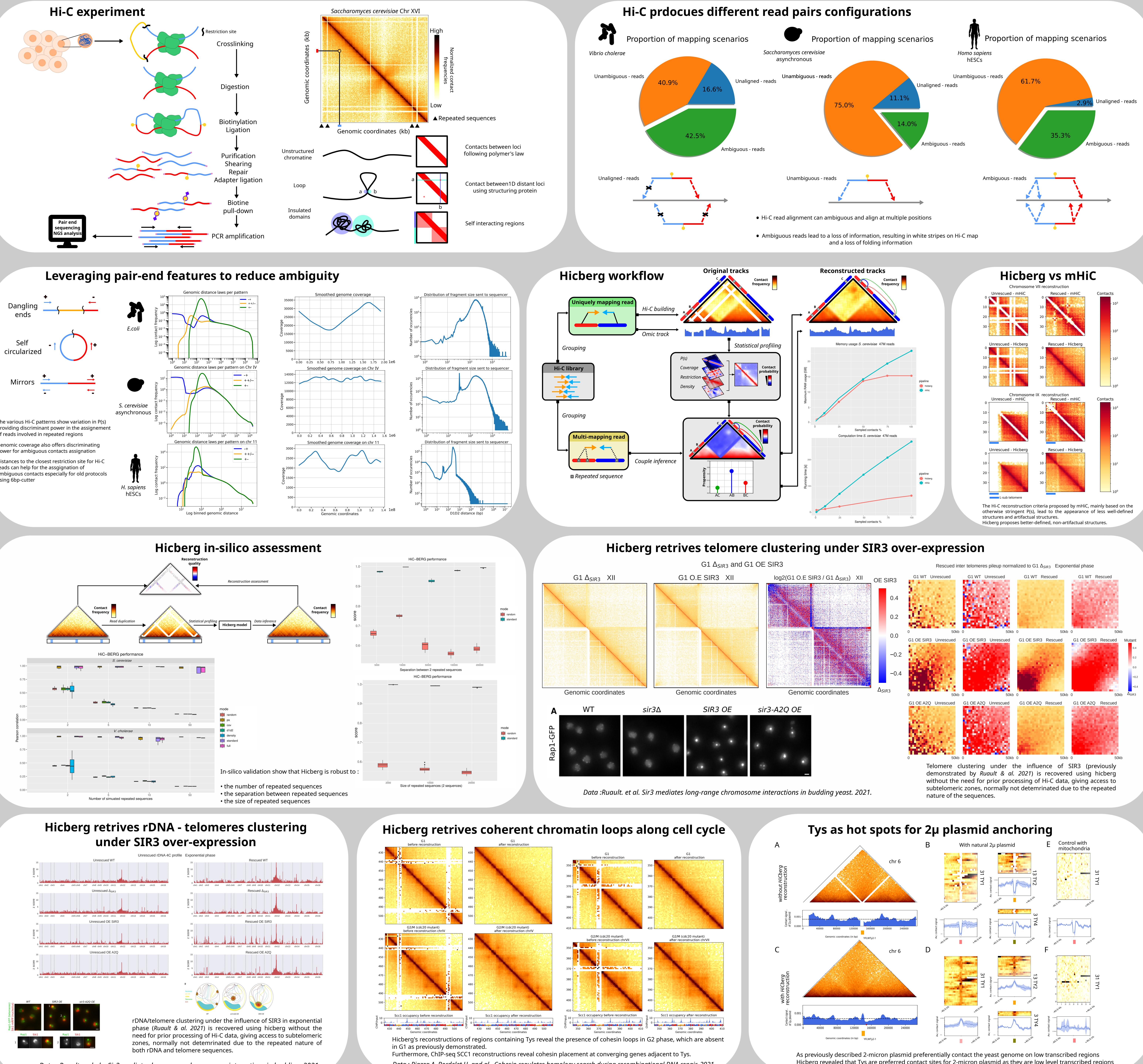


During their evolution, the genomes of micro-organisms can acquire quantities of different repeated elements such as retrotransposons, duplicated genes or tandem repeats. This type of sequence within genomes cannot be processed directly by NGS technologies because they generate short reads that cannot be located unambiguously on reference genomes. This information is filtered out by most current pipelines, leading to incomplete genomic tracks resulting in a significant loss of information on biological functions, processes and genomic structures involving repeated elements. To overcome these limitations we developed Hicberg, an algorithm that uses statistical inference and pseudo-random generators to predict the positions of repeated sequences' reads from different omics paired-end data (including Hi-C, Mnase-seq, ChIP-seq, ...), based multiple components relative to the polymeric behavior of the DNA and sequencing protocols' features, established on the non unambiguous part of the tracks. Thus read pairs belonging to repeated sequences can be assigned with robust confidence in genomes filling-in genomic tracks. After development and calibration on a controlled test bench Hicberg improves genomic data interpretability of various species, starting with microbial one such as *Saccharomyces cerevisiae*.

Reconstructions of Hi-C and ChIP-seq genomic tracks with Hicberg revealed how some retrotransposons in this model contribute to the positioning of cohesin, a molecular motor involved in the formation of chromatin loops. A new role of retrotransposon sequences as contact hot points for the elusive yeast 2 micron episomal molecule was also identified. Overall, these results underline the power of the approach to discover new novel molecular relationships, and the interest in applying this tool more widely to larger genomes with greater quantity of repeats. The proposed method can therefore provide an alternative visualization of genomic signals in a wide variety of biological conditions and allow a more comprehensive view of genome organization and plasticity. Importantly, existing dataset can be revisited using this approach to unveil overlooked features.



Conclusion & Perspectives

Hicberg paves the way to the exploration of functional 3D organization of structures involving repeated elements such as *V. cholerae* super integron, *S. cerevisiae* transposable elements, *P. falciparum* telomeric regions, etc.

Currently applied on genomes of ~ 100Mb in size carrying variable repeated sequences involved in pathogenicity (*V. dahliae*, *E. festucae*)

Hicberg allows working with more complete set of data that will also improve the quality of the already visible part of genomes giving more evenly distributed signals at the genome scale

Hicberg may improve Hi-C normalization procedure and specific patterns detections such as peaks (1D signals) and loops/domains (Hi-C).

References

- [1] Dekker J, Rippe K, Dekker M, Kleckner N. Capturing chromosome conformation. *Science*. 2002 Feb 15;295[5558]:1306–11.
- [2] Lieberman-Aiden E, van Berkum NL, Williams L, Imakaev M, Ragoczy T, et al. Comprehensive mapping of long-range interactions reveals folding principles of the human genome. *Science*. 2009 Oct 9;326[5950]:289–93.
- [3] Cournac A, Koszul R, Mozziconacci J. The 3D folding of metazoan genomes correlates with the association of similar repetitive elements. *Nucleic acids research*. 2016;44[1]:245–255.
- [4] Forcato M, Nicoletti C, Pal K, Livi CM, Ferrari F, Bicciato S. Comparison of computational methods for Hi-C data analysis. *Nature methods*. 2017;14[7]:679.
- [5] Cournac A, Marbouy M, Mozziconacci J, Koszul R. Generation and analysis of chromosomal contact maps of yeast species. *Yeast Functional Genomics: Methods and Protocols*. 2016;227–245.