# HiC-TE: a Nextflow pipeline to study repeat interactions in the 3D genome

Matej Lexa

ENBIK 2022

Němčice u Kolína

13.-15.6.2022

# Genome 3D organization
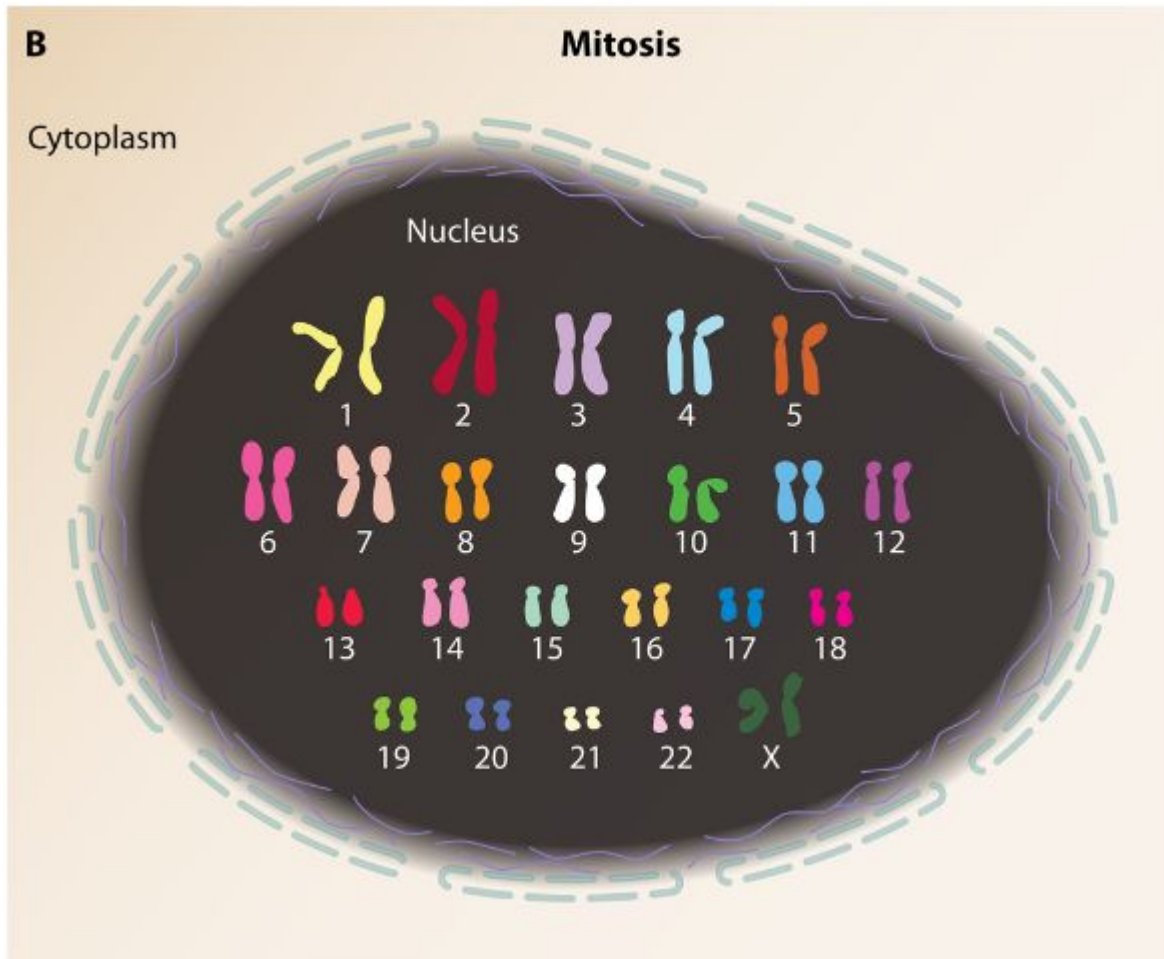
**An Overview of Genome Organization and How We Got There: from FISH to Hi-C**

James Fraser,[a] Iain Williamson,[b] Wendy A. Bickmore,[b] Josée Dostie[a]

Department of Biochemistry, and Goodman Cancer Research Center, McGill University, Montréal, Québec, Canada[a]; MRC Human Genetics Unit, MRC Institute of Genetics and Molecular Medicine, University of Edinburgh, Edinburgh, United Kingdom[b]

# Genome 3D organization



An Overview of Genome Organization and How We Got There: from FISH to Hi-C

James Fraser,[a] Iain Williamson,[b] Wendy A. Bickmore,[b] Josée Dostie[a]

Department of Biochemistry, and Goodman Cancer Research Center, McGill University, Montréal, Québec, Canada[a]; MRC Human Genetics Unit, MRC Institute of Genetics and Molecular Medicine, University of Edinburgh, Edinburgh, United Kingdom[b]
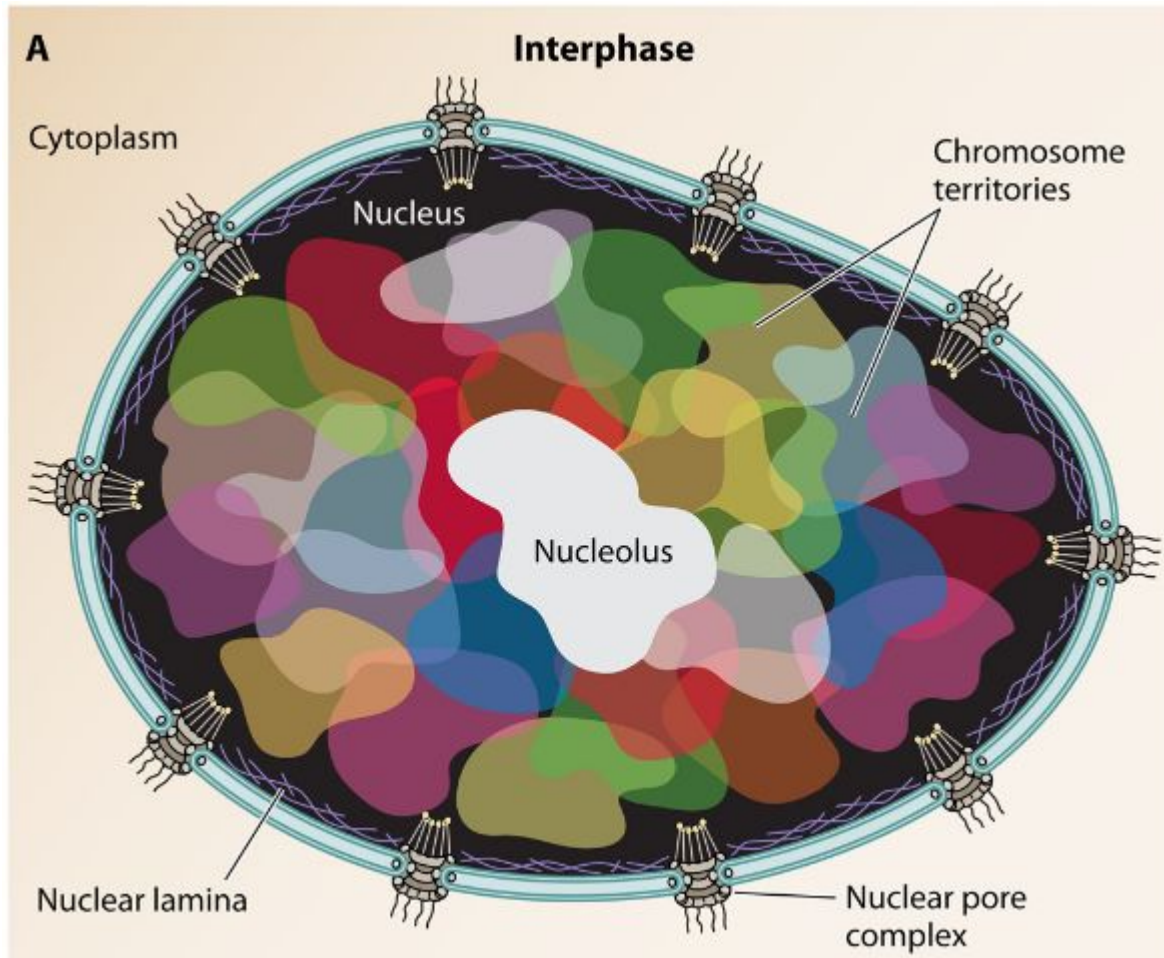
# Genome 3D organization



An Overview of Genome Organization and How We Got There: from FISH to Hi-C

James Fraser,[a] Iain Williamson,[b] Wendy A. Bickmore,[b] Josée Dostie[a]

Department of Biochemistry, and Goodman Cancer Research Center, McGill University, Montréal, Québec, Canada[a]; MRC Human Genetics Unit, MRC Institute of Genetics and Molecular Medicine, University of Edinburgh, Edinburgh, United Kingdom[b]
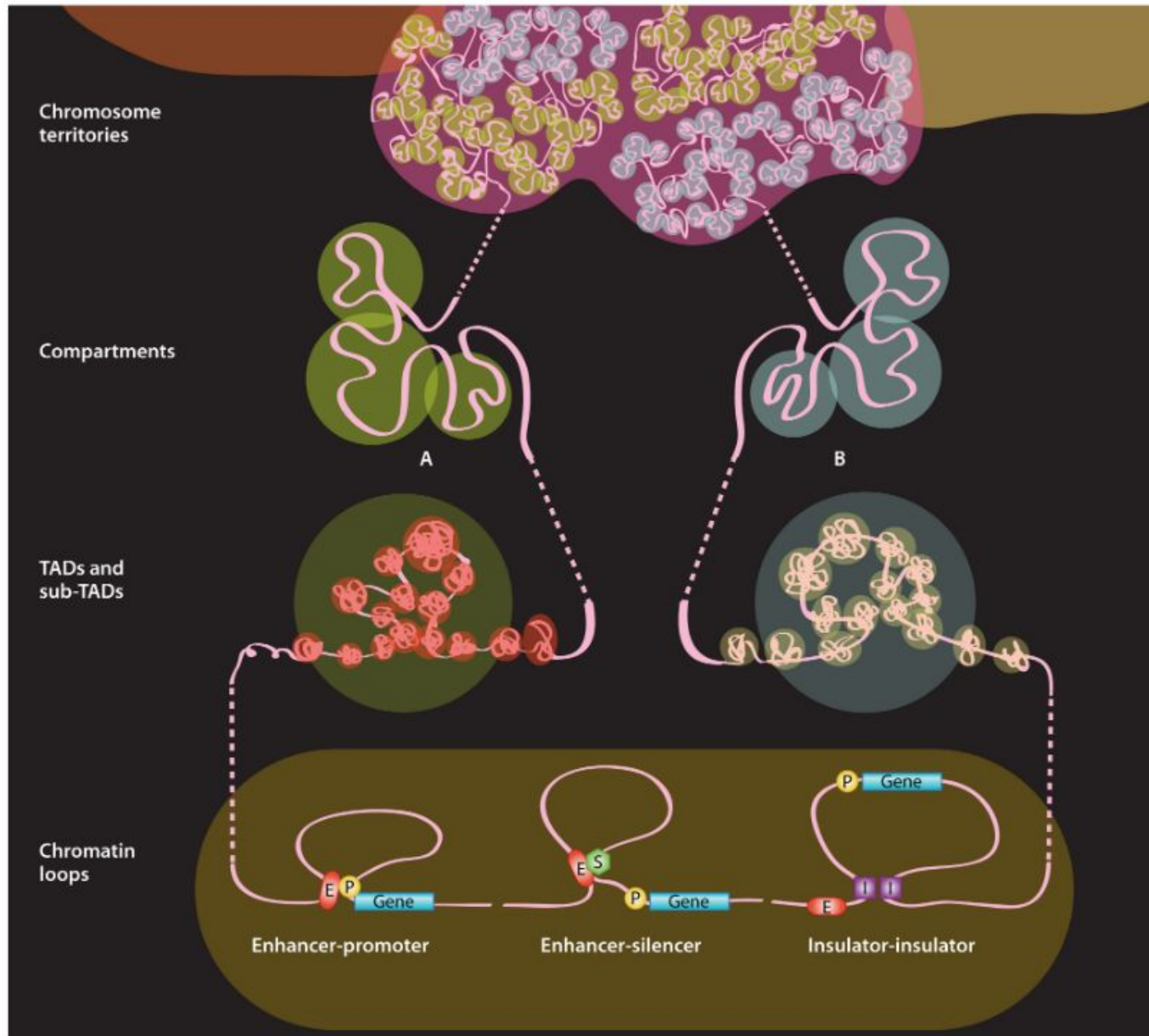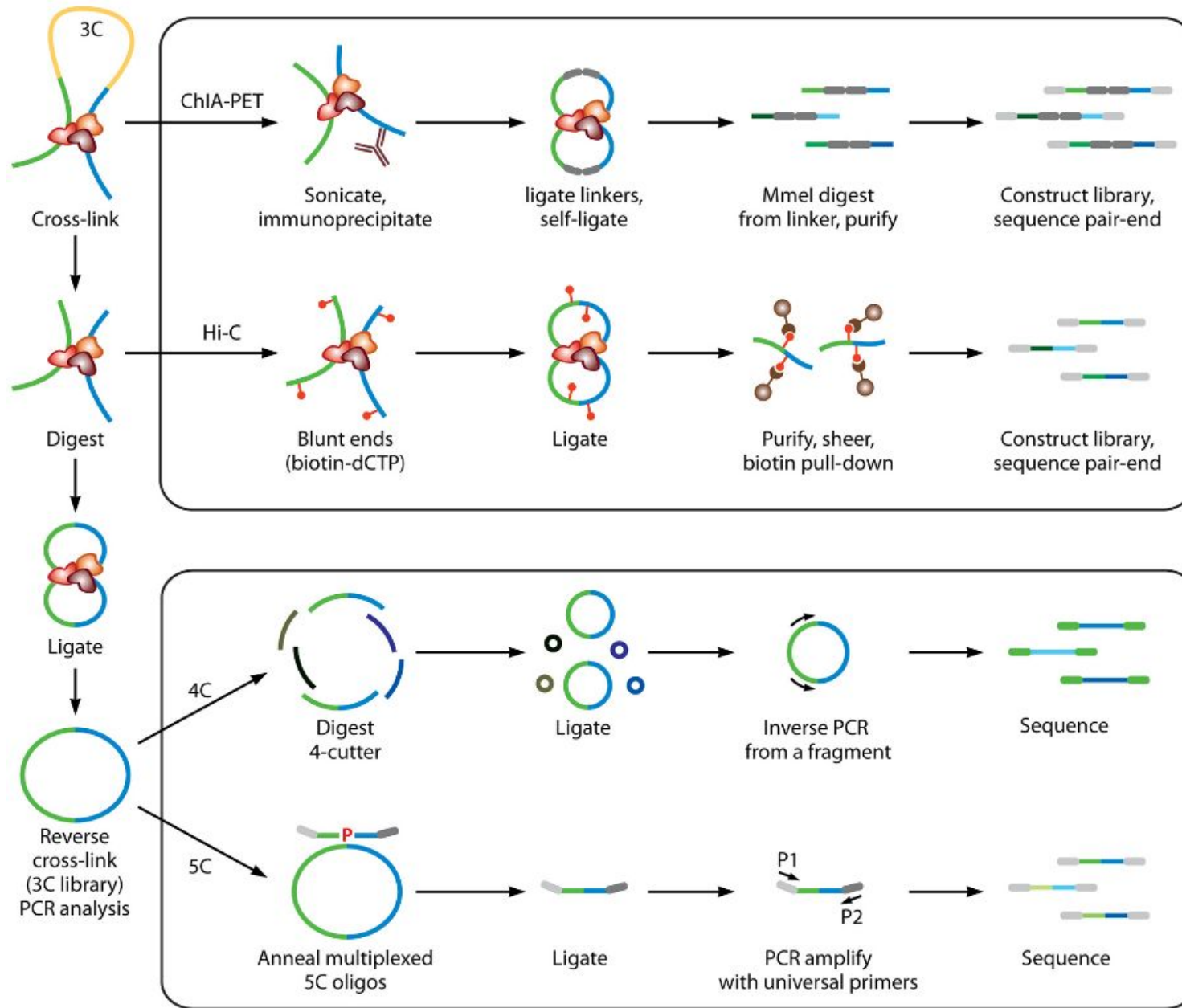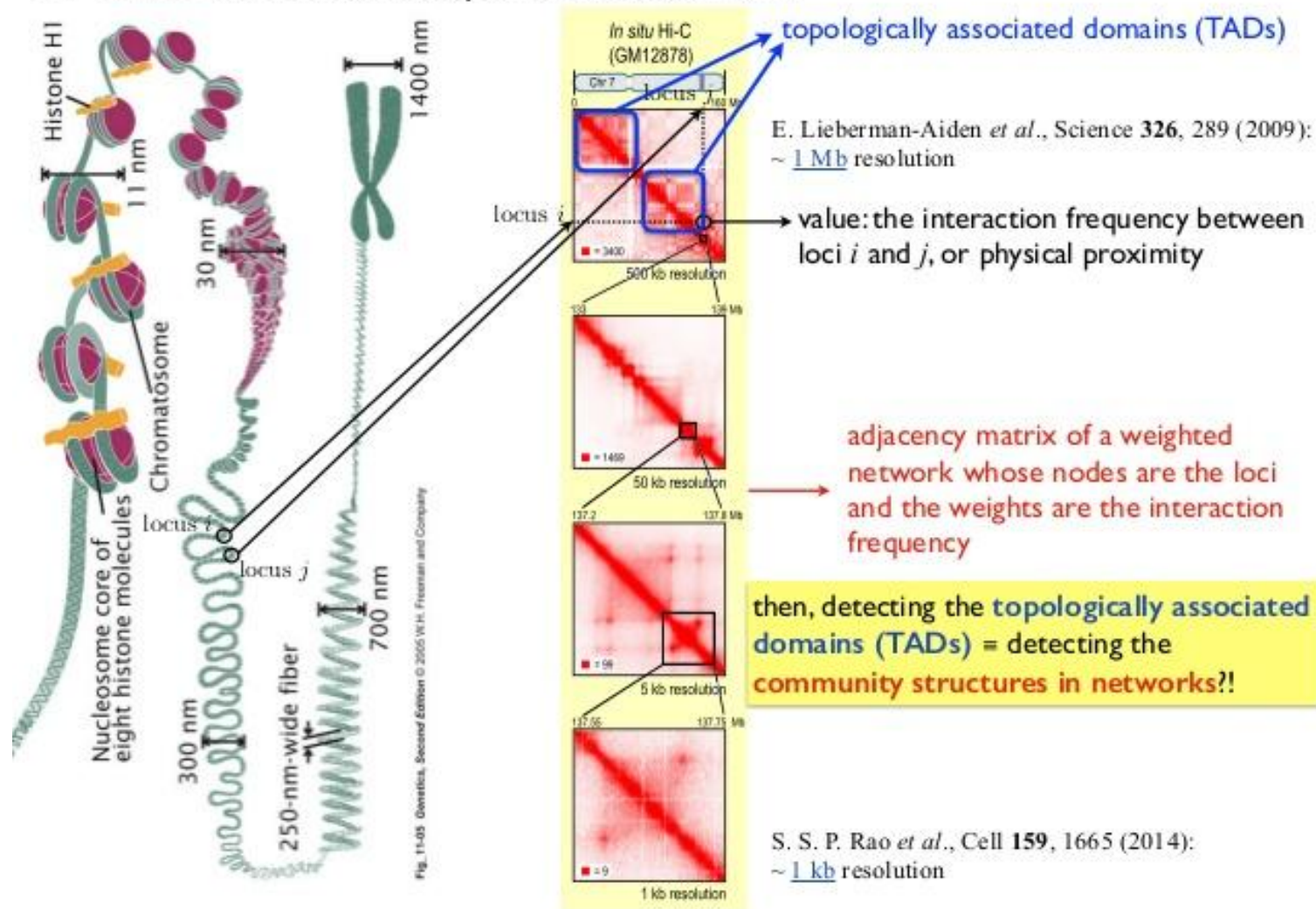
# Genome 3D organization

# 3D seq methods - 3C, HiC

# 3D seq methods - 3C, HiC



Hi-C: the interaction map of chromatin loci

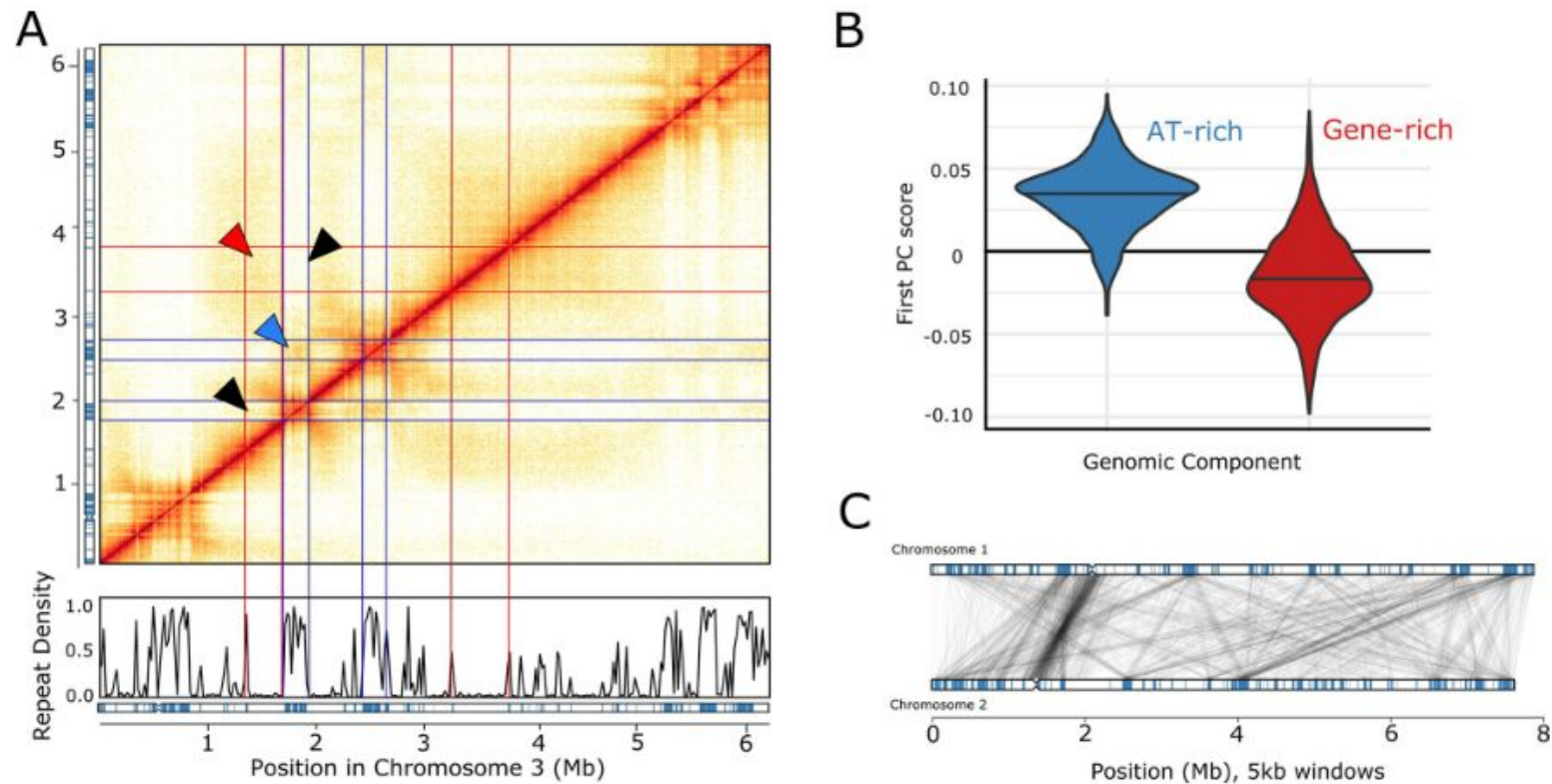# 3D Genome: do repeats matter?

# 3D Genome: do repeats matter?

## Repeat elements organise 3D genome structure and mediate transcription in the filamentous fungus *Epichloë festucae*

David J. Winter [1,2], Austen R. D. Ganley [3], Carolyn A. Young [4], Ivan Liachko [5], Christopher L. Schardl [6], Pierre-Yves Dupont [7], Daniel Berry [7], Arvina Ram [7], Barry Scott [2,7], Murray P. Cox [1,2]*

*"Our results reveal a genome in which very repeat-rich blocks of DNA with discrete boundaries are interspersed by gene-rich sequences that are almost repeat-free. In contrast to other species reported to date, the three-dimensional structure of the genome is anchored by these repeat blocks, which act to isolate transcription in neighbouring gene-rich regions. Genes that are differentially expressed in planta are enriched near the boundaries of these repeat-rich blocks, suggesting that their three-dimensional orientation partly encodes and regulates the symbiotic relationship formed by this organism."*

# 3D seq methods - 3C, HiC



Fig 4. Hi-C data reveals interactions within and among chromosomes. A. Each element of the matrix reflects the frequency of contacts between two genomic windows in an exemplar region of chromosome 3. Repeat density and the locations of AT-rich regions (blue) are plotted below the matrix, as in Fig 2. Red lines represent the boundaries of gene-rich regions, blue lines boundaries of AT-rich regions. Triangles highlight examples of interactions among specific genomic regions. Red, an interaction among gene-rich regions with high contact frequency (dense shading); blue, an AT-AT interaction with high contact-frequency (dense shading); black, interactions with low contact frequency between gene-rich and AT-rich regions (light shading). B. Distribution of first principal component scores estimated from Hi-C data for 5 kb regions entirely made up of AT-rich (blue) or gene-rich sequence (red). C. Inter-chromosomal contacts between chromosomes 1 and 2 are shown. All 5 kb windows sharing more than five Hi-C contacts are connected by a grey line. The AT-rich blocks in each chromosome are indicated (blue boxes).
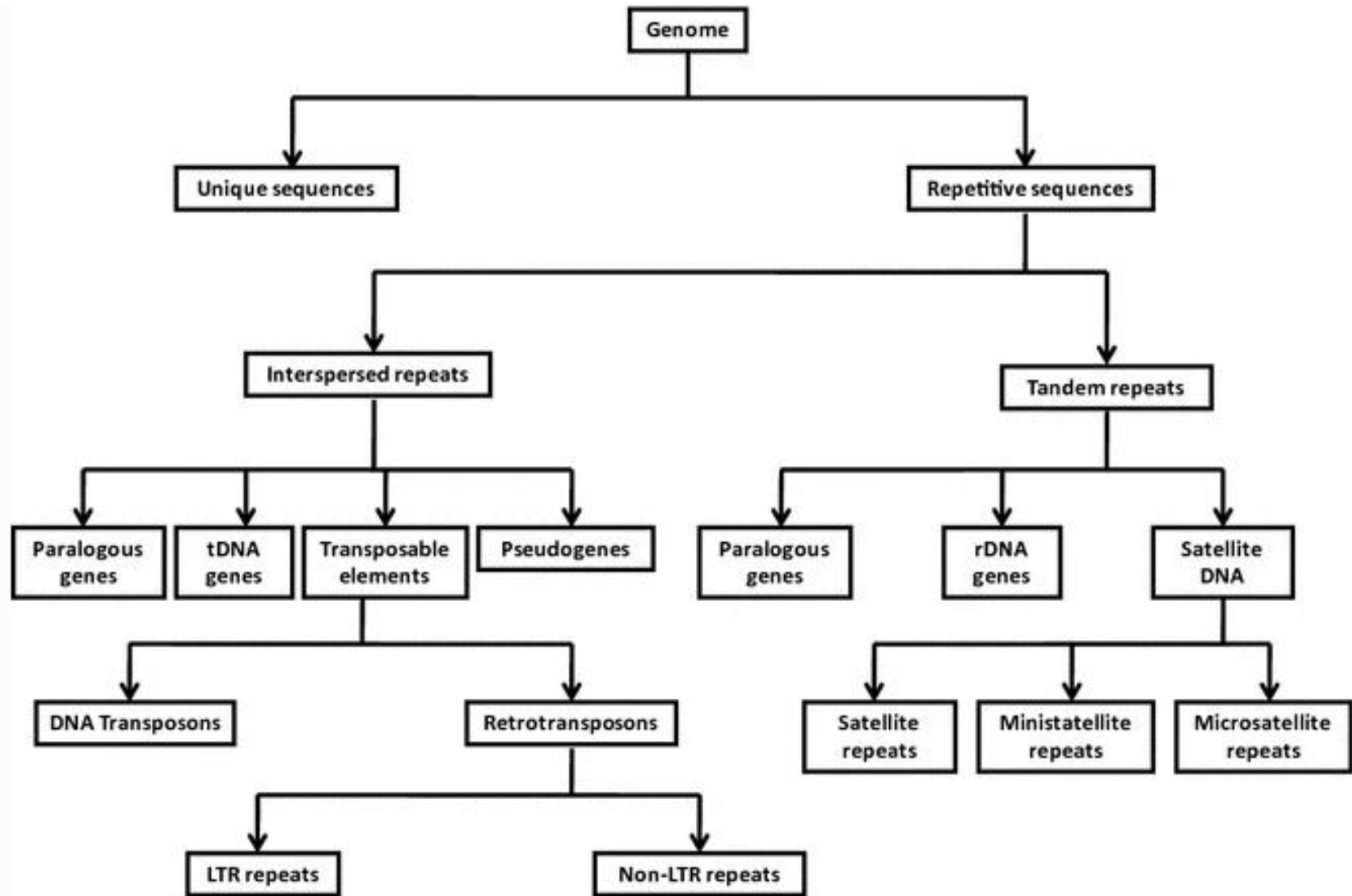
https://doi.org/10.1371/journal.pgen.1007467.g004

# Questions/hypotheses

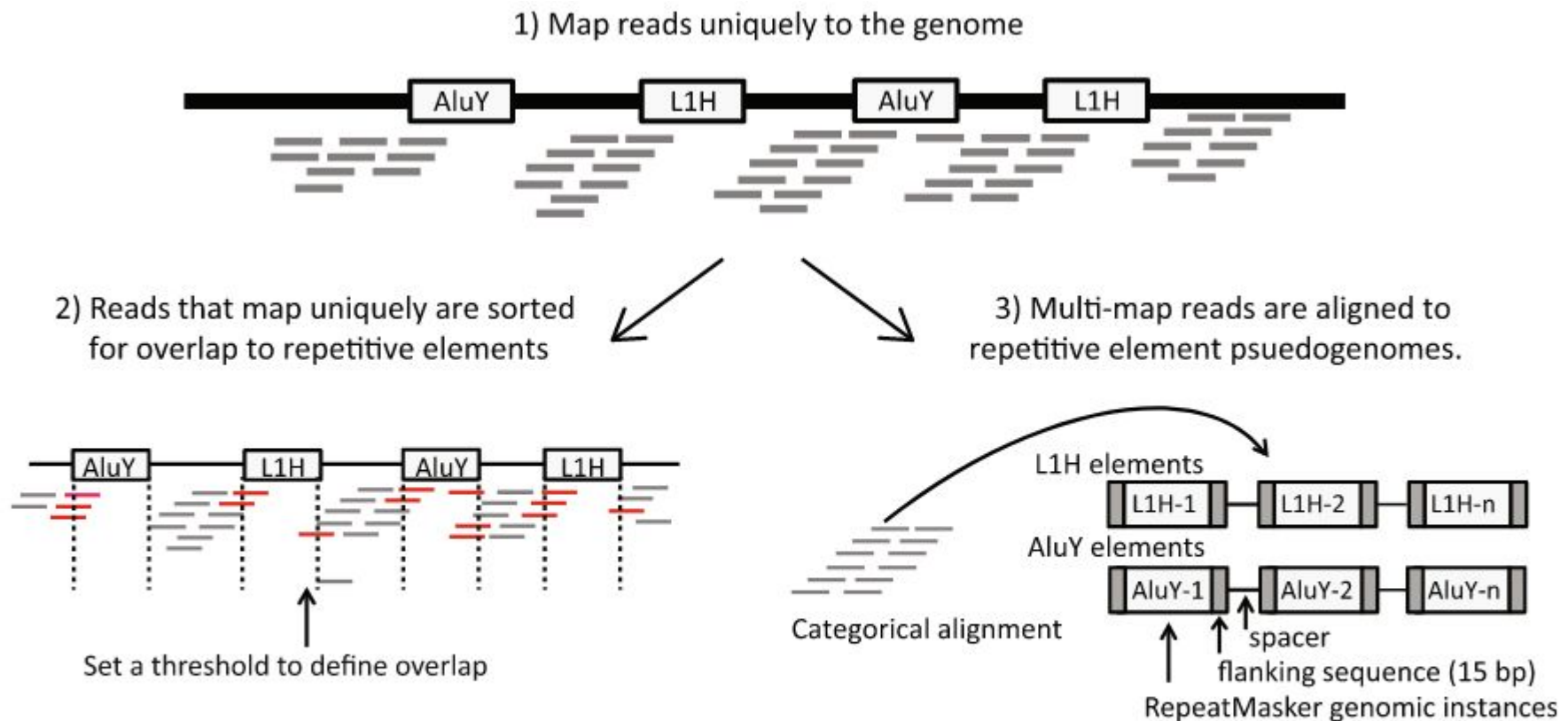Can repeat contribution to interphase genome 3D organization be seen in HiC data from plants?

Which families of repeats may have such function?

Where in the genome are these effects concentrated?

# Unique/multimapping reads

# Unique/multimapping reads



**Figure 1** *RepEnrich* **read mapping strategy.** Reads are mapped to the genome using the *Bowtie1* aligner. Reads mapping uniquely to the genome are assigned to subfamilies of repetitive elements based on their degree of overlap to *RepeatMasker* annotated genomic instances of each repetitive element subfamily. Reads mapping to multiple locations are separately mapped to repetitive element assemblies – referred to as repetitive element psuedogenomes – built from *RepeatMasker* annotated genomic instances of repetitive element subfamilies.

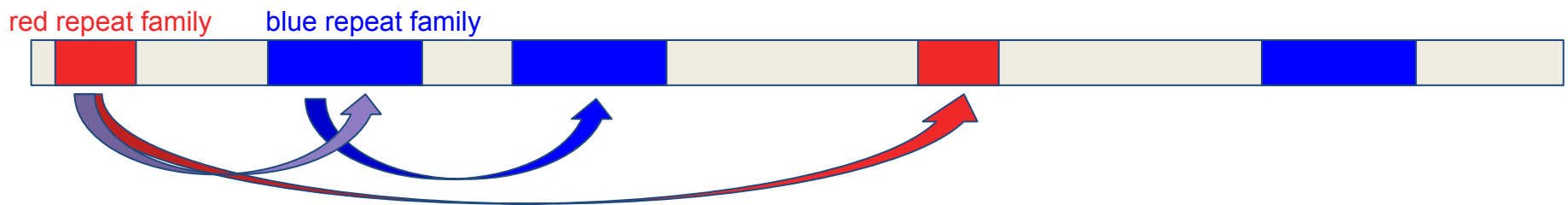Transcriptional landscape of repetitive elements in normal and cancer human cells

Steven W Criscione,[1] Yue Zhang,[1] William Thompson,[2,3] John M Sedivy,[1] and Nicola Neretti[1,3]

# Unique/multimapping reads



Cechova, 2020

# How to find contacts between repeat families?

- Mapping to the reference

red repeat family      blue repeat family



For red and blue repeat families, the Hi-C read pairs can be formed between:
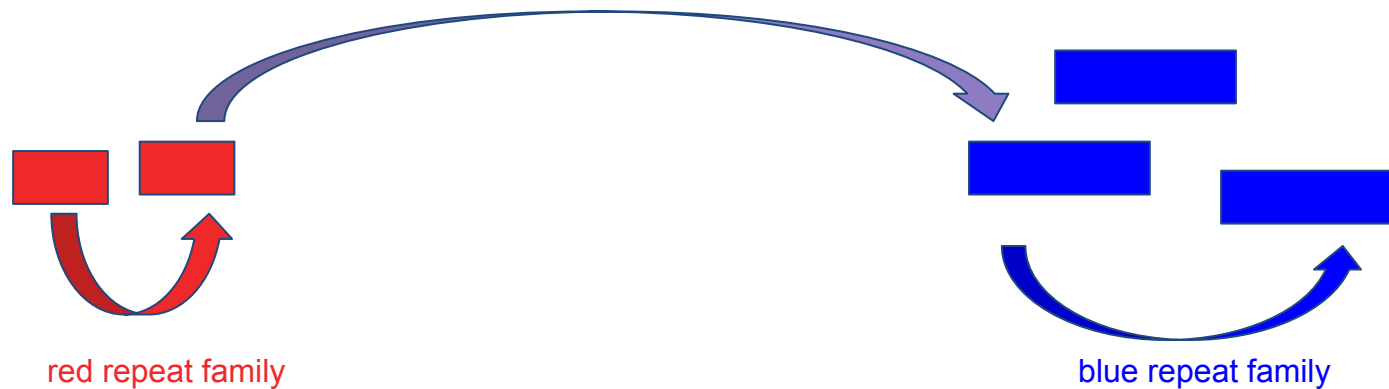- ➤ red and red family
- ➤ red and blue family
- ➤ blue and blue family

Mapping to the reference brings challenges:
- ❖ how to deal with multi-mapping reads?
- ❖ imperfect representation of repeats in reference genomes

# RepeatExplorer

- Find read pairs within and between clusters
- Label each cluster
- No reference necessary

red repeat family

blue repeat family

=> complementary method to the reference-based approach

# Do we see more contacts than expected?

Three methods: **observed versus expected counts of read pairs**

- joint probability
  - expected counts are p(A,B) = p(A).p(B)

- label permutation
  - each row represents a read pair with annotations of family1 and family2, family names are assigned to random rows of the table

- annotation shuffling
  - annotation of repeat families are shuffled along the reference genome

# HiC-TE pipeline

# HiC-TE pipeline

# HiC-TE pipeline

# HiC-TE pipeline

**Supplementary Table 2** - Hi-C tomato (*Solanum lycopersicum*) leaf mesophyll sequencing runs from project SRP110225 (Dong et al., 2017) used to test the Nextflow pipeline. The individual runs represent different biological and technical replicates (see batch and plant numbers).

| Sample_name | Batch | Plant | Gbp | SRA ID |
|---|---|---|---|---|
| SlMC_HiC_1.1.1 | 1 | 1 | 49.34 | SRR5748725 |
| SlMC_HiC_1.1.2 | 2 | 1 | 38.99 | SRR5748726 |
| SlMC_HiC_1.2.1 | 1 | 3 | 27.34 | SRR5748729 |
| SlMC_HiC_1.2.2 | 2 | 3 | 31.31 | SRR5748730 |
| SlMC_HiC_2.2.1 | 1 | 4 | 23.88 | SRR5748733 |
| SlMC_HiC_2.2.2 | 2 | 4 | 13.56 | SRR5748734 |

# HiC-TE pipeline

**Supplementary Table 1** - HiC-TE Nextflow pipeline performance on a 4-core 3.0GHz Intel Ubuntu box and in the cloud (MetaCentrum metacentrum.cz). Numerical values are averages of 12 runs excluding TE-greedy-nester reference annotation (is needed only once).

| Hardware platform | Time | RAM | Temp.files | Output |
|---|---|---|---|---|
| Linux 4-core Intel 3.0GHz | 18h | 12GB | 120GB | 35GB |
| Linux MetaCentrum 8CPUs | 6h | 240GB | 240GB | 35GB |

# HiC-TE pipeline output

## hic-te

A Nextflow workflow to analyze HiC data from SRA (NCBI Short Read Archive) for 3D contacts between repeat families. Slightly biased (but not limited to) towards LTR retroTEs and plant genomes.

*SYNOPSIS*

nextflow run [FILEBASE].nf -profile LIST[,LIST...] [PARAMS...]

To run the pipeline, the following parameters are mandatory:

*DATA*

**reads**
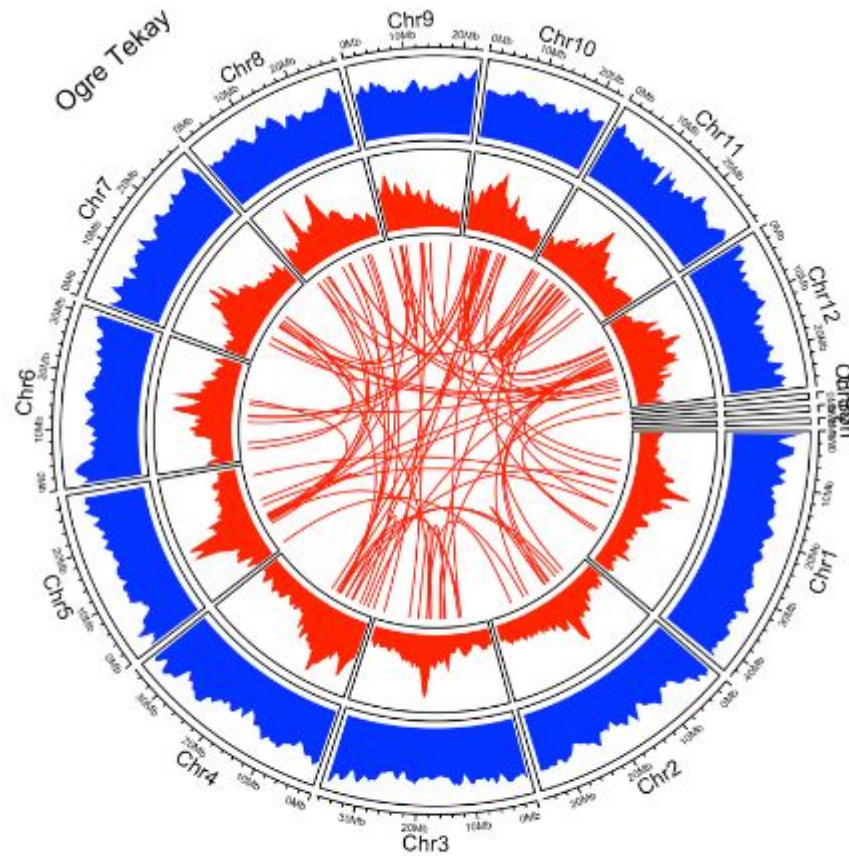Reads from a Hi-C experiment. The easiest way to provide this is by listing their SRA id.
--sra_run SRR14458670

**reference**
Reference genome corresponding to the organism the reads belong to. The reference should be in the fasta format.
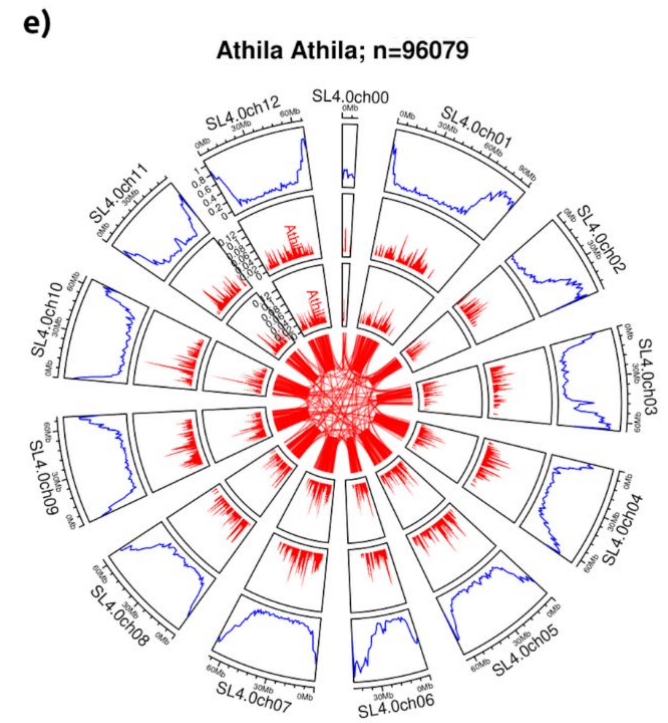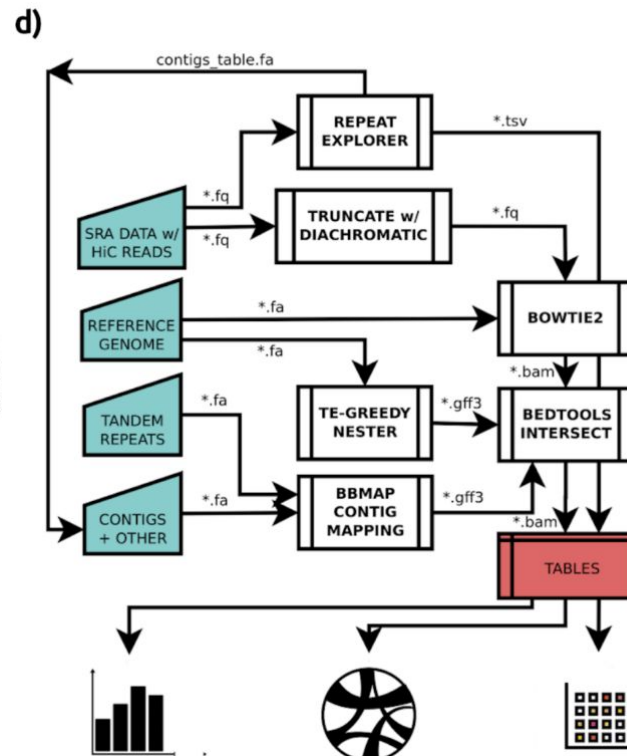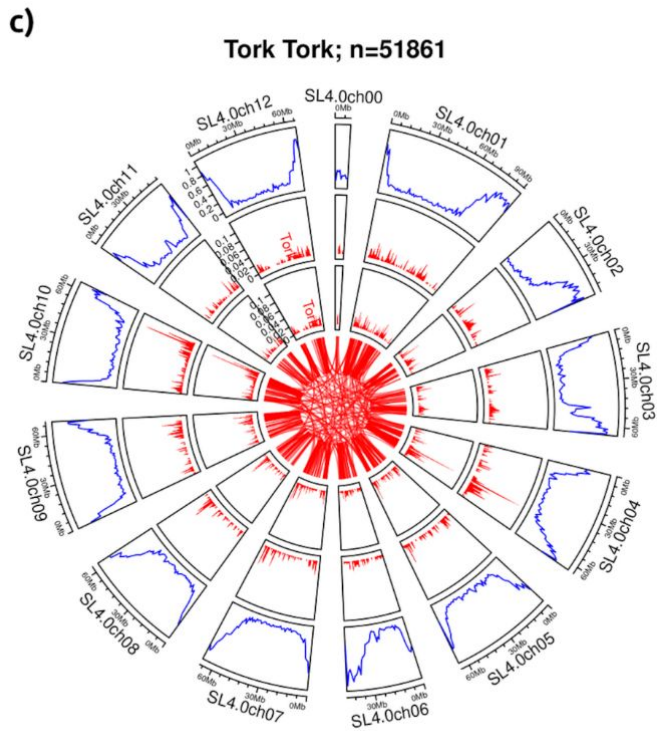--reference Athaliana_167_TAIR10.fa

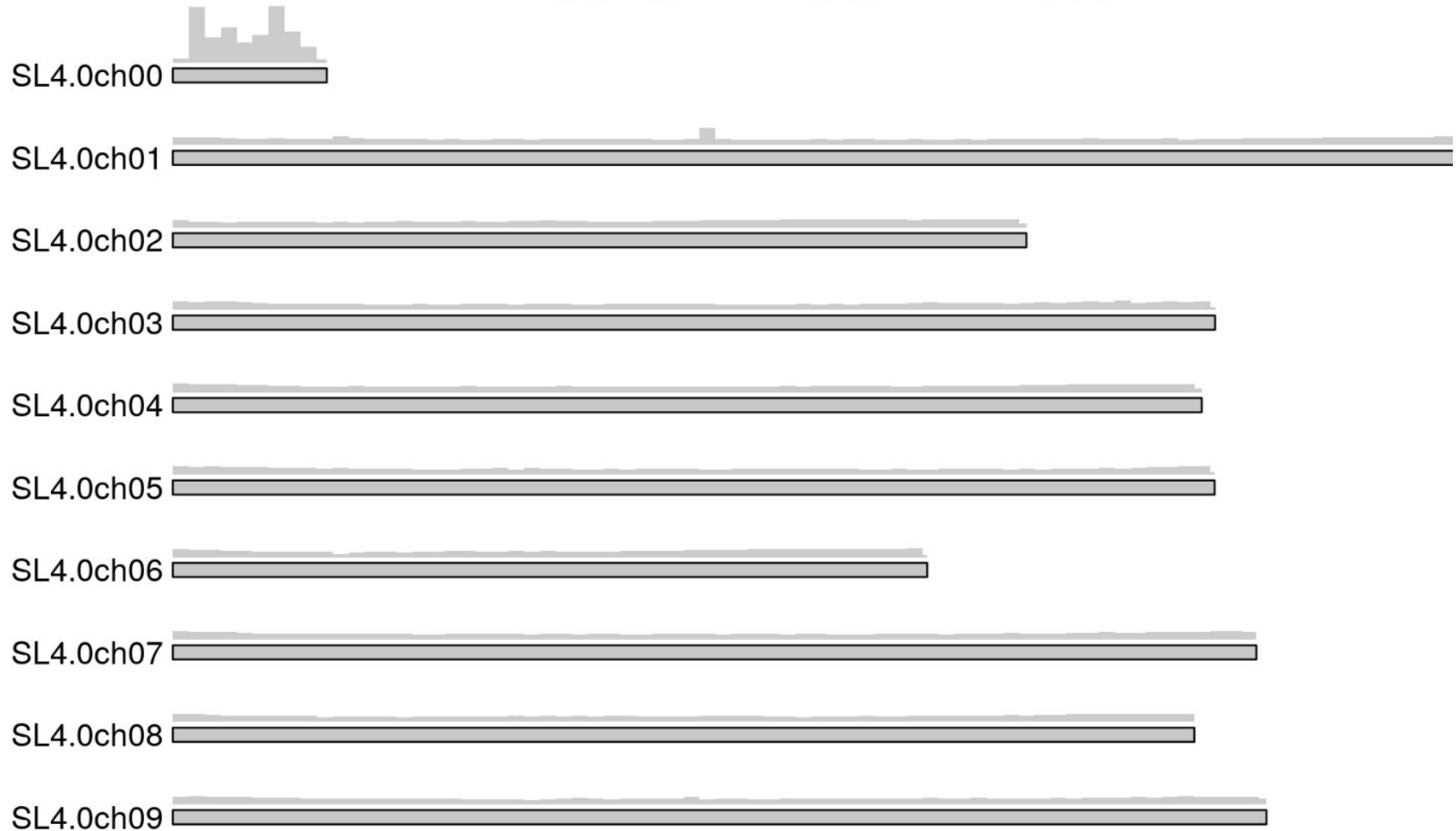# HiC-TE pipeline output



genes
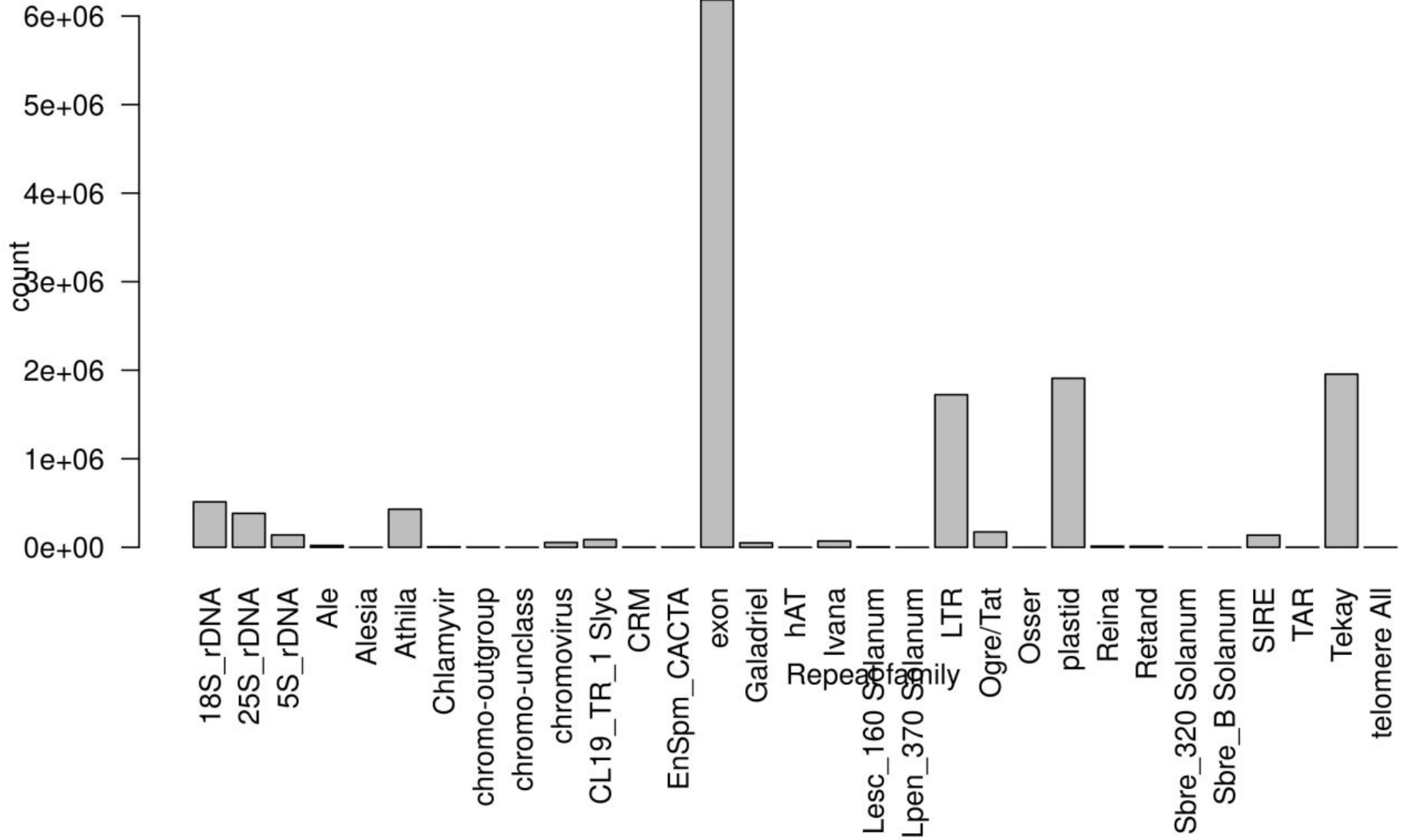repeats; TE + (micro)satellites

# HiC-TE pipeline output

# HiC-TE pipeline output
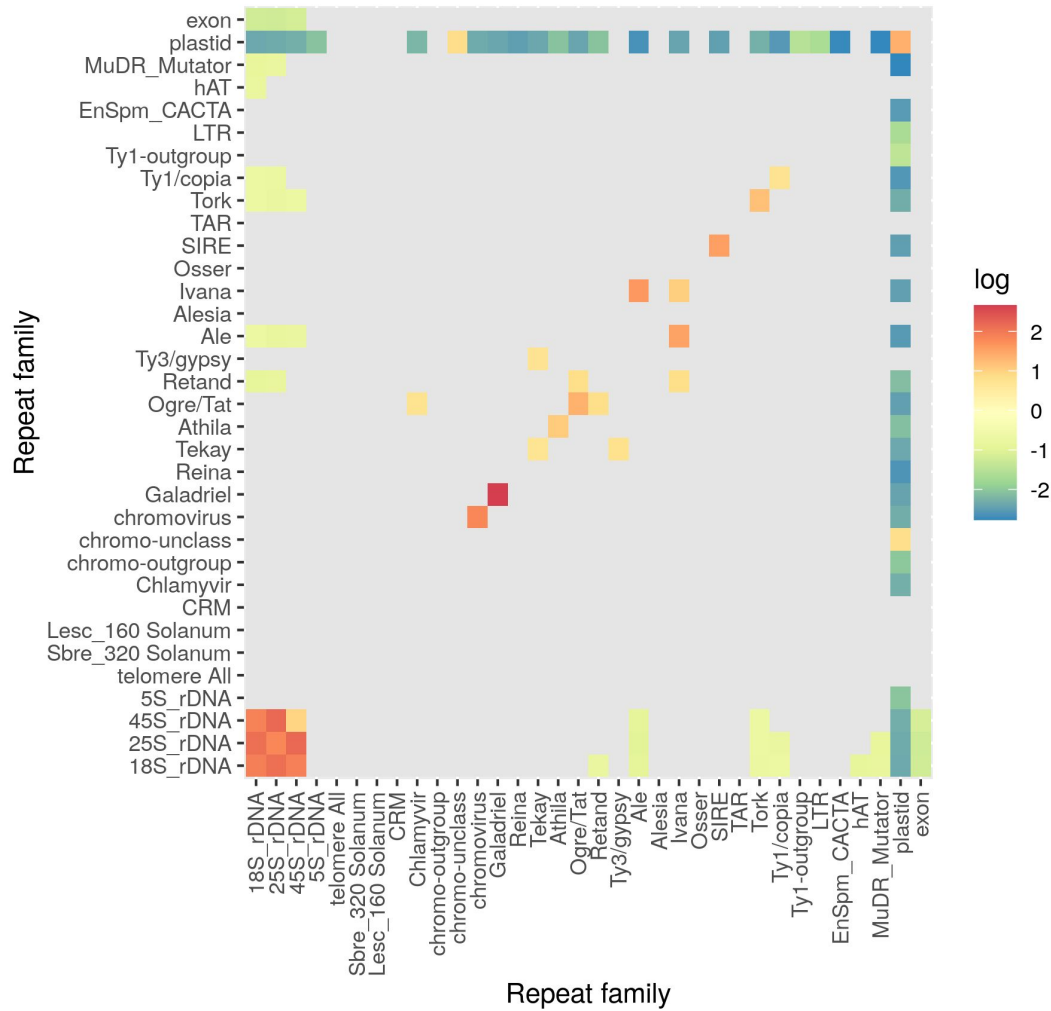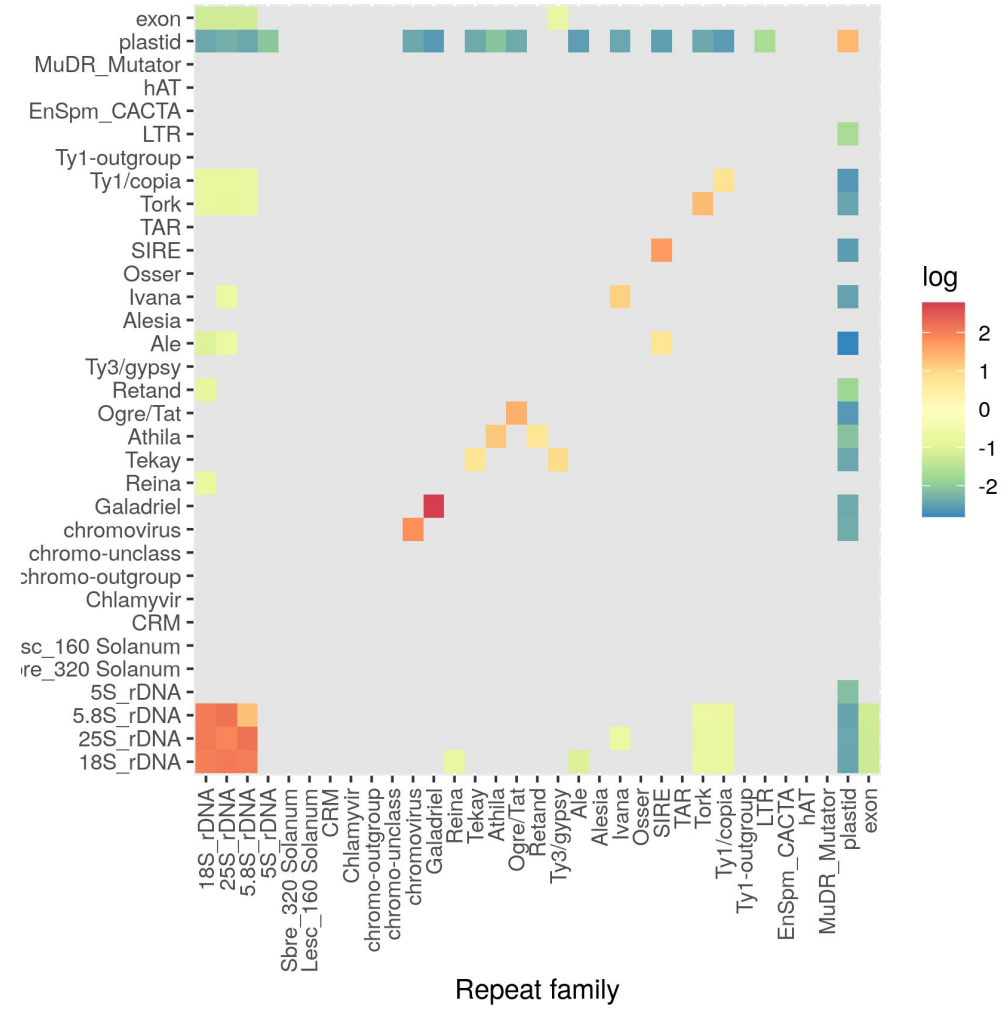


Slyc_SL40_SRR5748725_png/plotBamCoverage.png
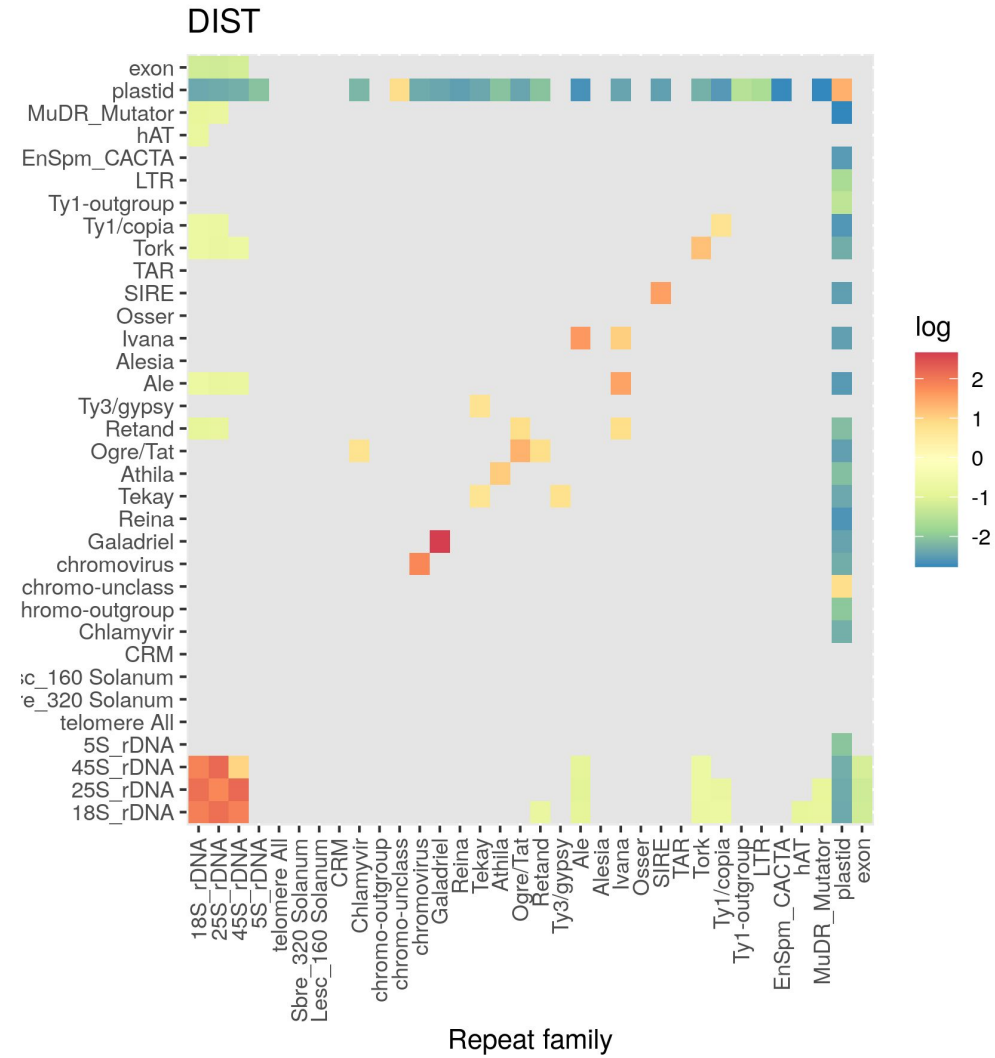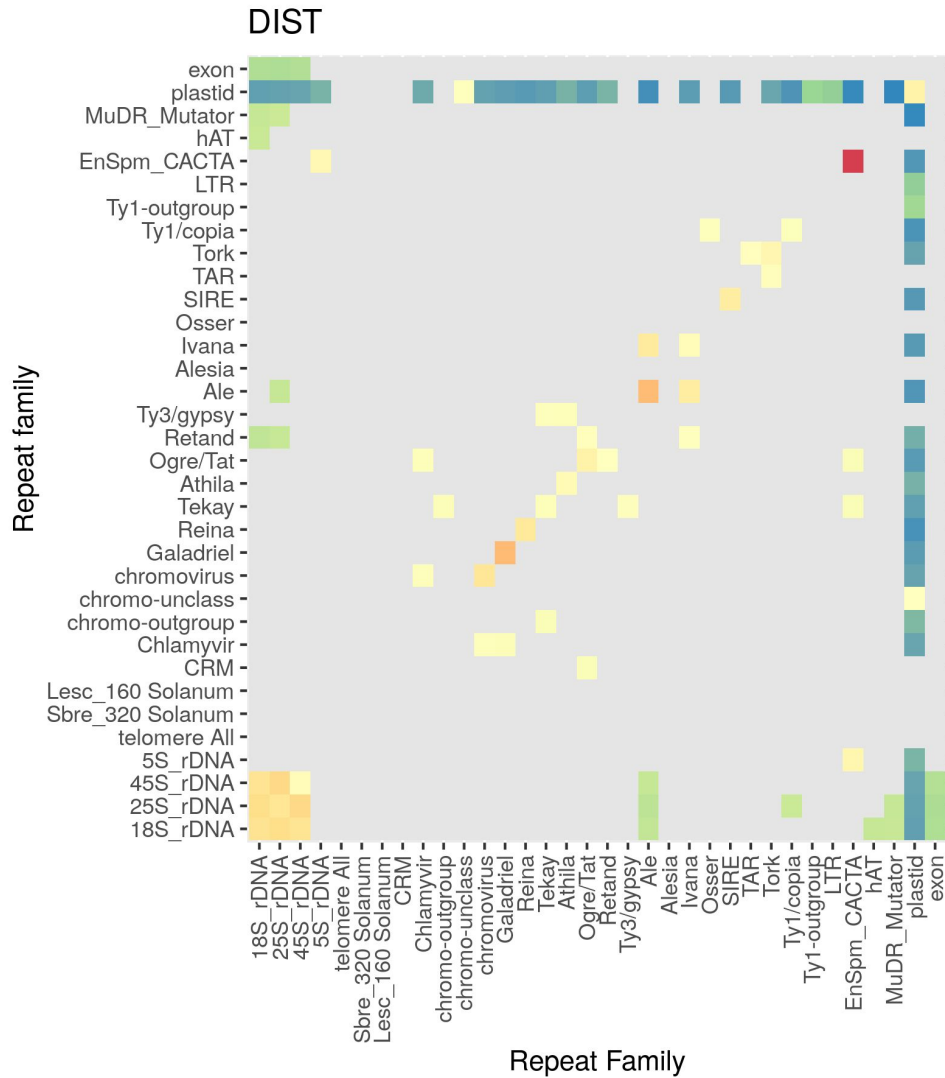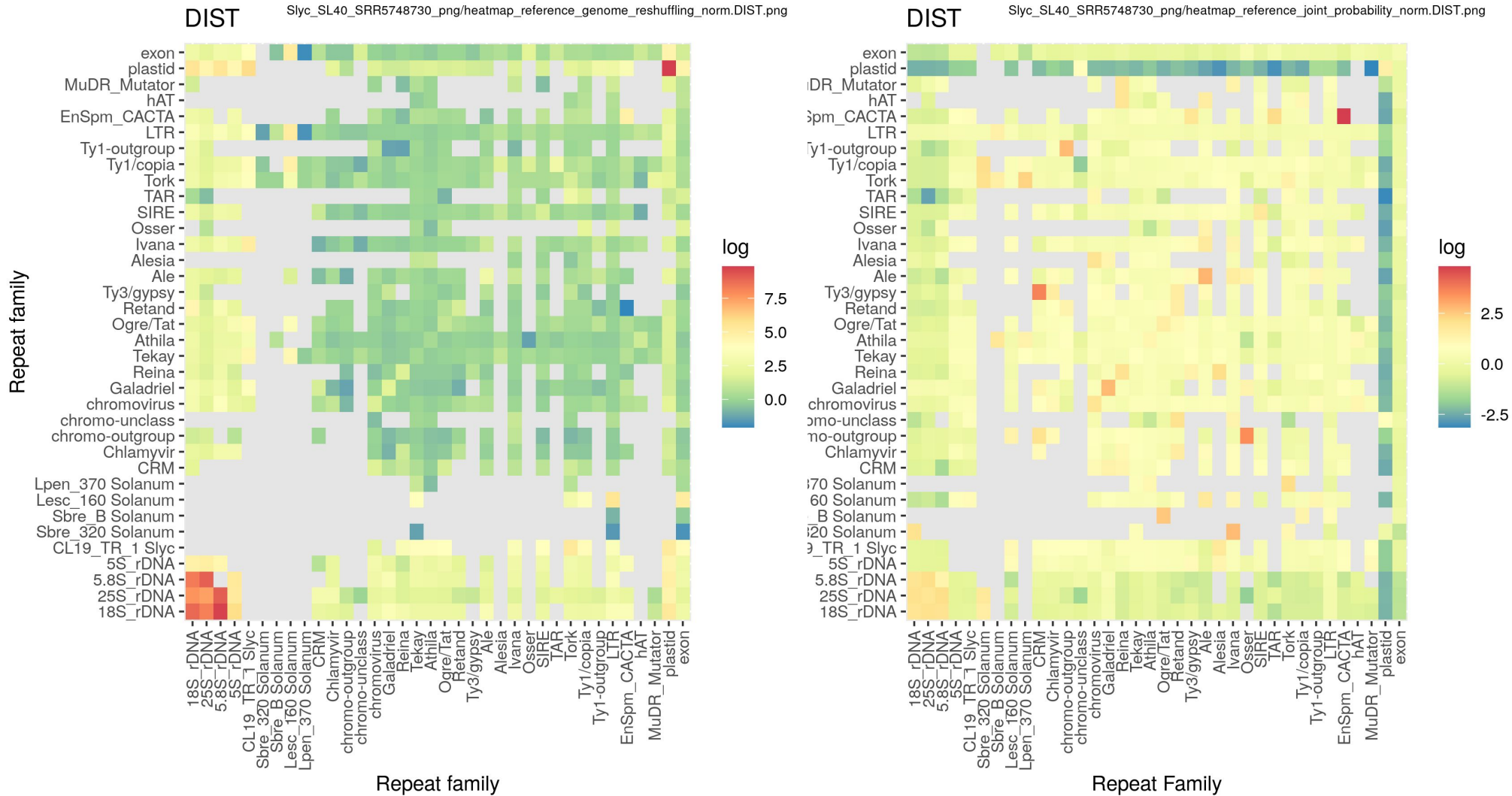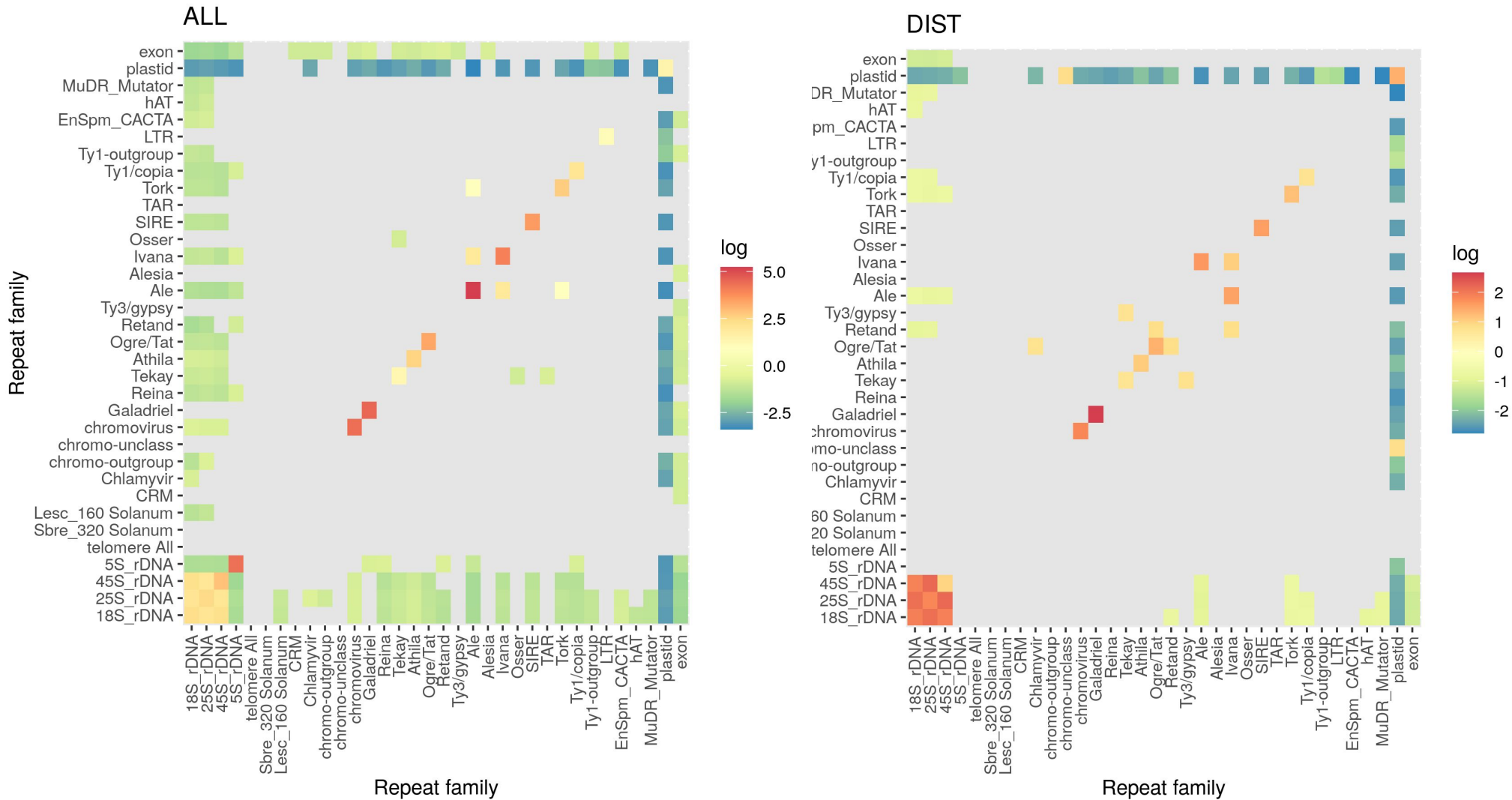
# HiC-TE pipeline output

# Different plant individual
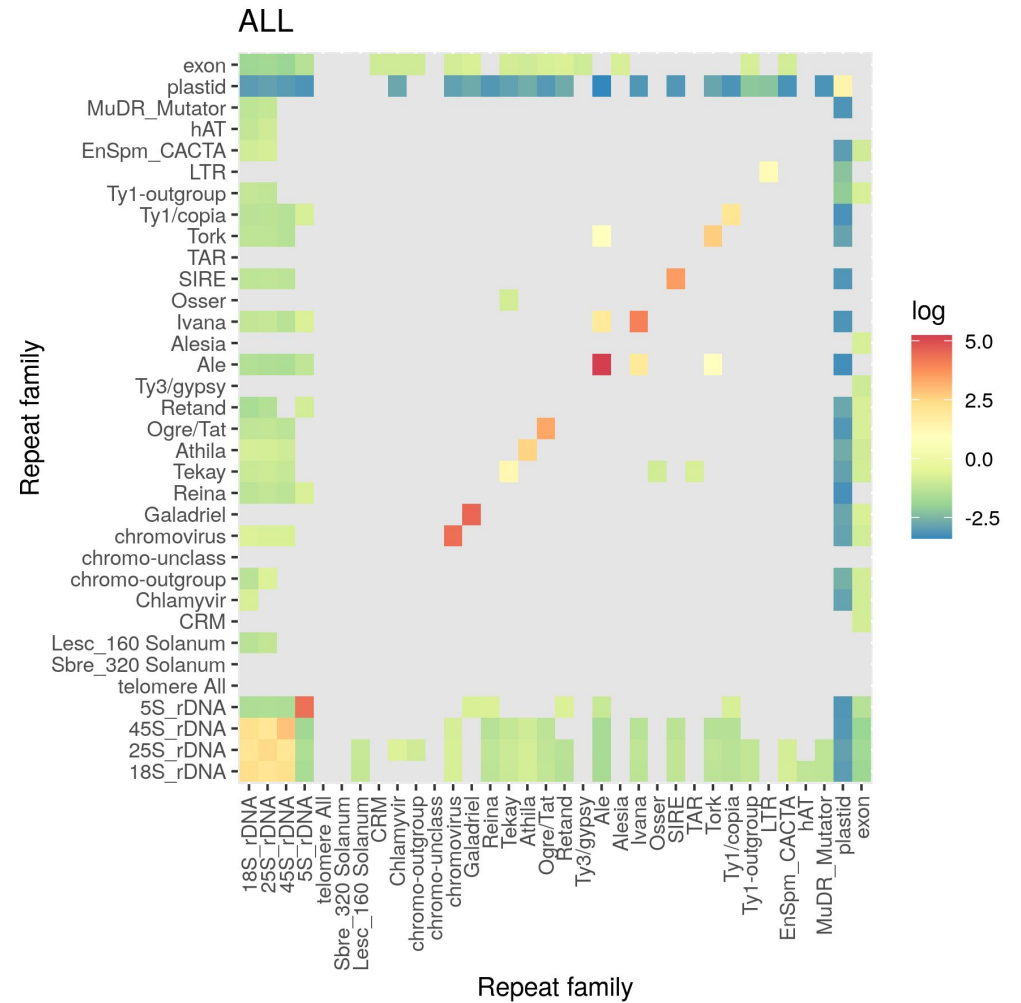
# Different normalization method

# Different normalization method

# All vs long-distance interactions



ALL

DIST

# RepeatExplorer versus reference

# Method matter (when dealing with repeats)

- General patterns stay the same but the individual variation is present depending on the individual, normalization method, repeat dispersion, density etc.

- RepeatExplorer clusters offer unique insights into repeat contacts not available with reference-based methods

- Hi-C TE pipeline can be used to uncover chromatin contacts between repetitive annotations in plant genomes

# HiC-TE pipeline contributors

Matej Lexa > hic-te

**H**  **hic-te** ⊕
Project ID: 18243 ☐

☆ Star | 1

○ **236 Commits**  ⅄ **2 Branches**  ⊘ **0 Tags**  ▤ **380.2 MB** Files  ▤ **380.3 MB** Storage

A workflow to analyze HiC data from SRA for 3D contacts between TE families.

| master ∨ | hic-te | History | Find file | ⬇ ∨ | Clone ∨ |

**Update Diachromatic.jar**
Son Hoang Nguyen authored 2 weeks ago

| Verified | 3bcaa3d9 ☐ |

▤ README   ⚖ MIT License

|  | Last commit | Last update |
|---|---|---|
| ▭ Docker/rtools | Add files for building docker image for rtools | 2 weeks ago |
| ▭ bin | Update Diachromatic.jar | 2 weeks ago |
| ▭ conf | uploading the test data for the arabidopsis… | 3 months ago |
| ▭ data | uploading the test data for the arabidopsis… | 3 months ago |
| ▭ doc | Updated source for hic-te flow diagram for… | 4 months ago |
| ▭ modules | Modified the script for the metacentrum e… | 3 months ago |
| ▭ tmp | Add singularity_runOption parameter | 4 months ago |

# HiC-TE pipeline contributors

ⓘ Matej Lexa, ⓘ Monika Cechova, Son Hoang Nguyen, ⓘ Pavel Jedlicka, ⓘ Viktor Tokan, Zdenek Kubat, ⓘ Roman Hobza, ⓘ Eduard Kejnovsky

💬 0 | ☑ 0 | 👥 0 | ⚙ 0 | 🖥 0 | 🎞 0 | 🐦 30

**Abstract**  Full Text  Info/History  Metrics  🗋 Preview PDF

## Abstract

The role of repetitive DNA in the 3D organization of the interphase nucleus in plant cells is a subject of intensive study. High-throughput chromosome conformation capture (Hi-C) is a sequencing-based method detecting the proximity of DNA segments in nuclei. We combined Hi-C data, plant reference genome data and tools for the characterization of genomic repeats to build a Nextflow pipeline identifying and quantifying the contacts of specific repeats revealing the preferential homotypic interactions of ribosomal DNA, DNA transposons and some LTR retrotransposon families. We provide a novel way to analyze the organization of repetitive elements in the 3D nucleus.

# Thank you for your attention

**INSTITUTE OF BIOPHYSICS**

Pavel Jedlička (@pj_naruto)
Viktor Tokan
Zdeněk Kubát
Roman Hobza
**Eduard Kejnovský**

**MASARYK UNIVERSITY**

**Matej Lexa** (@matej_lexa)
Monika Čechová (@biomonika)
Son Hoang Nguyen