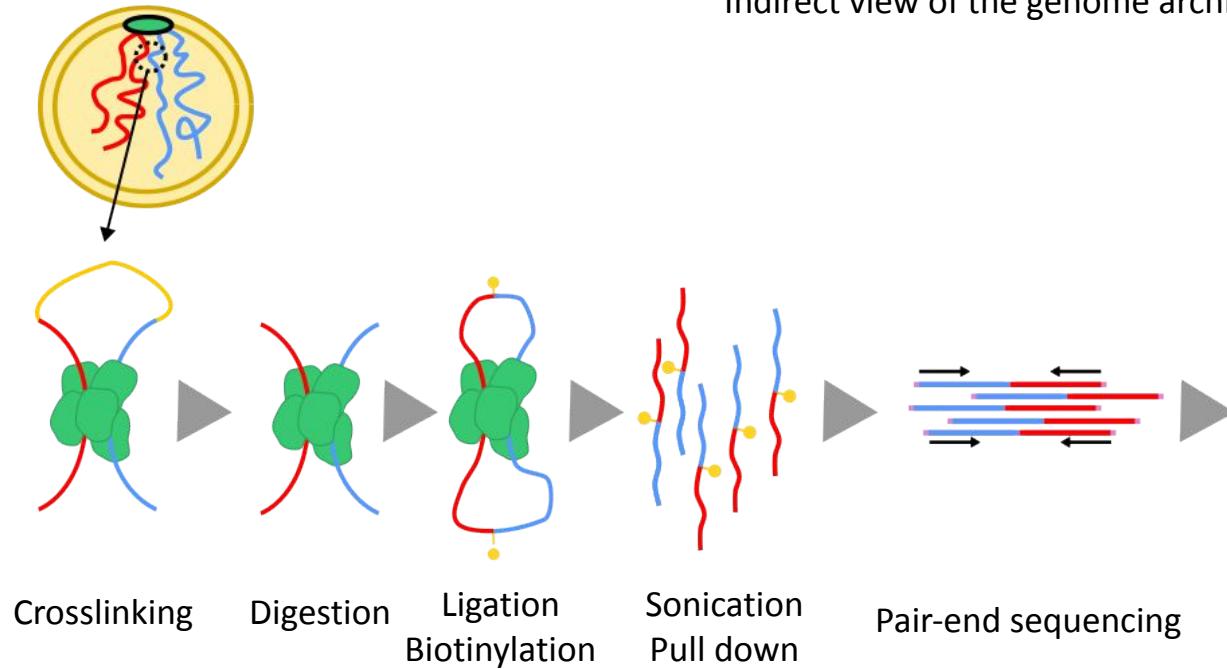
$$12.10^6 \text{ bp} = 8 \text{ mm}$$

# Chromosome folding influences or regulates dynamic processes?



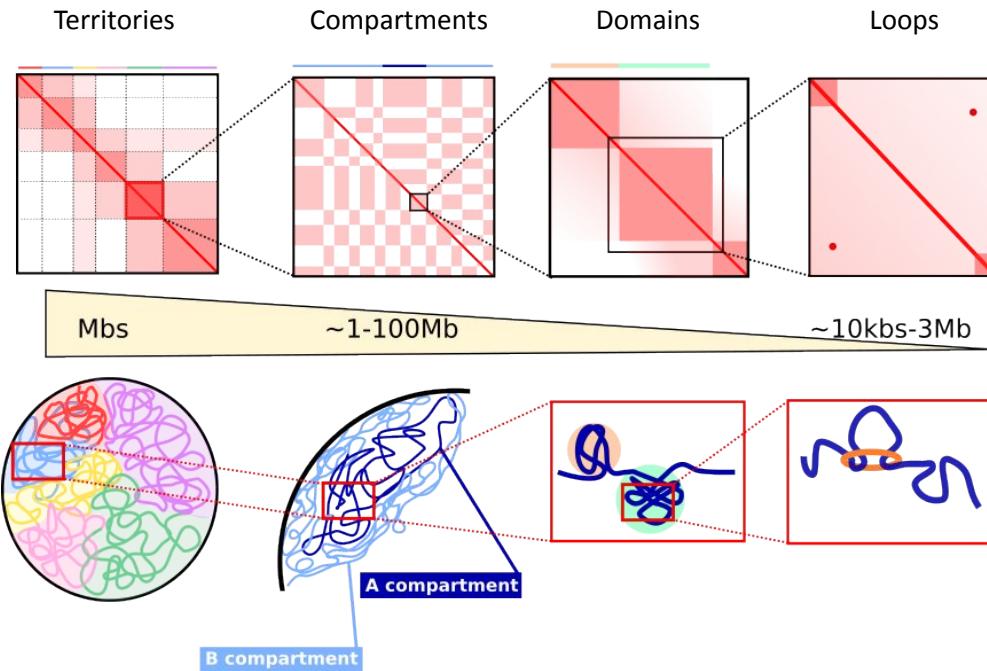
# Hi-C to capture genomes 3D structures

Indirect view of the genome architecture



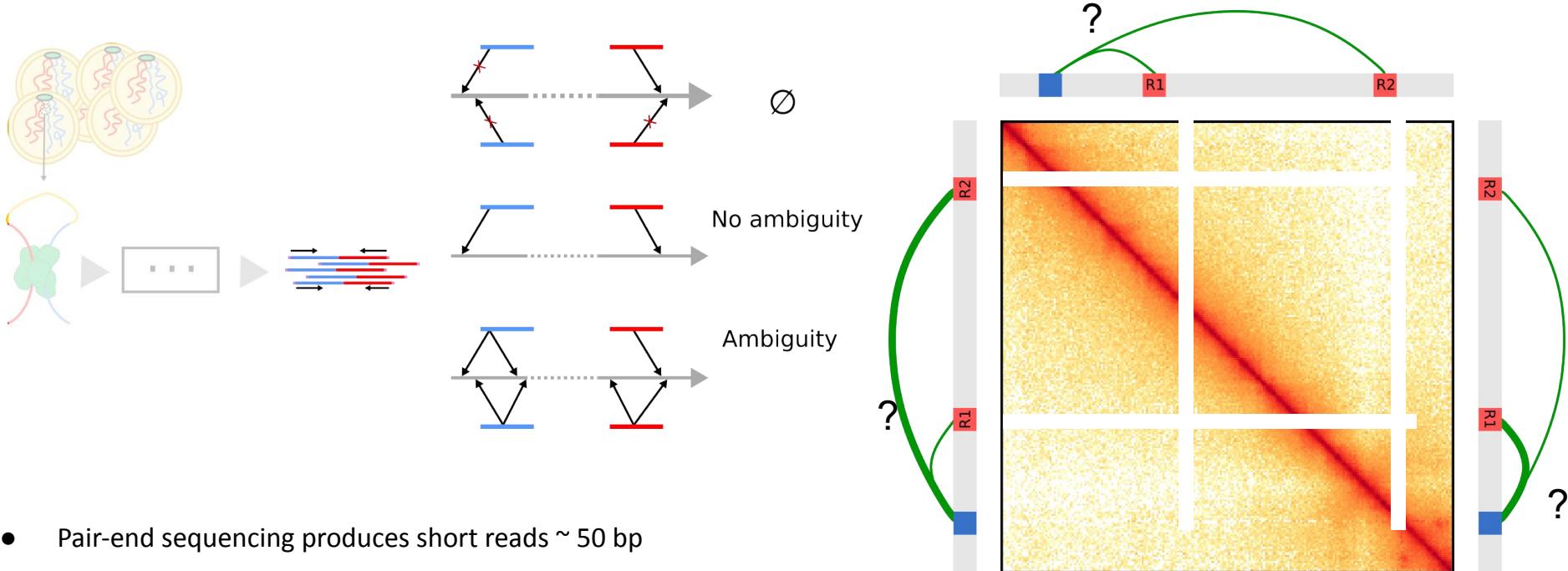
# Chromosome Configurations: Influencing Genome Organization

How these different spatial structures correlates with biological processes?



- Cell cycle phase dependence
- Response to stresses
- Quiescence / replicating?
- Viral infection?

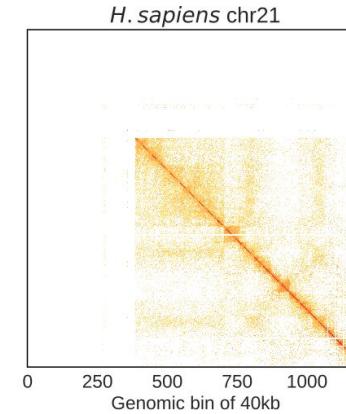
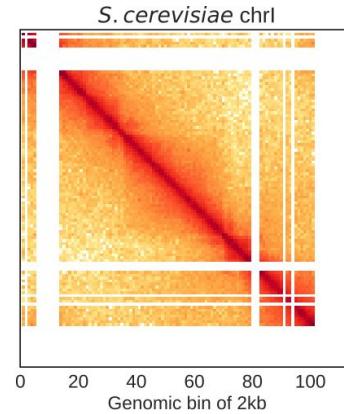
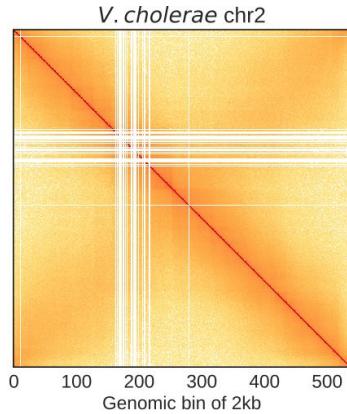
# Limitations of Traditional Hi-C



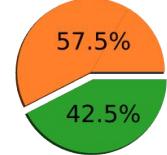
- Pair-end sequencing produces short reads  $\sim 50$  bp
- Genomes carry repeated regions  
→ Some reads can have multiple genomic positions

Statistical approach for ambiguous contacts prediction?

# Thesis Objectives



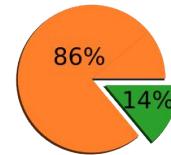
*V. cholerae*



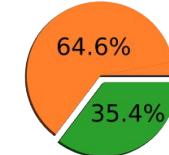
■ Non repeated sequences  
■ Repeated sequences



*S. cerevisiae*



*H. sapiens*



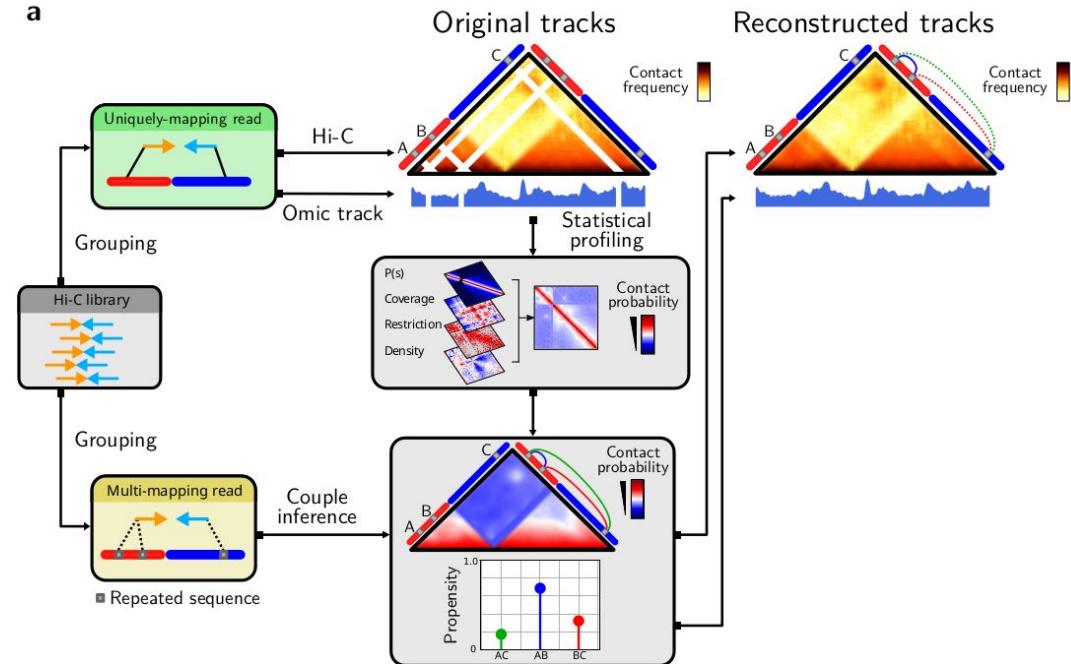
Can a statistical modeling approach be used to infer missing information in Hi-C data and accurately reconstruct the 3D interactions of repeated sequences?

## **Hicberg: Reconstruction of genomic signals from repeated elements**

---

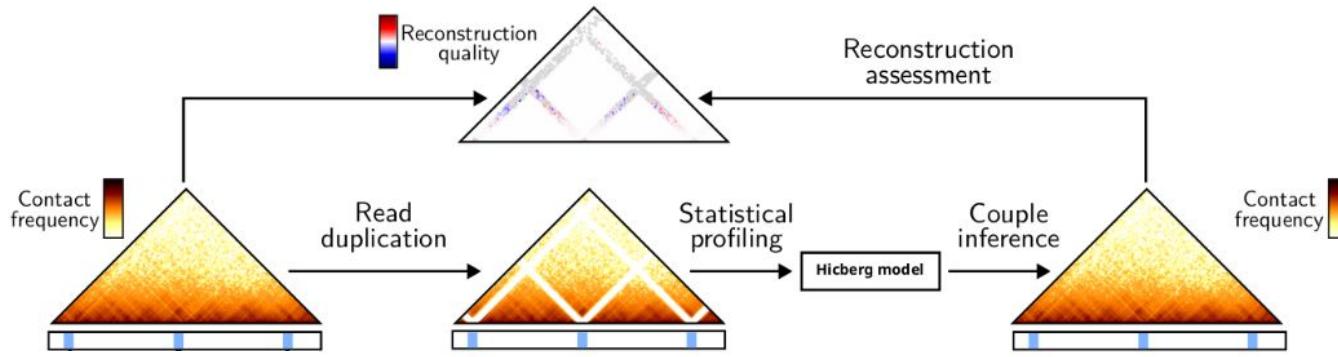
# Hicberg: a novel tool for (3D) genomics

a



- **Hicberg:** Reconstruction of Hi-C data, and more!
- **On-the-fly statistical profiling:** extraction of statistical tendencies from unambiguous data.
- **Probability Mass Function estimation:** Guides accurate placement of ambiguous reads.

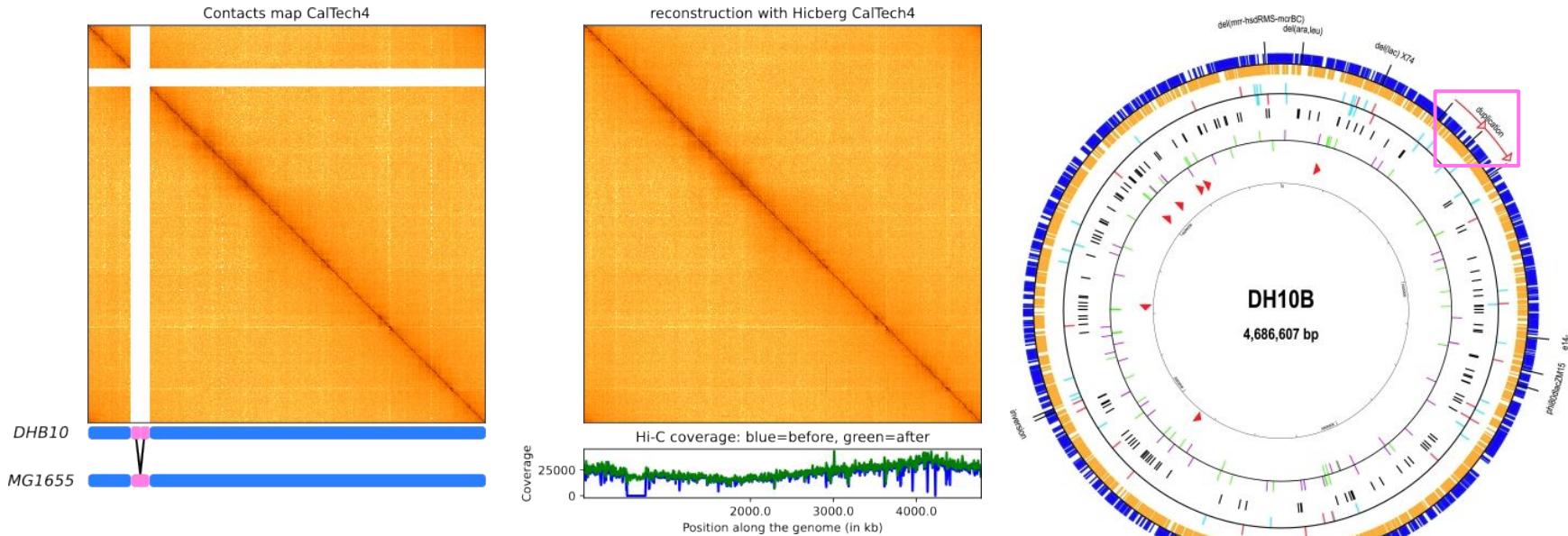
# Benchmarking & In-silico Evaluation



- Simulated repeated sequences by duplicating unambiguous reads.
- Robust across various repeat sizes, numbers, and spacing.
- Polymeric behavior component is key for accurate inference.

Does Hicberg perform as well in real-context?

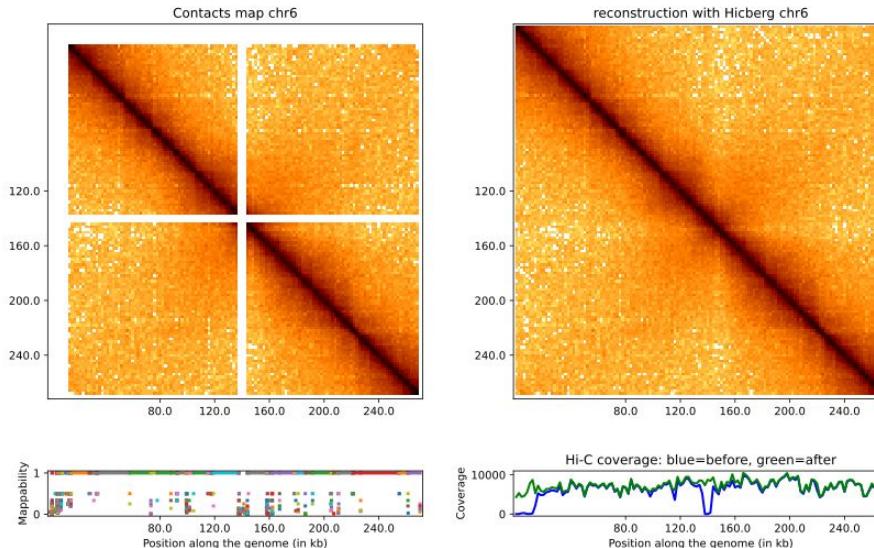
# Biological validation



- Accurate reconstruction of the complex 3D structure of *E. coli* carrying tandem duplication.
- Hicberg captured expected DNA polymer behavior, demonstrating versatility and power.

# Biological validation

---



- Accurate reconstruction of the complex 3D structure of *S. cerevisiae* chr VI carrying repeats.
- Hicberg captured expected DNA polymer behavior, demonstrating versatility and power.
- Hicberg recovered expected mean coverage despite low theoretical mappability

# Unveiling Hidden Connections: Hicberg in action

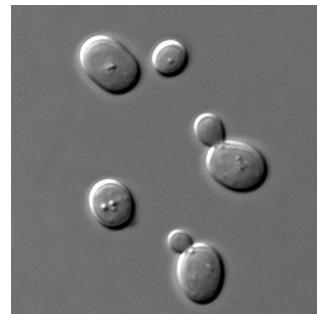
---

How do repetitive elements shape the 3D genome and influence its function?

---

- 1) *Study of Ty1 Retrotransposons: 3D Organization and Cohesin Loading*

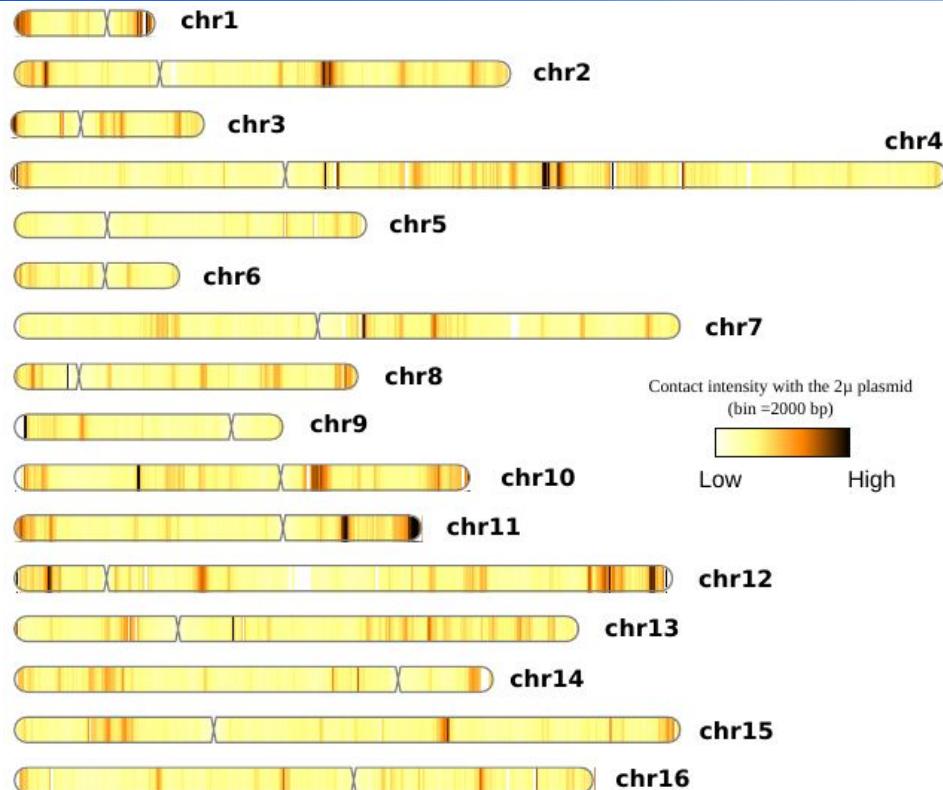
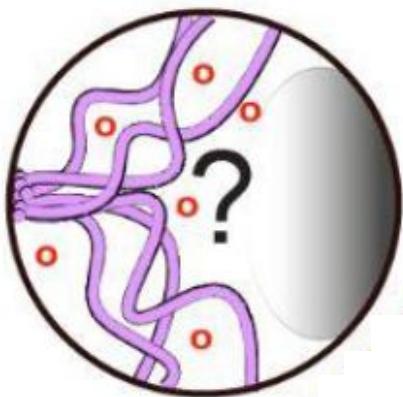
- 2) *Study of rDNA Contact Dynamics in Response to Stress*



## **Beyond Parasitic Elements: Ty1 Retrotransposons as Organizers of Nuclear Architecture**

---

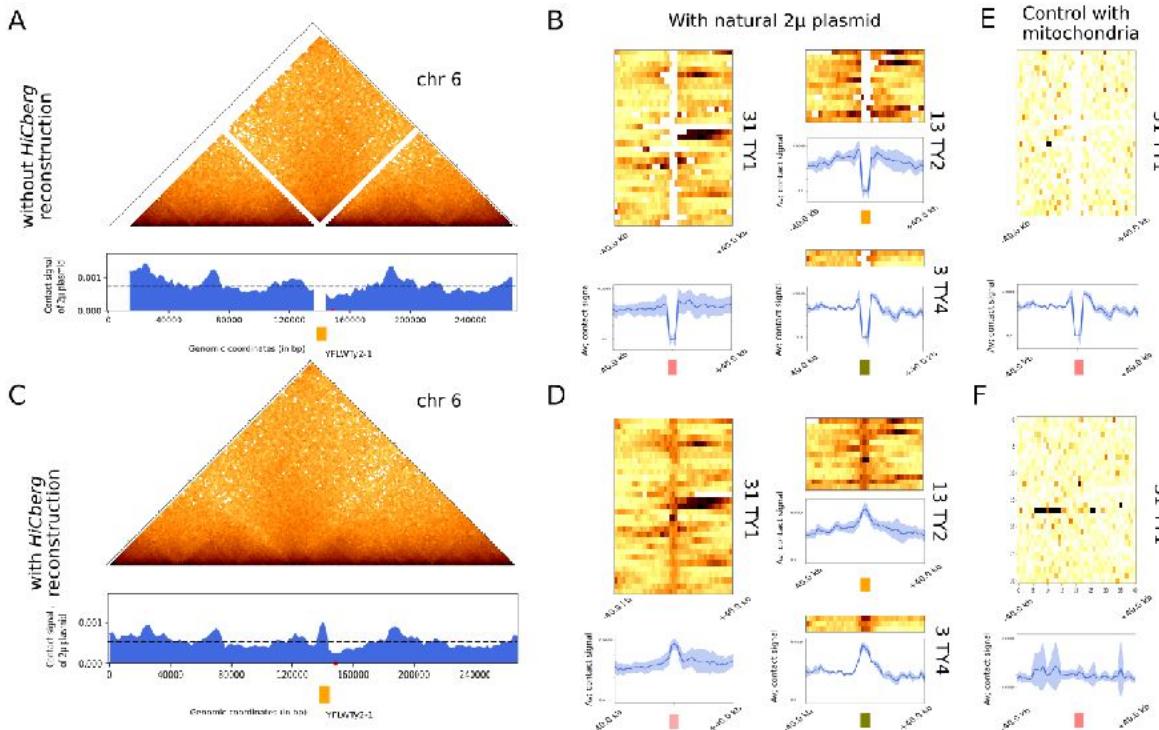
# Application : *S. cerevisiae* TYs as new 2 $\mu$ preferencial contact spots



Presence of natural 2 $\mu$  plasmid in *S. cerevisiae*

Detection of ~ 75 hotspots of contact

# Application : *S. cerevisiae* TYs as new 2μ preferencial contact spots

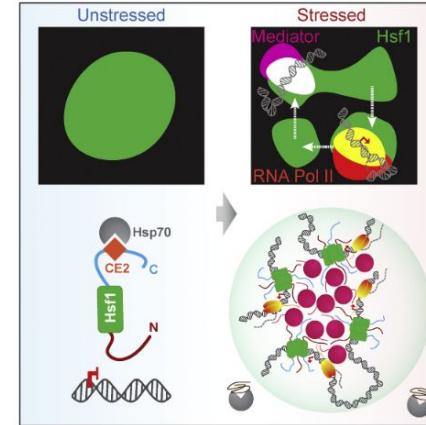
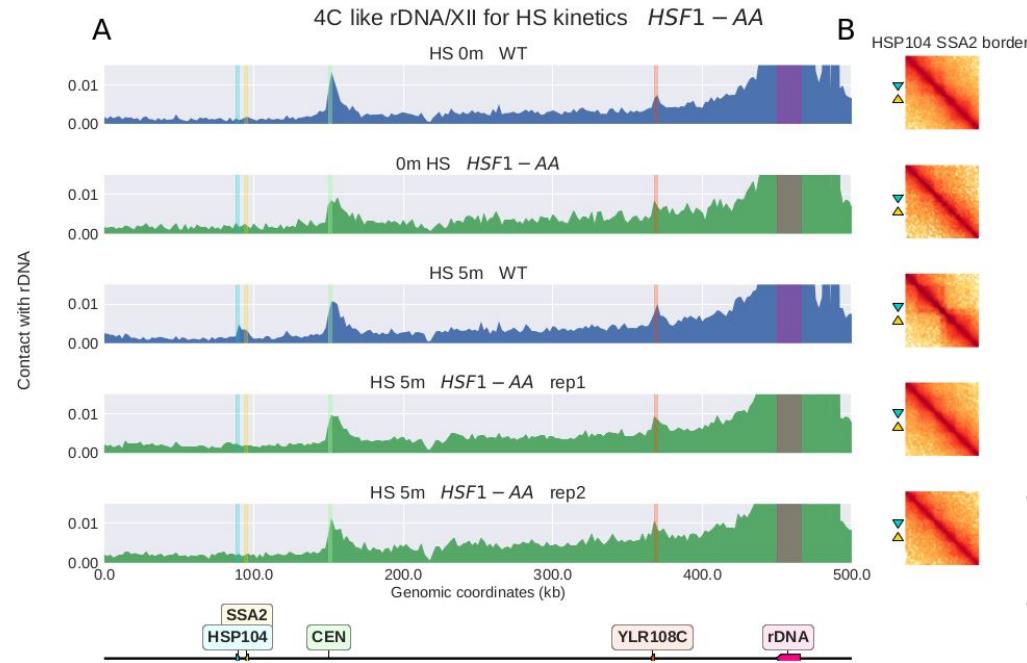


Retrotransposons (TYs) appear to be new hotspots of 2μ plasmid

## **Unveiling the rDNA Network: Key Players and Choreographers of its Spatial Organization**

---

# Application : *S. cerevisiae* TYs as new 2μ preferencial contact spots



- **Hsf1 Mutant:** Transient rDNA-HSP104/SSA2 interaction and border pattern are lost.
- **Dual Requirement:** Both transcription and Hsf1 are essential for this stress-induced structure.
- **Mechanism:** Supports Chowdary's hypothesis of micro-aggregates, not loop extrusion, for contact formation.

- Hicberg showed **good performances** both in **simulated** and in **biological** data
- Several genomic features have been unveiled in *S. cerevisiae*:
  - 2μ plasmid **hotspots of contacts** of yeast retrotransposons
  - **New rDNA interaction during heat shock**, involving *Hsf1* protein complex

→ First link about spatial organization of rDNA and **stress**

---



## Régulation Spatiale des Génomes

Romain Koszul  
**Axel Cournac**  
Martial Marbouthy  
Agnès Thierry  
Hélène Bordelet  
Jacques Serizay  
**Pauline Larrous**  
Devon Conti  
Justine Groseille  
Manon Perrot  
Maëlys Delouis  
Corina Pascal

### Anciens membres

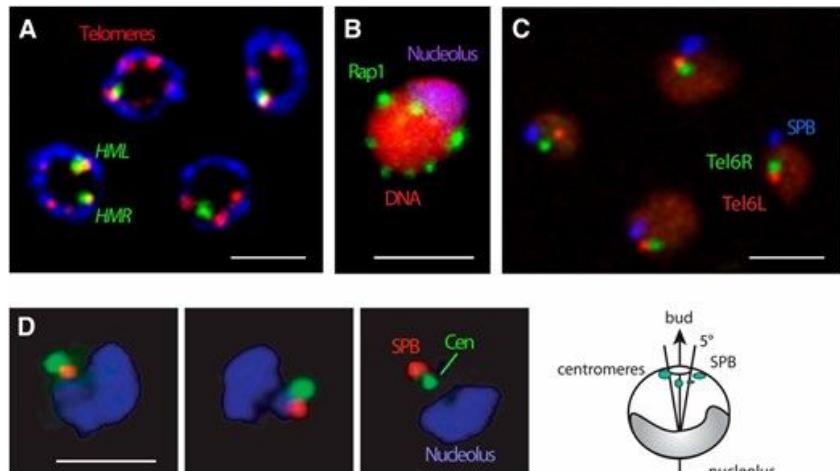
Cyril Matthey-Doret  
Amaury Bignaud  
Léa Meneu  
Samuel Ortion

### Collaborators

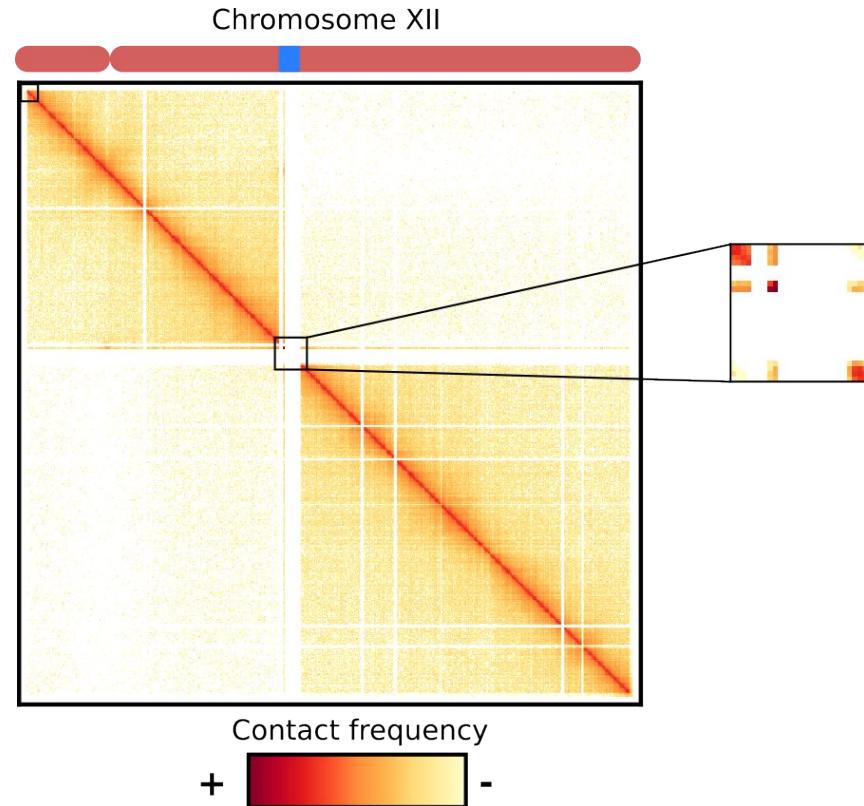
Bouk Wim De Jong



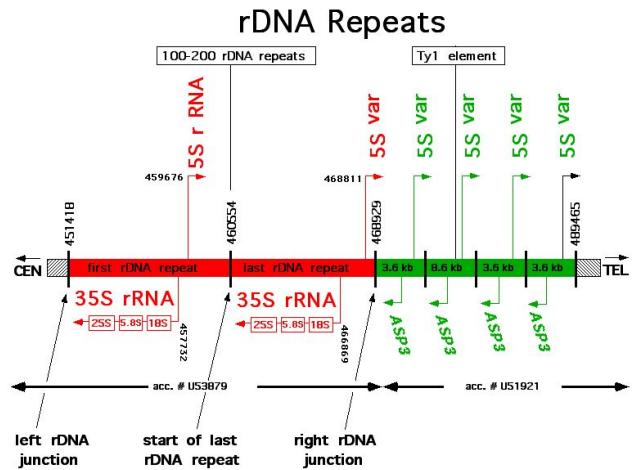
# rDNA: A Landmark in the Nuclear Landscape



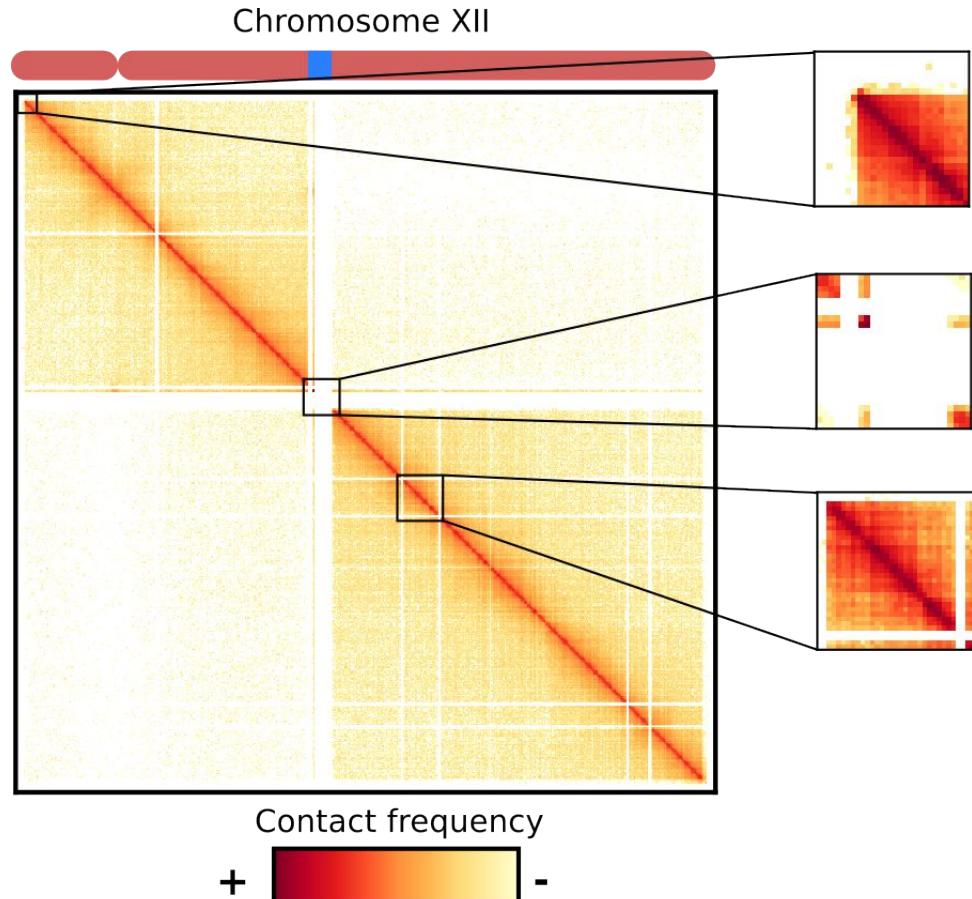
- rDNA as an organizational feature within the nucleus
- Ribosomes biogenesis
- Often exhibits interactions with nuclear landmarks (telomere, centromeres)



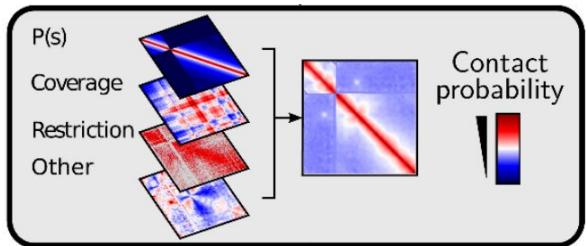
# The Challenge of Repeated Sequences



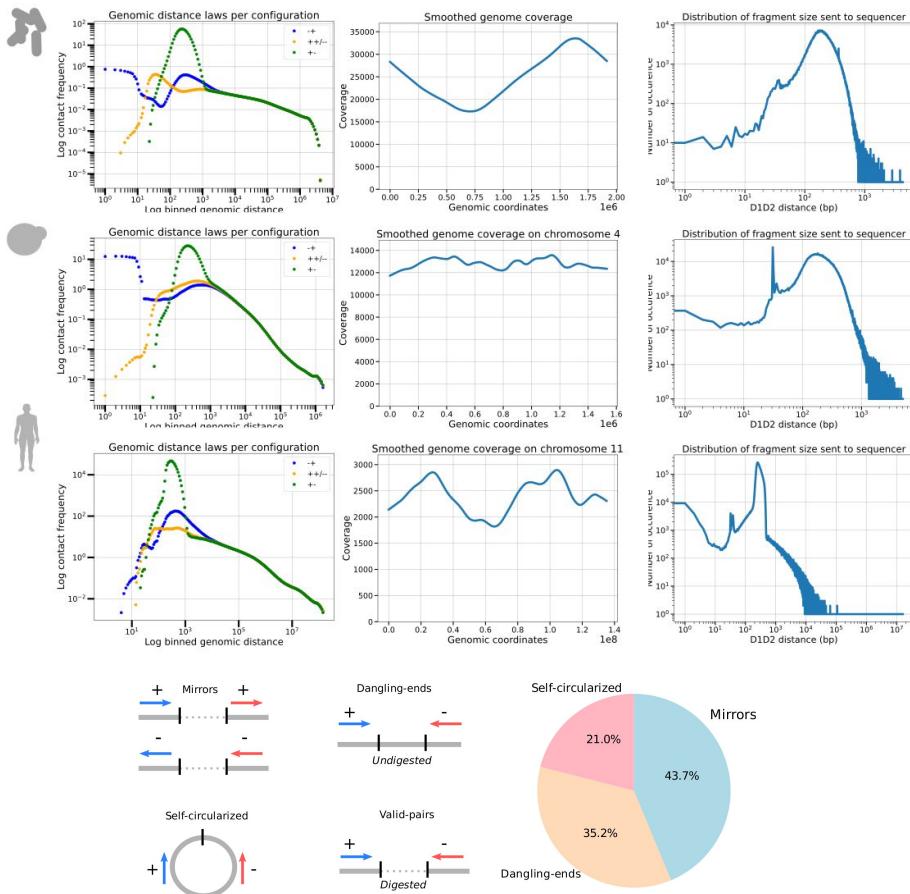
Why are repetitive sequences, particularly challenging to study using traditional Hi-C?



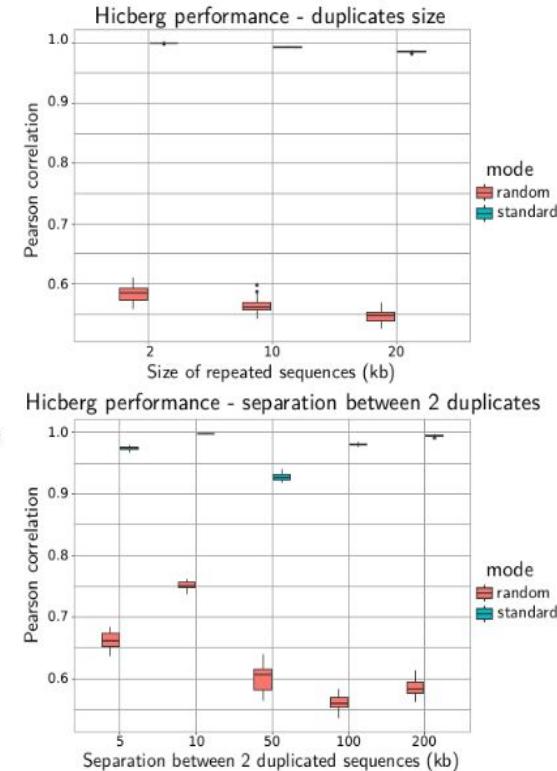
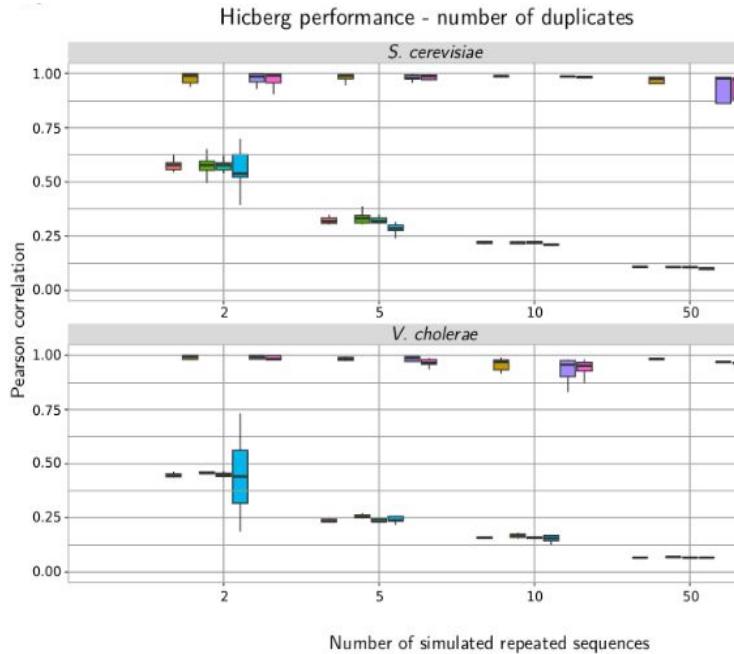
# Hicberg: a novel tool for (3D) genomics



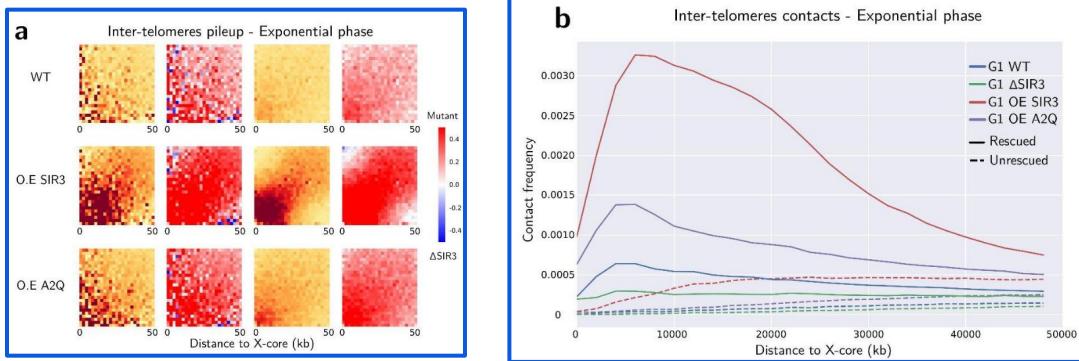
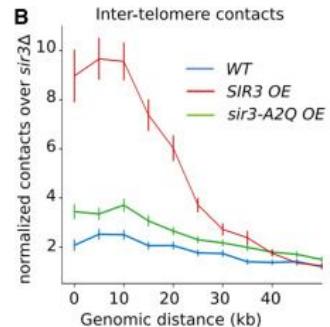
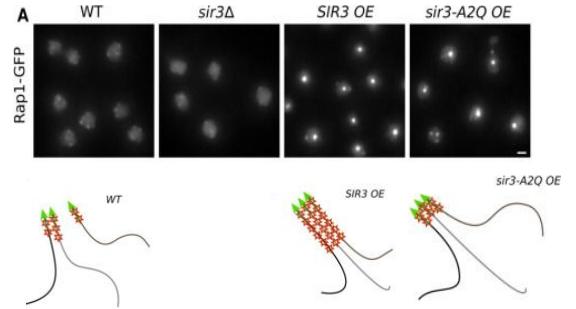
- **Polymer Behavior:** Captures DNA flexibility with 3 species/condition-specific sub-laws.
  - **Coverage Bias:** Accounts for uneven read distribution; higher coverage increases interaction probability.
  - **Restriction Site Proximity:** Accounts for fragment length biases.



# In-silico Evaluation



## Biological validation

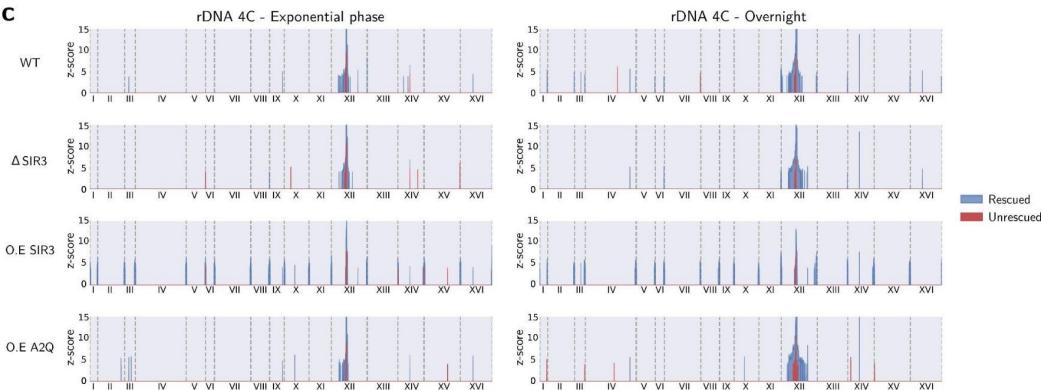
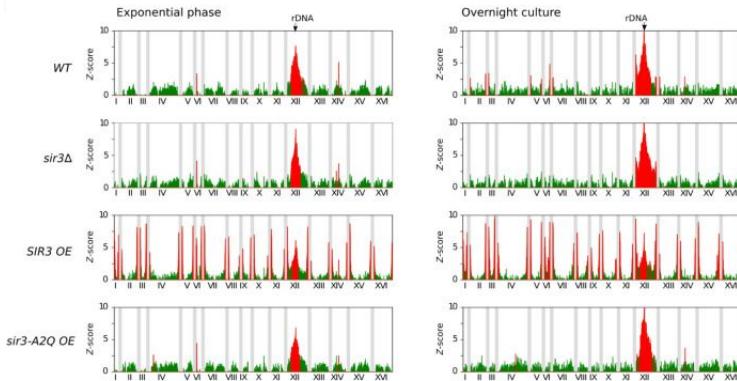
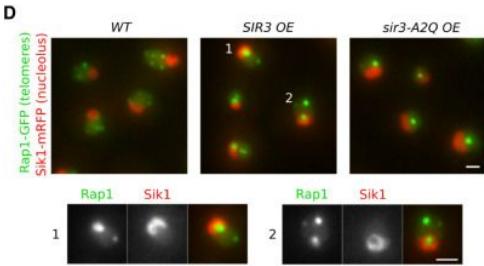


Hicberg to re-analyze yeast Hi-C data on telomere clustering under SIR3 over-expression studied with a diffusion-based method (Serpentine)

**Hicberg reconstruction:** Accurately reproduced  
Ruault's findings:

- SIR3 overexpression → telomere hyper-clustering.
  - Altered SIR3 (A2Q) → reduced clustering.

# Biological validation



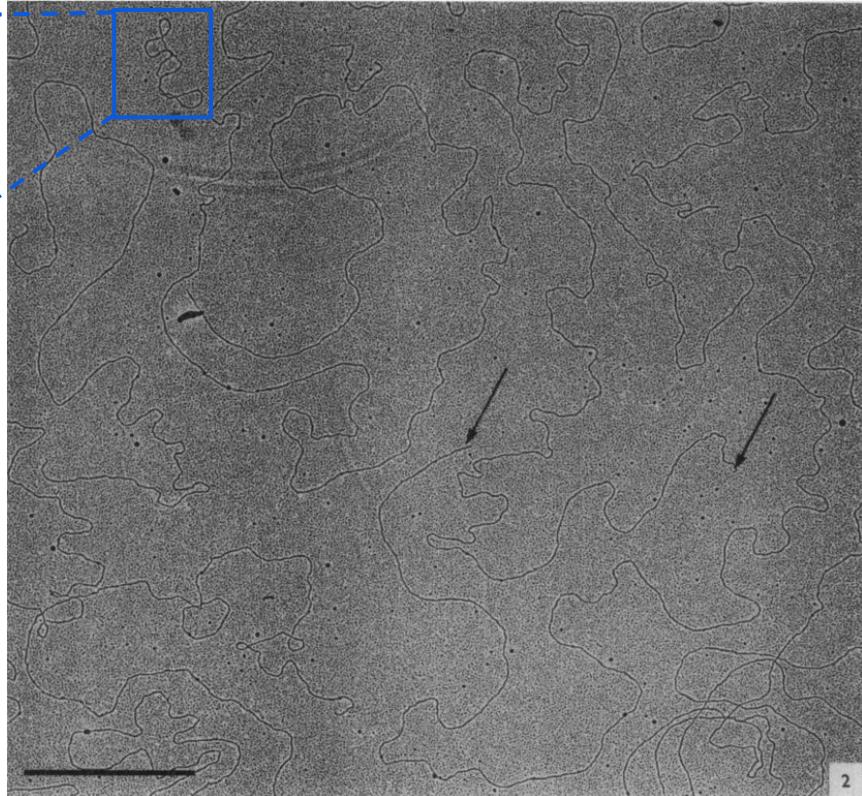
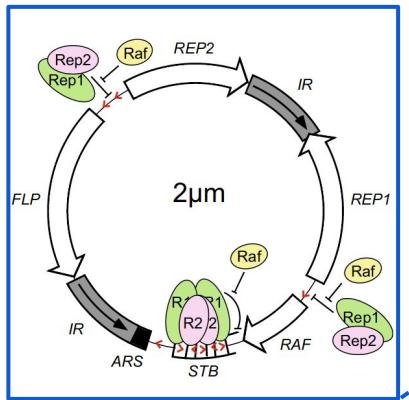
**Hicberg reconstruction:** Reproduced Ruault's findings:

- *SIR3* overexpression → increased rDNA-telomere clustering (exponential and overnight cultures).
- Altered *SIR3* (A2Q) → partially reduced clustering.

**Hicberg accurately captures complex, condition-specific 3D interactions in repeat rich regions, further validating its biological relevance.**

# Discovery of 2μ plasmid in *Saccharomyces cerevisiae*

27



Observation in the 1970's of **enriched fraction of small circular DNA** found as different molecules species

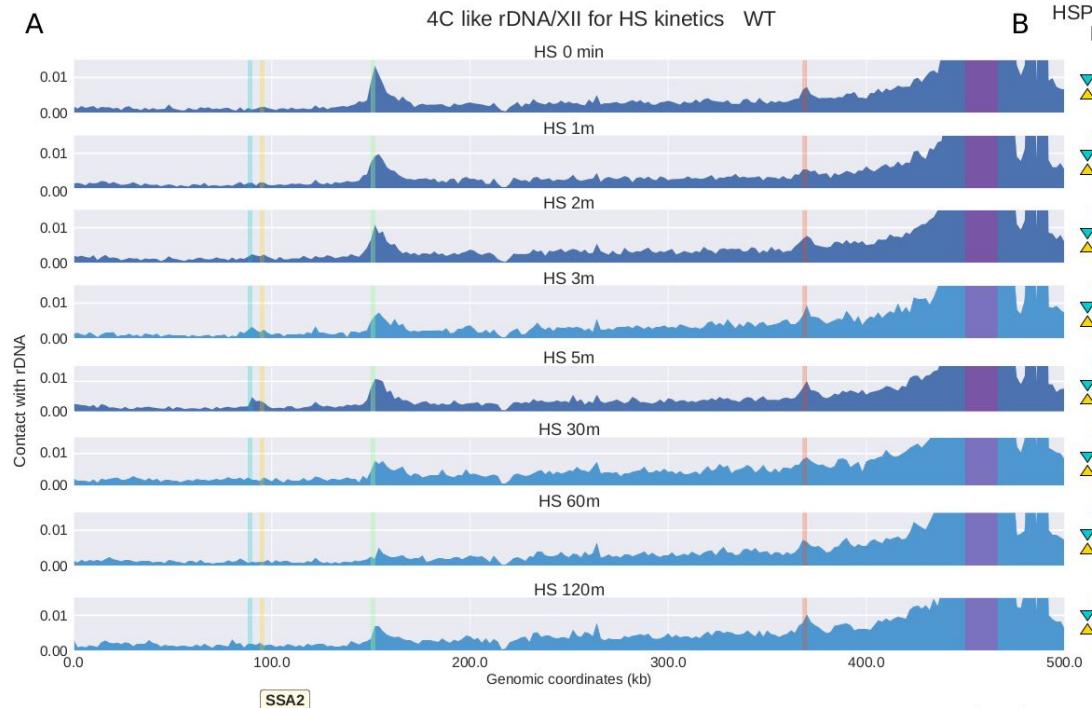
Predominant class of with **mean contour length of  $1.88 \pm 0.11 \mu\text{m}$**

Conserved through generations

Specific locations in *S.cerevisiae* genome (long genes, low transcription spots)

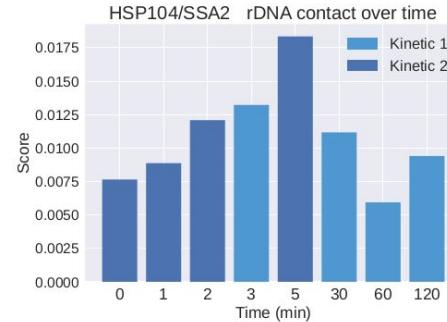
# rDNA-HSP104/SSA2: A Transient Tango During Heat Shock

A



B

HSP104 SSA2 border



- **rDNA in Heat Shock:** Hicberg reveals transient interaction with HSP104/SSA2 genes.
- **Dynamic contact:** Forms and dissolves during heat shock, suggesting a role in stress response.
- **Border pattern:** Emerges around HSP104/SSA2, indicating potential transcriptional activation.

- Hicberg showed good performances both in simulated and in context data
  - Hicberg outperformed mHiC on the definition of reconstructed structures, and number of retrieved Hi-C reads
  - Several genomic features have been unveiled (2μ hotspots on TYs, new rDNA interaction during heat shock, transcriptionally induced aggregates)
  - First link about spatial organizationof rDNA and stress
- 

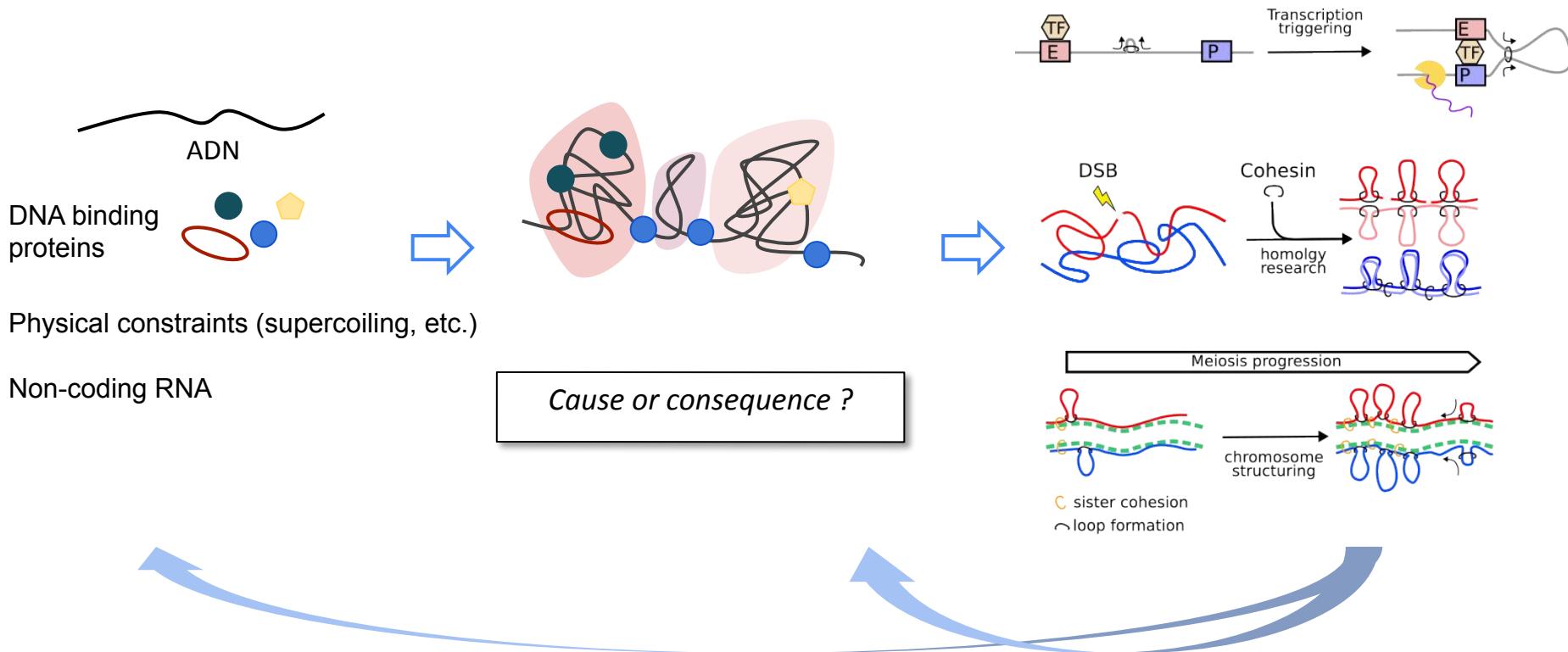
## Expand Applications:

- Explore diverse stress conditions in yeast.
- Apply to other organisms and datasets.
- Investigate various "OMICs" data.

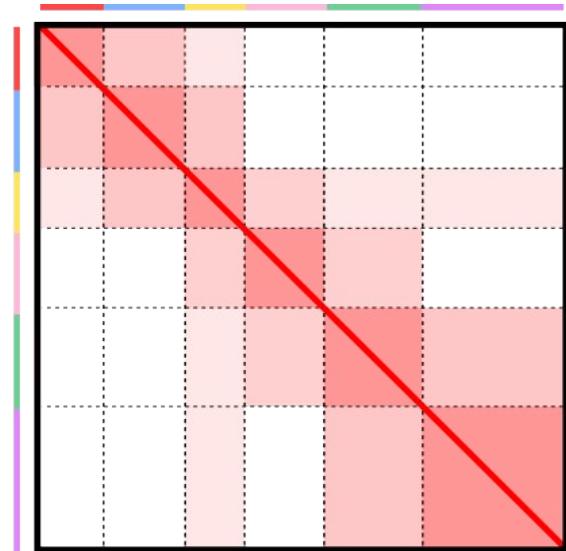
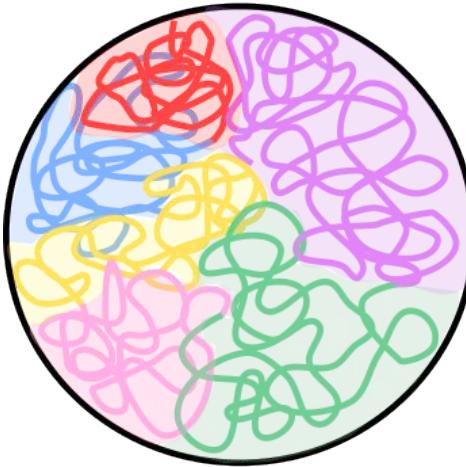
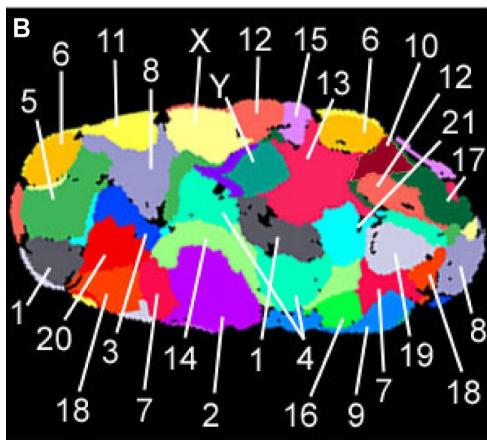
## Enhance Performance:

- Optimize for speed and efficiency.
- Facilitate routine use across species.

# Chromosome folding influences or regulates dynamic processes?



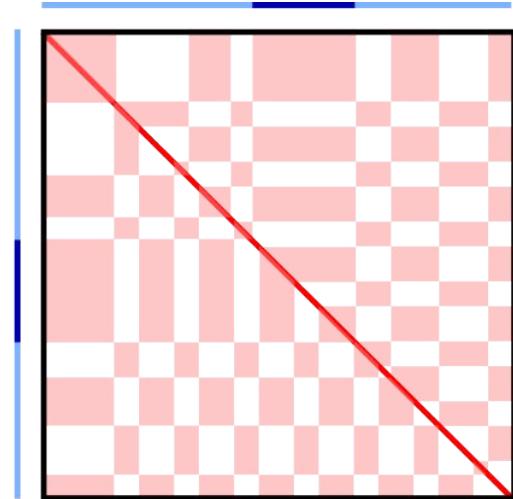
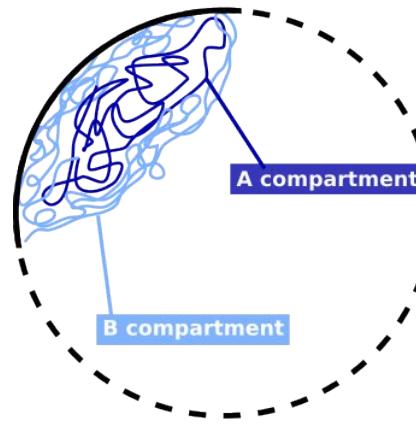
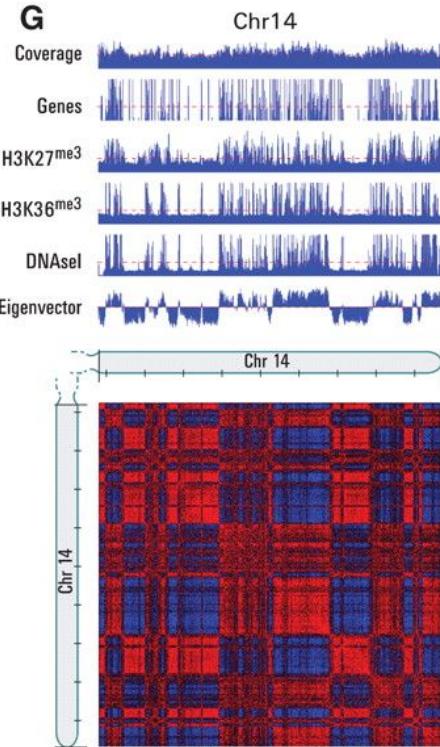
# Chromosome Location in Genomes



- Size/shape vary → chromosome, cell type
- Dynamic state (cellular processes, cycle, ...)
- Gene rich toward interior

- Response to biological processes:
- Efficiency of DNA replication and repair
  - Prevent inter-chromosomes translocations

# From chromosomes to compartments



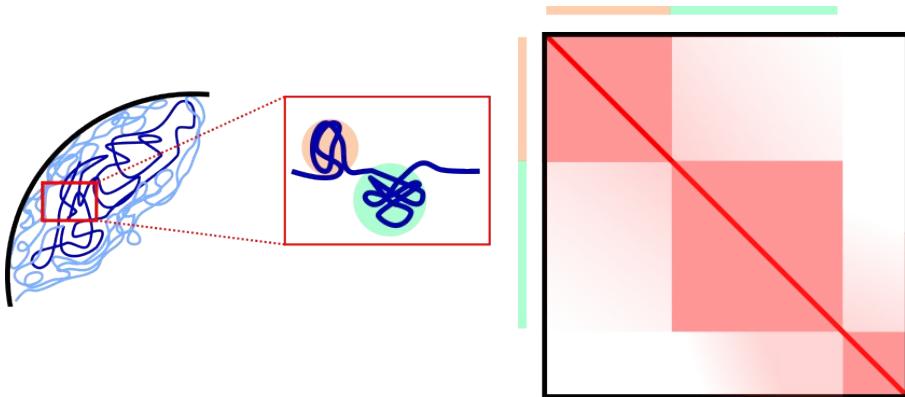
A compartments:

- Open chromatin
- Gene rich
- Acetylation marks
- Interior of nucleus

B compartments:

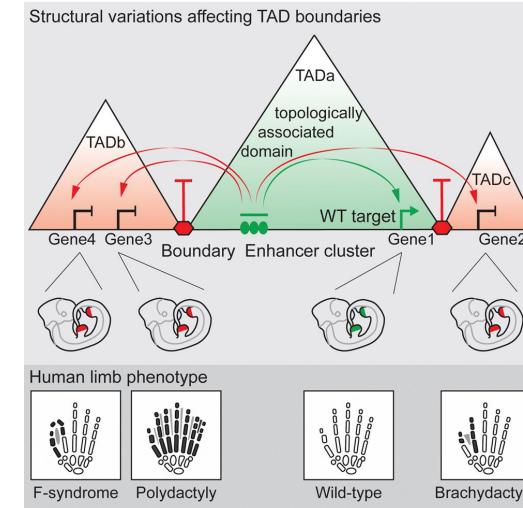
- Dense chromatin
- Gene devoid
- Methylation marks
- Toward lamina

# Topological Associated Domains



High level of self-interaction:

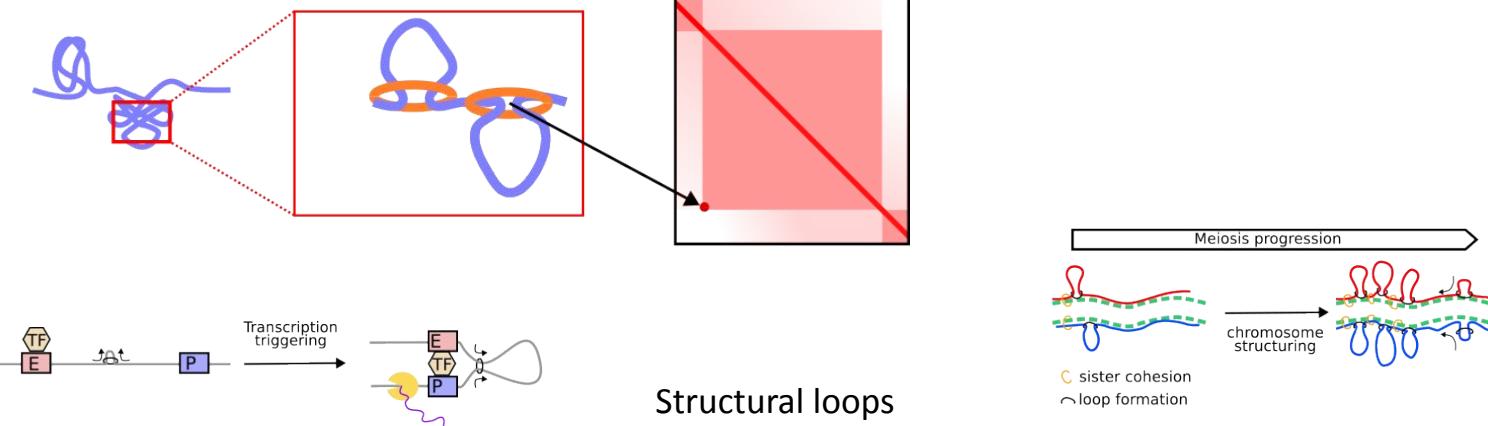
- Boundary regions insulator elements : separators
- Enriched for CTCF or architectural proteins



- Regularity units: appropriate enhancer/promoter
- Boundary disruption = inappropriate gene activation  
→ Disorders (development, cancer, ...)

**How do distant regulatory elements communicate with their target genes to control gene expression?**

# Chromatin loops: Bringing Elements Together



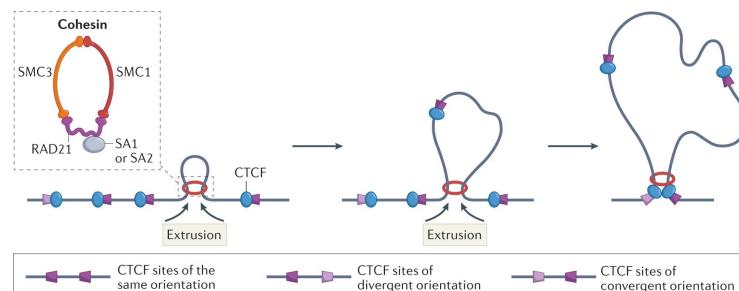
## Regulatory loops

- Modulate gene expression
- Small range, gene expression precise control
- Dynamic: TF, chromatin remodelers, signaling

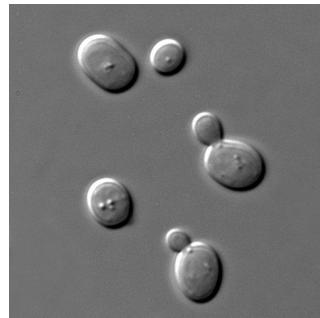
**Dual role, different scales but shared molecular mechanisms**

## Structural loops

- Establish, maintain territories
- Long range, overall genome structuration
- More stable → slower timescale (cellular cycle)



# *S. cerevisiae* as a model for 3D genomics studies



- ~ 12 millions bp
- ~ 6000 genes
- Haploid and diploid states
- Easy to handle

