

Comparison between ChromStruct4 and TADbit

Flagship Project InterOmics – WP1 CNR-ISTI

Claudia Caudai¹, Emanuele Salerno¹, Monica Zoppè², and Anna Tonazzini¹

National Research Council of Italy

¹Institute of Information Science and Technologies, and ²Institute of Clinical Physiology, Pisa, Italy

`claudia.caudai@isti.cnr.it`

Abstract In this report performances of ChromStruct4 and TADbit have been compared. These are two methods for the inference of chromatin three-dimensional conformations starting from Chromosome Conformation Capture data. TADbit and ChromStruct4 have been tested against the same data sets from real HI-C experiments. With comparative experiments, also robustness of ChromStruct4 against biases have been tested.

1 Introduction

TADbit [1] is a very popular method with good performances and with reasonable calculation times. It takes in input the HI-C contact data [2] from Illumina sequencing. This method removes biases through the ICE algorithm [3], produces normalized contact matrices and identifies topological association domains (TADs) [4]. TADbit models the chromatin fibre as a set of non-connected consecutive beads; all the beads have the same radius and are subject to constraints. The method identifies, in an empiric way, three different distances of balance: for consecutive beads, for couples of beads whose contact frequencies exceed a certain threshold (they are often in contact so they are supposed to be close to each other) and for couples of beads whose contact frequencies do not exceed a certain threshold (they are rarely in contact, so they are supposed to be far from each other). TADbit produces many three-dimensional conformations by optimizing harmonic functions associated with Euclidean distances between pairs of beads, using as equilibrium distances the constraints previously described. Finally, the method selects configurations best fitting the input data through a Monte Carlo method.

ChromStruct4 [5, 6] takes in input HI-C contact matrices. The method doesn't clean up data because it uses a data filtering method which makes it robust against biases. ChromStruct4 identifies topological domains (TADs) using its own algorithm and models the chromatin fibre as a bead-chain in which consecutive beads are connected to each other. The system is treated as a kinematic

chain and motion is performed by quaternions. HI-C data derive from millions of cells, so ChromStruct4 is aimed to sample the space of possible solutions. At every run, this method produces a three-dimensional chromatin conformation through a Simulated Annealing method, adopting a scoring function oriented to rewarding fit with input data and discouraging geometrical constraint violations.

2 Tests against Hi-C real data

Both methods have been tested against the following data sets:

- Long arm of human Chromosome 1: from 150.3*Mb* to 179.4*Mb* (binned at 100*kb*) [GEO: GSE18199] [2]
- *Caulobacter Crescentus* CB15 (binned at 10*kb*) [GEO: GSE45966] [3]

TADbit takes in input fastq sequencing data and produces contact matrices, both raw and normalized with ICE. The raw matrix produced by TADbit for the long arm of Chromosome 1 is different from the contact matrix of Lieberman-Aiden [2], even though they both derive from the same HI-C data, also the raw *Caulobacter* matrix produced by TADbit is different from the contact matrix of Imakaev [3], even though same HI-C data have been used.

Observing the heatmaps of contact matrices produced by TADbit (Figures 2-3 and Figures 5-6) and values inside contact matrices (Figures 7-10), it is noteworthy that contact frequencies are different before and after normalization with ICE (two or three orders of magnitude), but the differences in heatmaps are small. Probably, the heatmap MICROSOFT EXCEL's algorithm attributes to raw and normalized values the same color.

We also observe that:

- Matrices produced by TADbit for Chromosome 1 are more concentrated around the main diagonal than those of Lieberman-Aiden
- Matrices produced by TADbit for the *Caulobacter Crescentus* have many rows and columns full of zeroes.

Probably, the differences between data produced by TADbit and Lieberman-Aiden and between data produced by TADbit and Imakaev depend on differences in the methods used in building the contact matrices from sequencing data.

3 Comparison of conformations produced by TADbit and ChromStruct4

In order to compare the two methods, the contact matrices produced by TADbit, both normalized and raw, have been used. Aims of this comparison are:

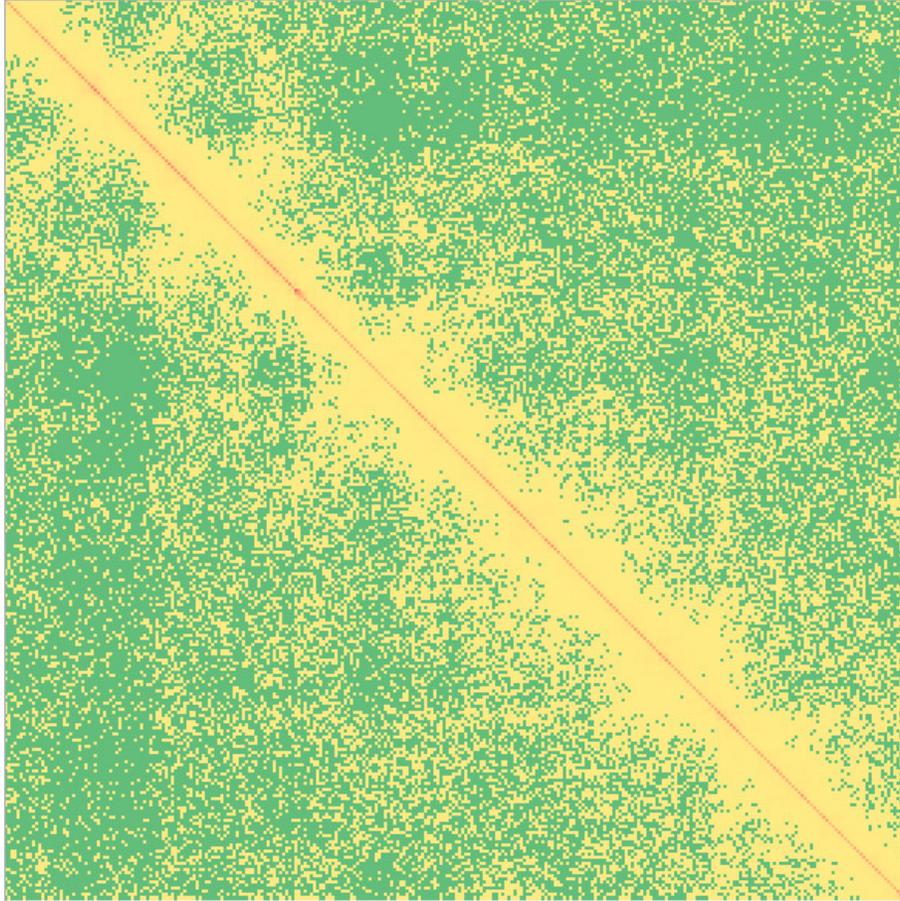


Figure1. Heatmap of the contact matrix of the long arm of Chromosome 1 produced by Lieberman-Aiden on HI-C data.

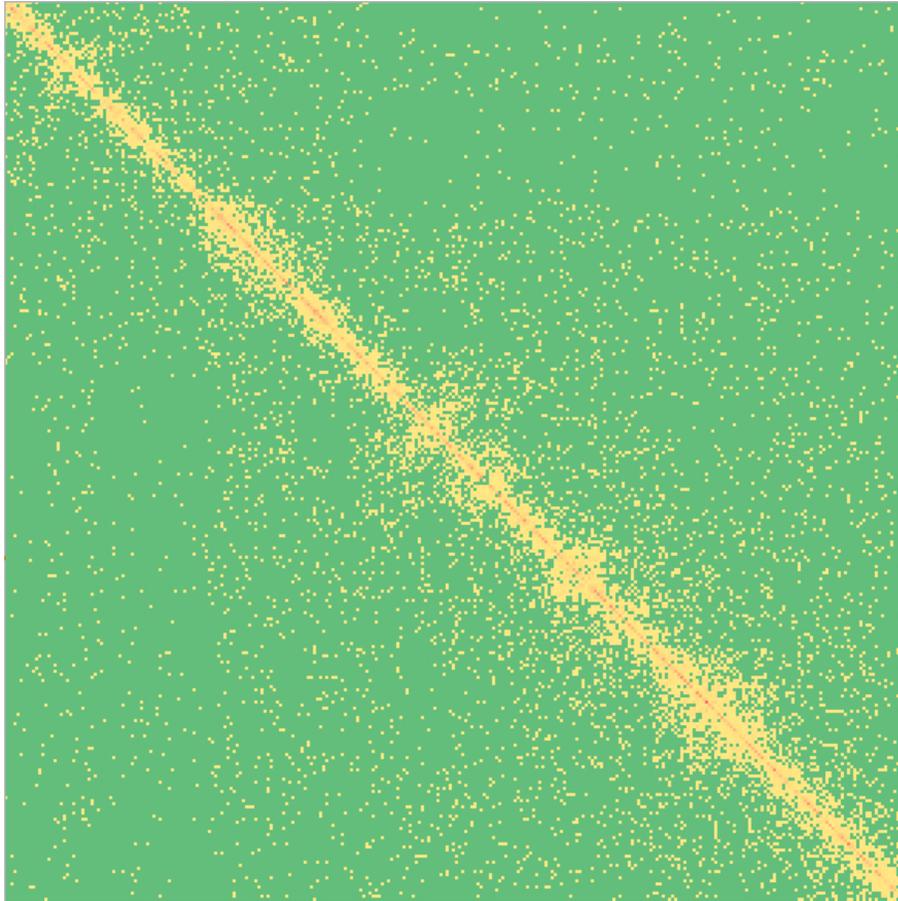


Figure2. Heatmap of the raw contact matrix of the long arm of Chromosome 1 produced by TADbit.

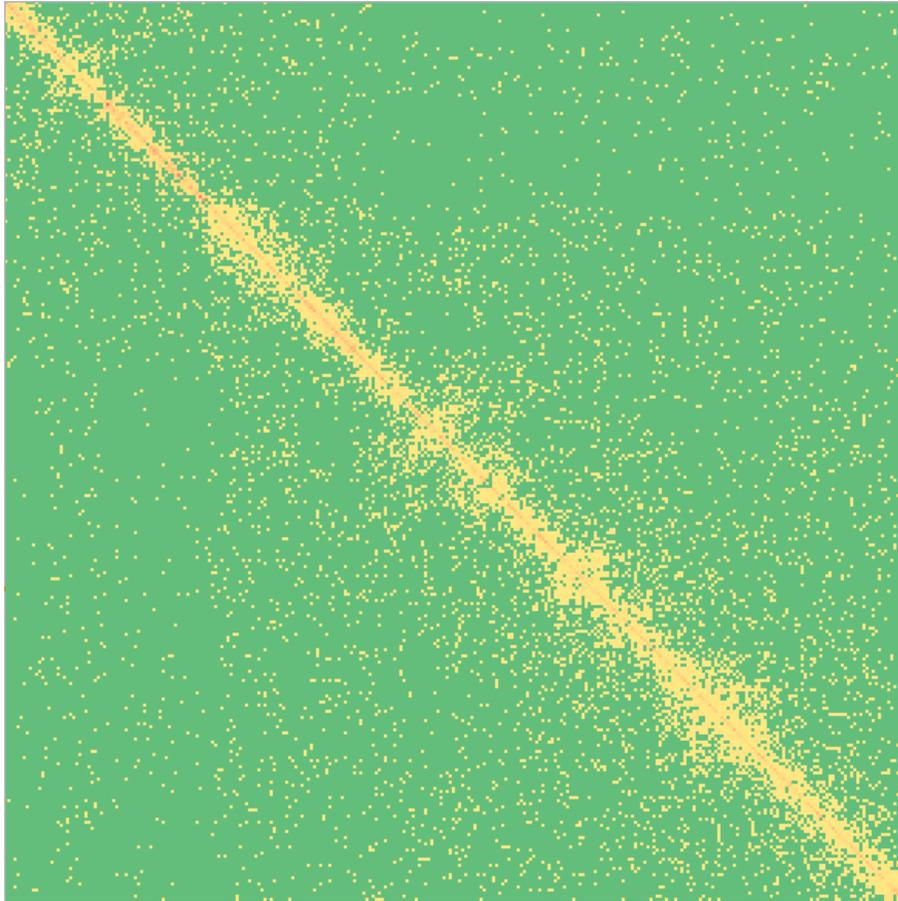


Figure3. Heatmap of the normalized (ICE) contact matrix of the long arm of Chromosome 1 produced by TADbit.

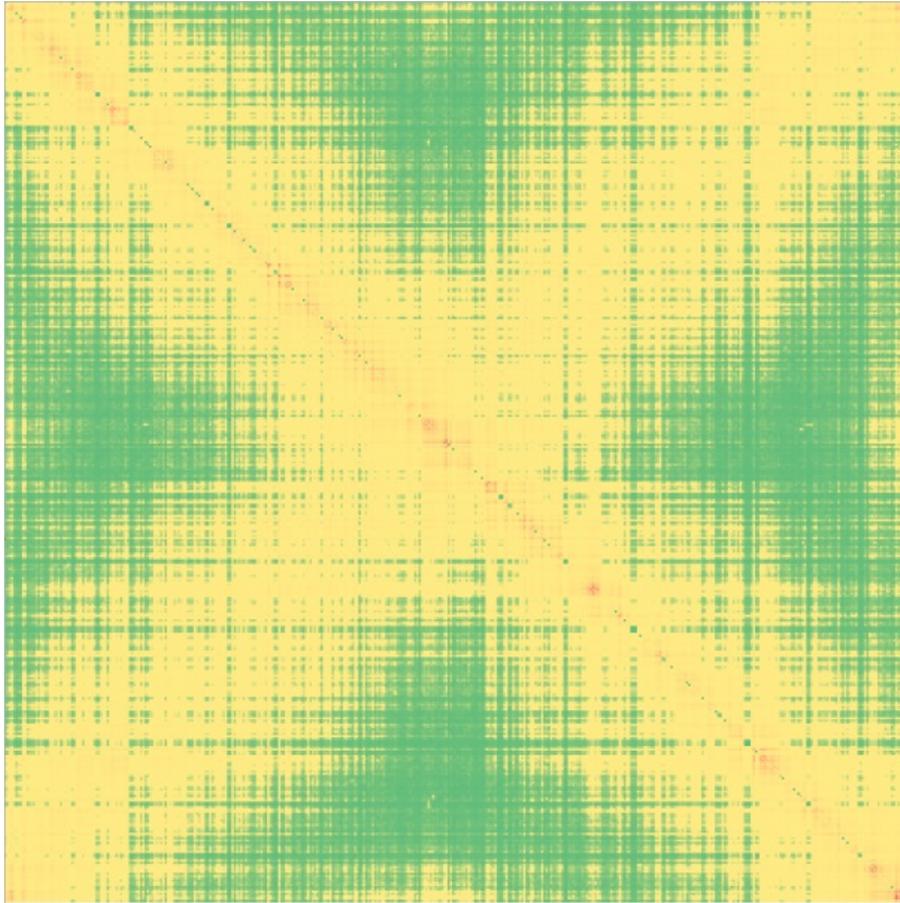


Figure4. Heatmap of the contact matrix of *Caulobacter Crescentus* produced by Imakaev on HI-C data.

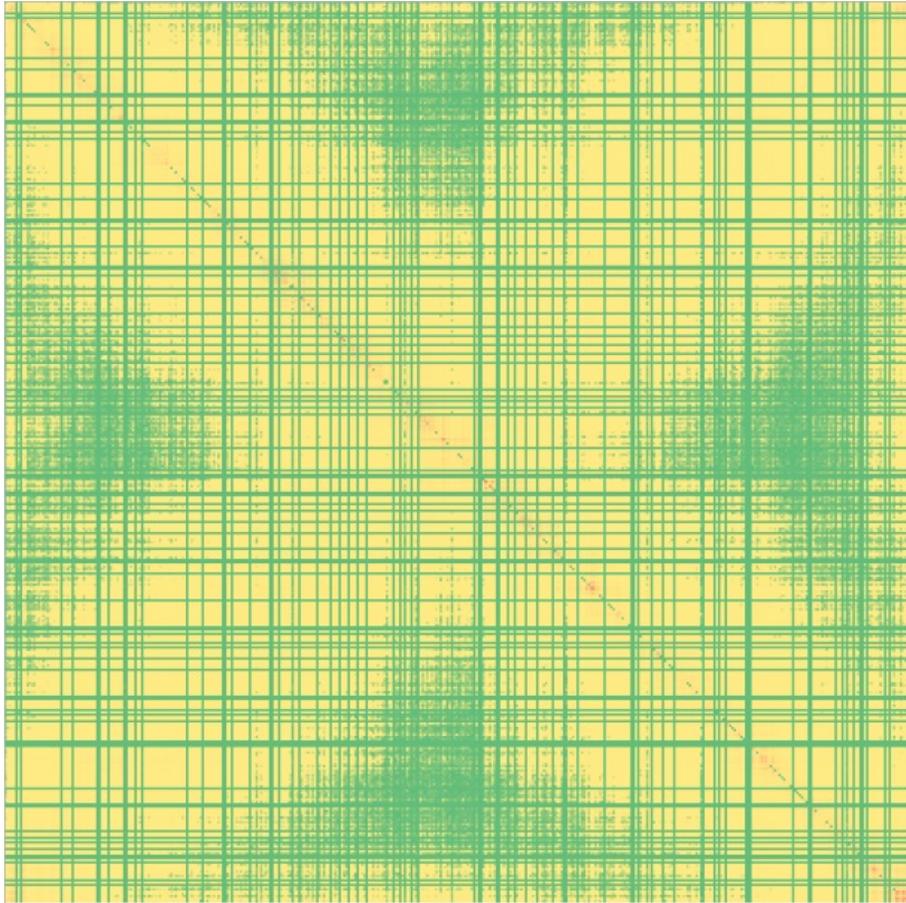


Figure5. Heatmap of the raw contact matrix of *Caulobacter Crescentus* produced by TADbit.

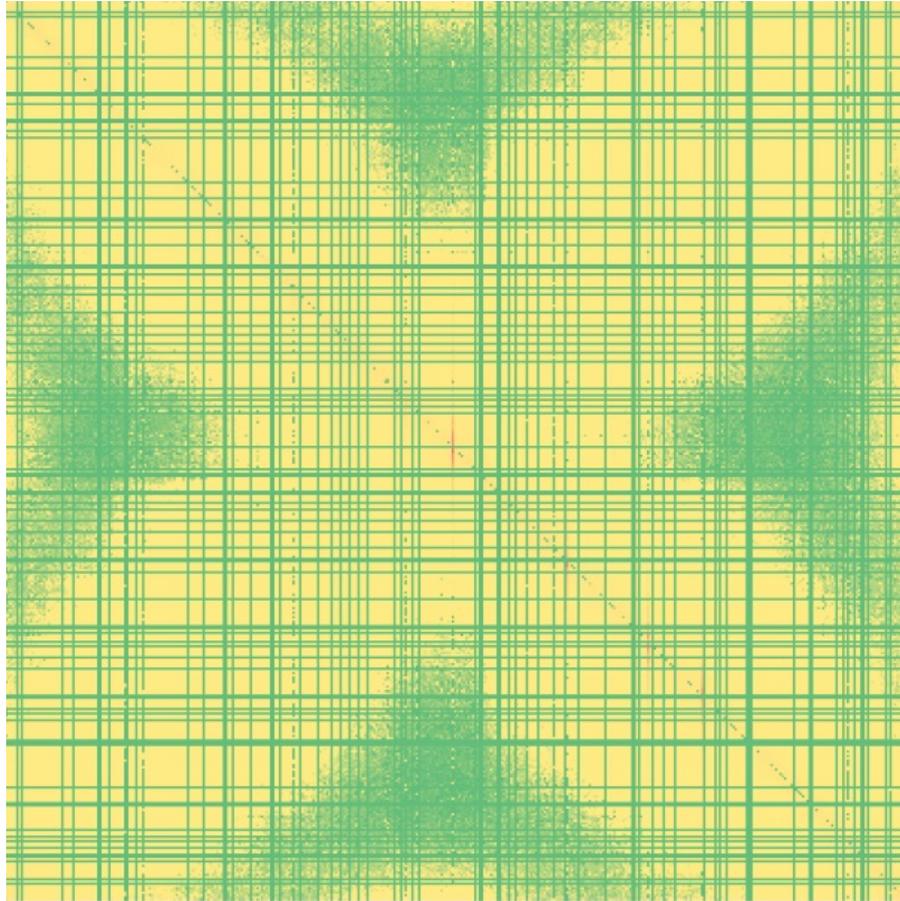


Figure6. Heatmap of the normalized (ICE) contact matrix of *Caulobacter Crescentus* produced by TADbit.

1	0,0	12,6107258575	16,1402320232	16,3287736996	12,8401805061	0,0	11,1466695648	0,0	6,59653709632	6,8927590615	4,68806855886	5,12020949857	5,058190274	6,8100314345
2	12,6107258575	12,8049494735	22,2228420865	21,5130936008	19,8911892096	0,0	12,1965429333	0,0	9,26780811065	9,02195724865	8,14990203915	6,87828285173	5,79531658946	6,2600310615
3	16,1402320232	22,2228420865	11,955529315	25,556297925	20,9400544956	0,0	14,8176378357	0,0	9,8554030667	7,51515361206	7,7323312621	7,4606145024	6,37861643611	5,89866142413
4	16,3287736996	21,5130936008	25,556297925	19,8911892096	25,2152171769	0,0	14,8510233444	0,0	10,0177953208	8,0691878087	6,58501145274	7,3694840687	6,78984405842	5,76392282656
5	12,8401805061	19,8911892096	20,9400544956	25,2152171769	7,25344058699	0,0	13,0446886862	0,0	10,9915743996	9,86760318623	7,38425164396	7,32067967162	7,84936622198	5,34937483506
6	0,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0
7	11,1466695648	12,1965429333	14,8176378357	14,8510233444	13,0448886862	0,0	0,0	0,0	15,3705613272	12,8875555335	10,671495687	13,0220593826	9,7893846136	9,0743577178
8	0,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0
9	6,59653709632	9,26780811065	9,8554030667	10,0177953208	10,9815743996	0,0	15,3705613272	0,0	14,1478396421	22,6254613455	16,4410879524	14,8641517871	14,2384296071	12,192724533
10	6,8927590615	9,02195724865	7,91515361206	8,0691878087	9,8676031862	0,0	12,8875555335	0,0	22,6254613455	15,4828480888	19,0686935218	16,9514915222	15,8360584842	13,1621050046
11	4,68806855886	5,12020949857	1,7732312621	6,8501145274	7,3645314368	0,0	16,4410879524	0,0	16,4410879524	15,0686935218	0,0	18,9663049566	20,6213537253	13,361341305
12	5,12020949857	6,87828285173	7,4696145024	7,3694840687	7,32067967162	0,0	13,0220593826	0,0	14,8641517871	16,9514915222	18,9663049566	0,0	21,7387107278	16,4605449608
13	5,058190274	5,79531658946	6,3786164361	6,78984405842	7,84936622198	0,0	9,7893846136	0,0	14,2384296071	15,8360584842	20,613537253	21,7387107278	0,0	11,3493102094
14	6,8100314345	4,2600310615	5,99868142413	5,76392282656	7,34937483506	0,0	9,0743577178	0,0	12,192724533	13,1621050446	13,361341245	16,065449608	11,3493102094	6,2083036585
15	4,5159158761	5,5649110429	5,95016074179	5,6757940639	6,837593127	0,0	8,91918520005	0,0	11,483807571	12,1328480847	13,9618335317	15,1064462526	15,6209462884	12,57223993
16	3,390565789	5,1166223602	6,90997447024	5,2851224157	6,5860215556	0,0	7,02475126289	0,0	10,2766213571	11,8124833931	13,3250181892	13,3473667208	14,0068346132	19,4470181815
17	3,48740072567	4,475010190416	4,28946468067	4,34911055322	5,34484962129	0,0	6,26783287199	0,0	7,34421063661	8,06895054174	8,70852118667	8,8854427261	10,0165681933	12,0897029744
18	2,95653707347	4,1540248469	4,2491108803	4,5159426147	4,4724458924	0,0	6,2101248621	0,0	7,82375269859	8,38410124976	9,08978866873	9,44019782058	9,8025060865	12,4879773817
19	2,4741002324	3,2151927887	3,07111501053	3,493690742	3,44889693841	0,0	4,3078878104	0,0	4,8445919793	5,01702630342	5,79858992808	6,1128748454	5,37743674291	5,8486576889
20	1,4318041622	3,7899018256	3,92944204273	3,1691810135	3,73776896228	0,0	3,4115294203	0,0	4,5493898024	4,4202243923	5,5995881406	5,3844312919	5,39759172241	7,35781759886
21	2,50407074225	3,3049375161	3,26122473004	2,2214544045	3,09424070326	0,0	3,65182575529	0,0	4,0093928574	4,12732174649	4,86404804651	5,10975284965	5,4409552277	12,1265166597
22	2,13807415626	2,8890222681	2,91816352634	2,78476578056	2,82573963938	0,0	3,65800754725	0,0	3,86889702502	3,9287727996	4,75913047352	4,58981613227	5,11484841753	5,84754463833
23	2,21710801542	3,41622488202	3,0372954285	2,9073174368	2,93075730235	0,0	3,37464689108	0,0	3,85363082346	4,0738501053	4,89909491706	4,57558886553	5,04830958444	5,2017441454
24	4,70772445887	2,73536921292	3,81416031383	2,5800333545	2,8781617106	0,0	5,5962371353	0,0	3,4098931215	3,53497714582	4,2326101289	4,047778139	4,90075539174	5,5513090356
25	1,8415152666	2,43792234987	2,63171625643	2,30978072186	2,3805671066	0,0	3,14602469967	0,0	3,0119785659	3,30438657607	4,0674356382	3,21181970681	3,99428028459	4,1022938921
26	0,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0
27	1,35550359776	1,61791521337	1,91851974755	2,00793918825	2,03862215466	0,0	2,41374447132	0,0	2,41374446051	2,34057926243	3,42281905177	3,00567942969	2,77397563257	3,76717560748
28	1,6572855124	1,97127649566	1,90978196053	1,76971274628	1,81882856185	0,0	2,34554397229	0,0	2,2847266971	2,2984241521	2,35830009566	2,8189189242	2,75242057476	2,61466119823
29	1,3493030402	2,09433046421	1,650729231	1,6551397508	1,7107828461	0,0	1,93234498444	0,0	2,091048086	2,4252300383	3,0480128474	2,4373933754	2,5069578411	3,0441932127
30	1,0841811477	1,94648485994	1,85328935182	1,80276885099	2,1375148918	0,0	1,93744153111	0,0	2,0741882691	2,0068161227	2,3492495569	2,4174687113	2,27711429086	2,89426304588

Figure10. Contact matrix of normalized (ICE) data of *Caulobacter Crescentus* produced by TADbit.

1. Verifying that the outputs of ChromStruct4 do not depend on data normalization, because ChromStruct4 uses a filtering method that makes it robust against biases.
2. Comparing final configurations produced by TADbit and ChromStruct4 from the two methods are based on different approaches. If the final outputs are similar, the features characterizing such configurations could actually mark biological characteristics captured by both methods.

The following results have been produced:

- ChromStruct4: 100 configurations of Chromosome 1 starting from the raw matrix
- ChromStruct4: 100 configurations of Chromosome 1 starting from the normalized matrix
- TADbit: 100 configurations of Chromosome 1 starting from the normalized matrix
- ChromStruct4: 100 configurations of *Caulobacter* starting from the raw matrix
- ChromStruct4: 100 configurations of *Caulobacter* starting from the normalized matrix
- TADbit: 100 configurations of *Caulobacter* starting from the normalized matrix

Figures 11-16 demonstrate how the configurations produced by ChromStruct4 are more varied in shape than those produced by TADbit, which are very homogeneous. Conformations of both Chromosome 1 and *Caulobacter* produced by TADbit are very similar to each other, this behavior may be due to the algorithm leaving few degrees of freedom or to the biases given by harmonics governing TADbit's evolution.

Figures 17-22 show the boxplots of Euclidean-distance vs genomic-distance for all output sets. These graphs help us to understand the variability of the

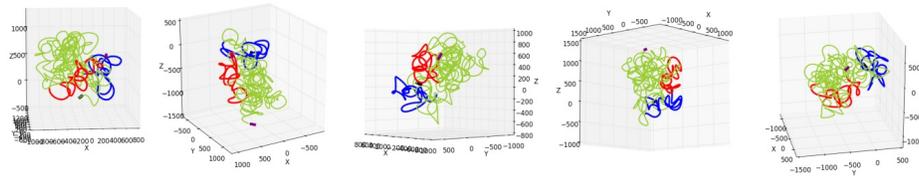


Figure11. Configurations of Chromosome 1 produced by ChromStruct4 from the raw contact matrix.

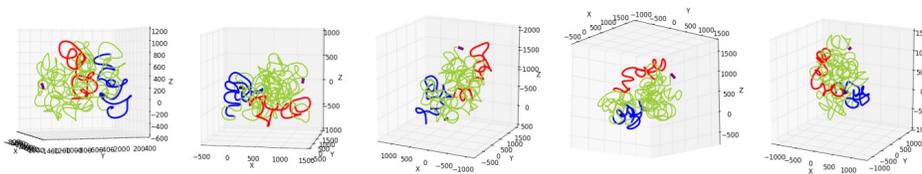


Figure12. Configurations of Chromosome 1 produced by ChromStruct4 from the normalized contact matrix.

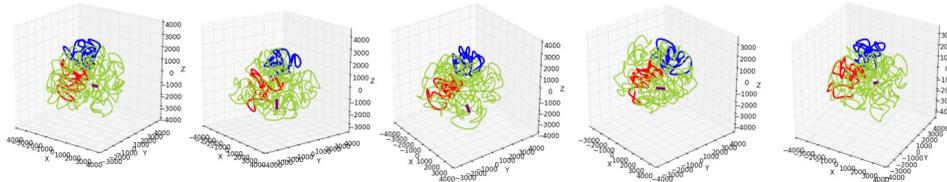


Figure13. Configurations of Chromosome 1 produced by TADbit from the normalized contact matrix.

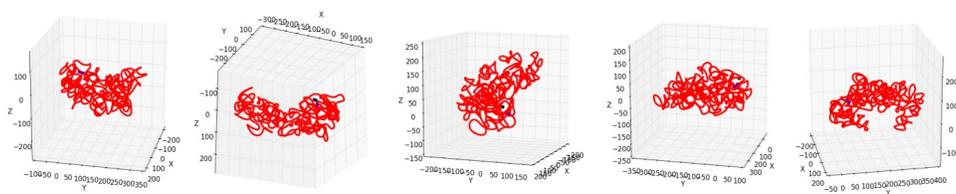


Figure14. Configurations of Caulobacter produced by ChromStruct4 from the raw contact matrix.

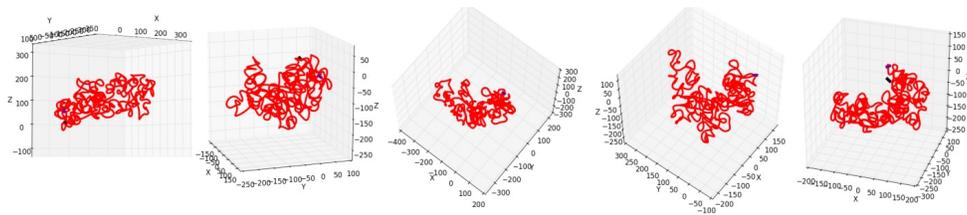


Figure15. Configurations of Caulobacter produced by ChromStruct4 from the normalized contact matrix.

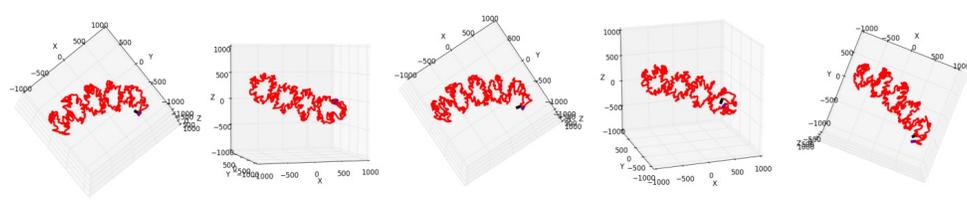


Figure16. Configurations of Caulobacter produced by ChromStruct4 from the normalized contact matrix.

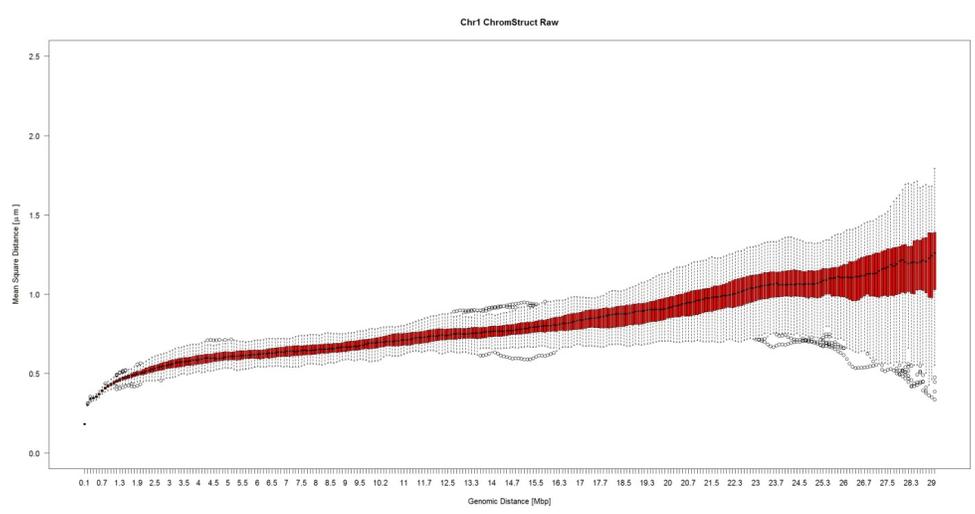


Figure17. Boxplot of 100 conformations produced by ChromStruct4 on raw data of Chromosome 1.

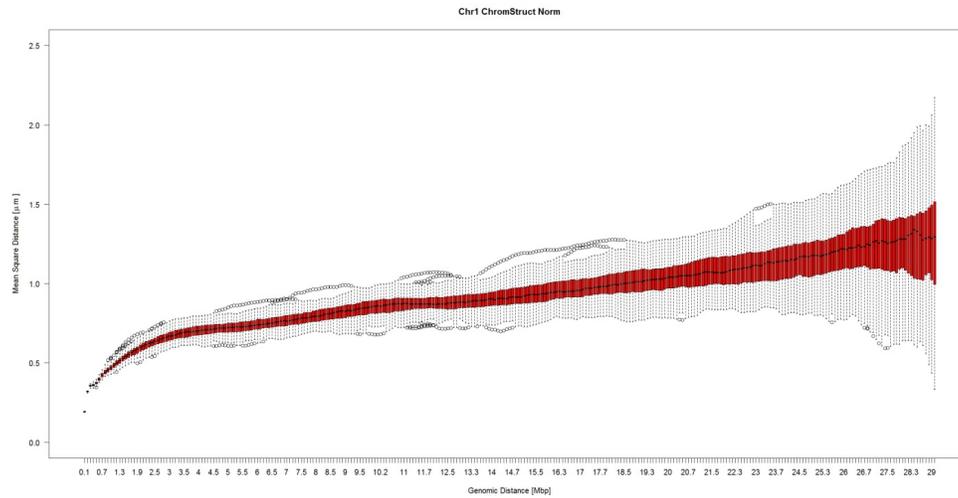


Figure18. Boxplot of 100 conformations produced by ChromStruct4 on normalized data of Chromosome 1.

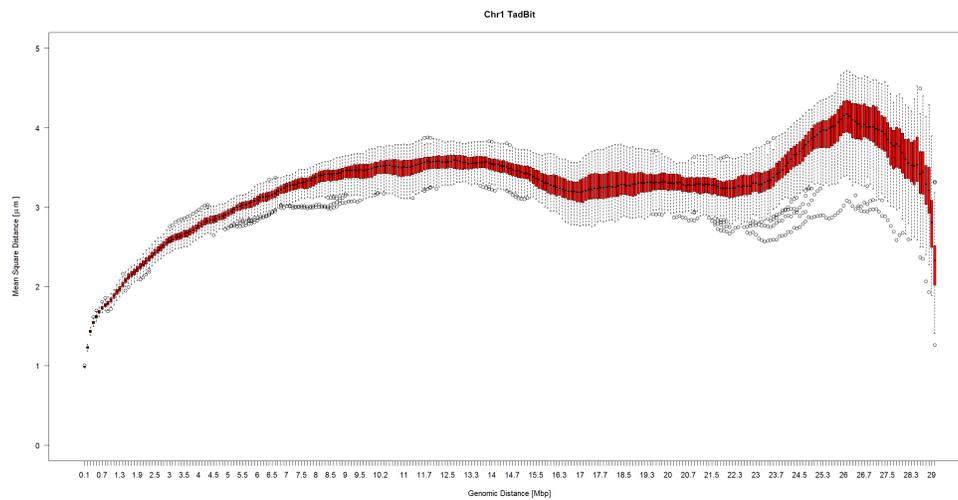


Figure19. Boxplot of 100 conformations produced by TADbit on normalized data of Chromosome 1.

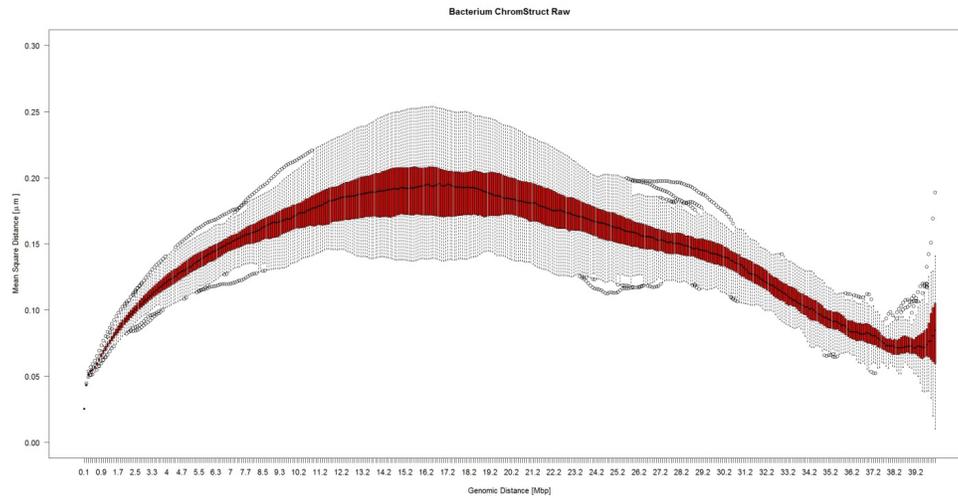


Figure20. Boxplot of 100 conformations produced by ChromStruct4 on raw data of Caulobacter.

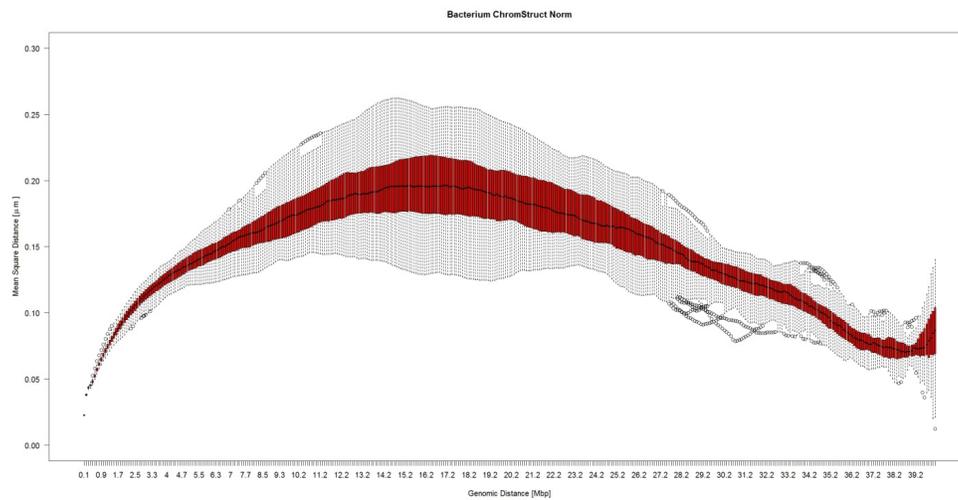


Figure21. Boxplot of 100 conformations produced by ChromStruct4 on normalized data of Caulobacter.

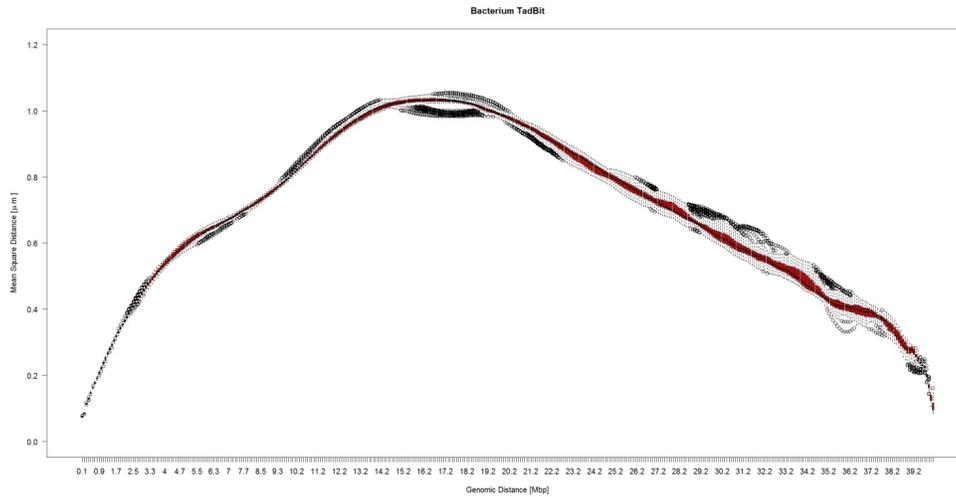


Figure22. Boxplot of 100 conformations produced by TADbit on normalized data of Caulobacter.

solutions with respect to their genomic organization. As can be seen, the configurations of Caulobacter and Chromosome 1 produced by TADbit are much less variable than those produced by ChromStruct4 (boxplots are less expanded). Other important observations can be made about size:

- The maximum size of the portion of Chromosome 1 [150.3Mb,179.4Mb] modelled by ChromStruct4 is about 2.5 microns, while the conformations produced by TADbit on the same data reach nearly 5 microns.
- The size of the Caulobacter conformations produced by ChromStruct4 are smaller, never exceeding 0.3 microns, than those produced by TADbit, always more than 1 micron.

Synthetic contact matrices have been produced for both ChromStruct4 and TADbit. Threshold of contact between two beads has been assumed as the sum of their radii (calculated by ChromStruct4) multiplied by 1.2.

Figures 23-28 show that the conformations produced by TADbit are much less varied than those produced by ChromStruct4, the contacts are infrequently scattered and almost always concern the same bin pairs. Moreover, for Caulobacter, the red areas (i.e. high frequencies) are concentrated around the main diagonal.

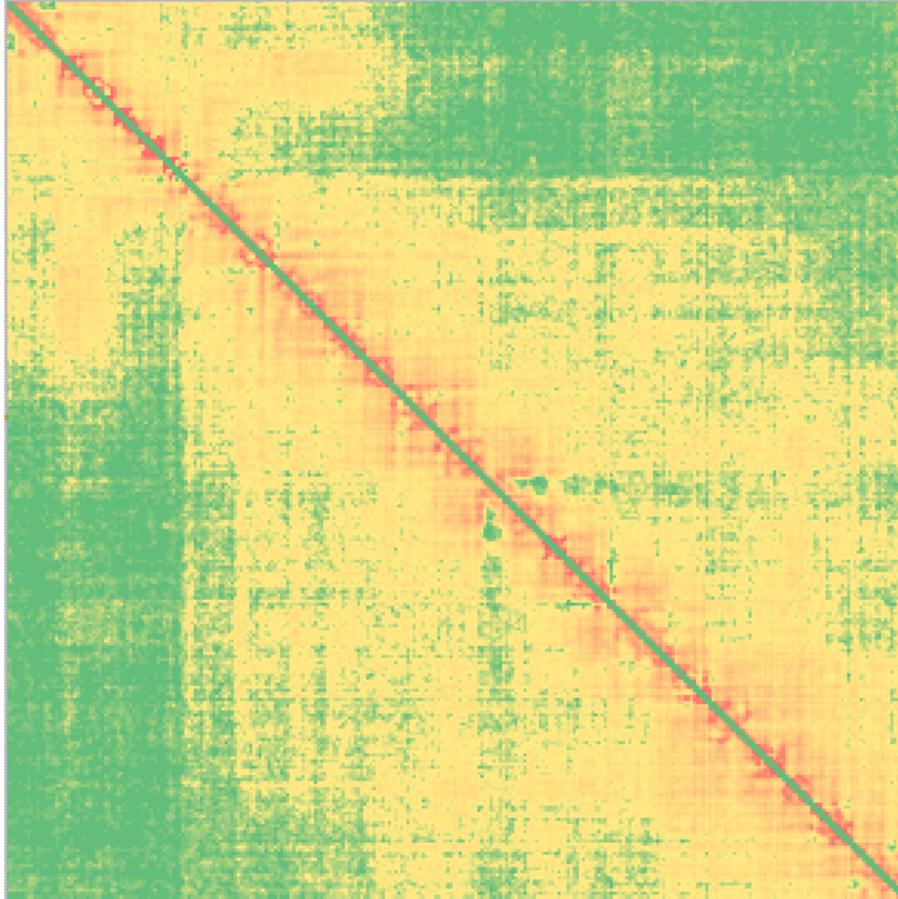


Figure23. Synthetic contact matrix of 100 conformations produced by ChromStruct4 on raw data from Chromosome 1.

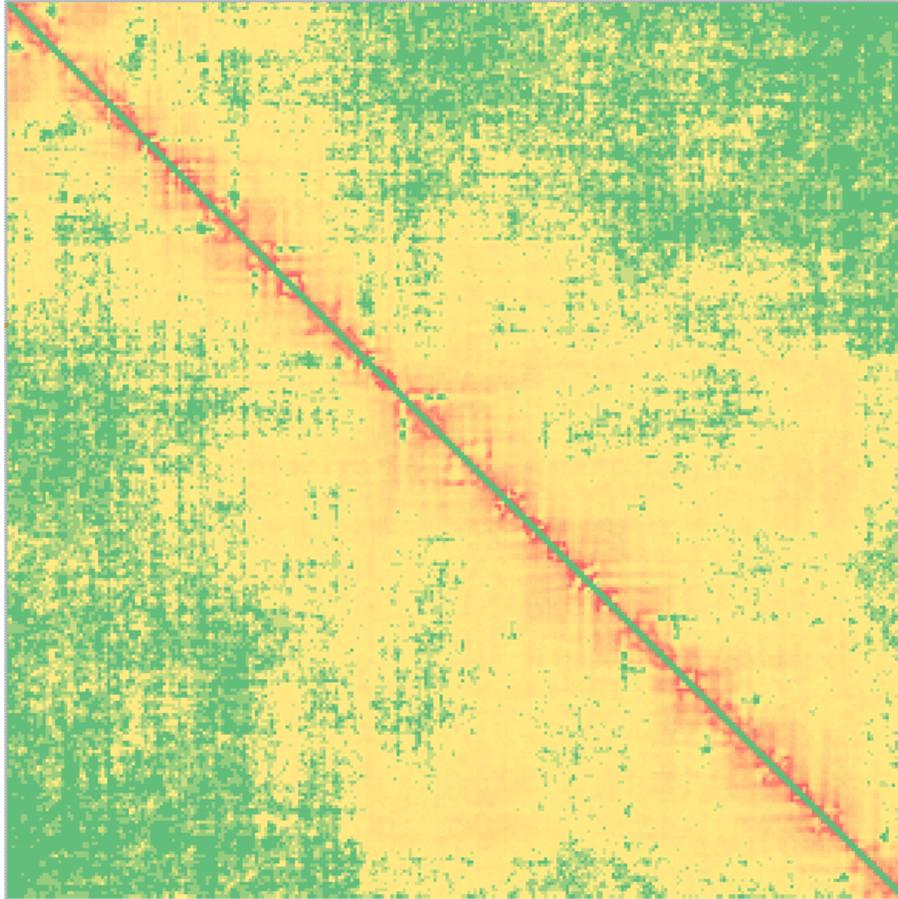


Figure24. Synthetic contact matrix of 100 conformations produced by ChromStruct4 on normalized data from Chromosome 1.

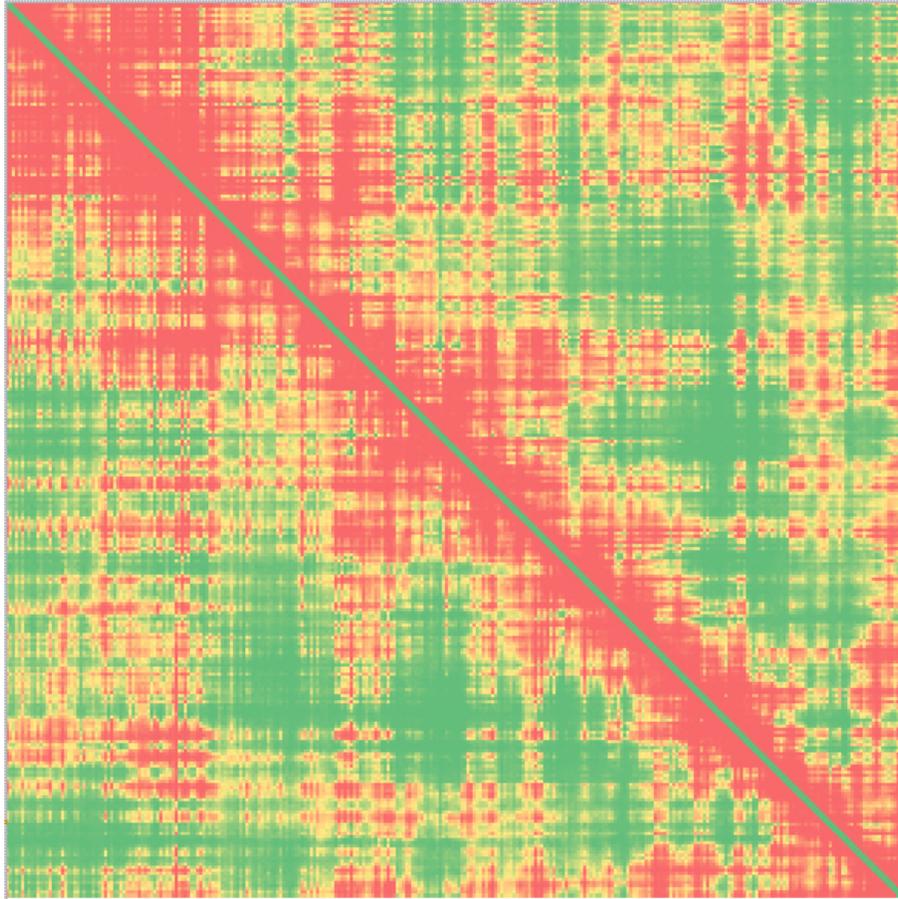


Figure25. Synthetic contact matrix of 100 conformations produced by TADbit on normalized data from Chromosome 1.

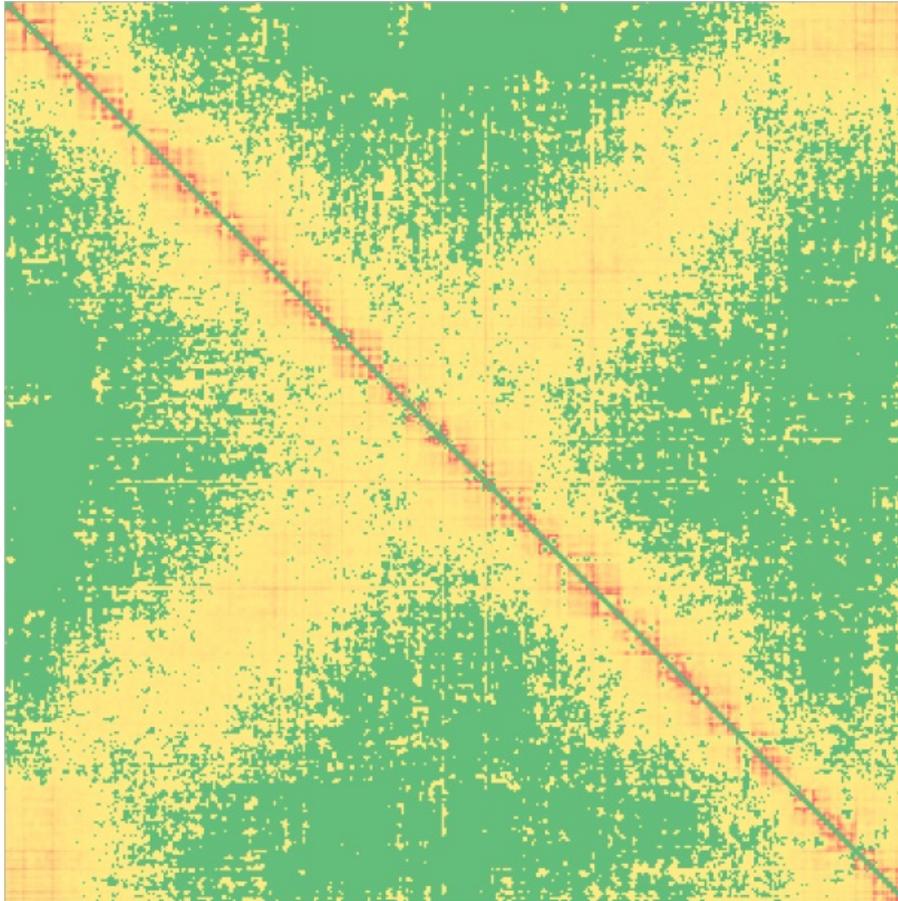


Figure26. Synthetic contact matrix of 100 conformations produced by ChromStruct4 on raw data from Caulobacter.

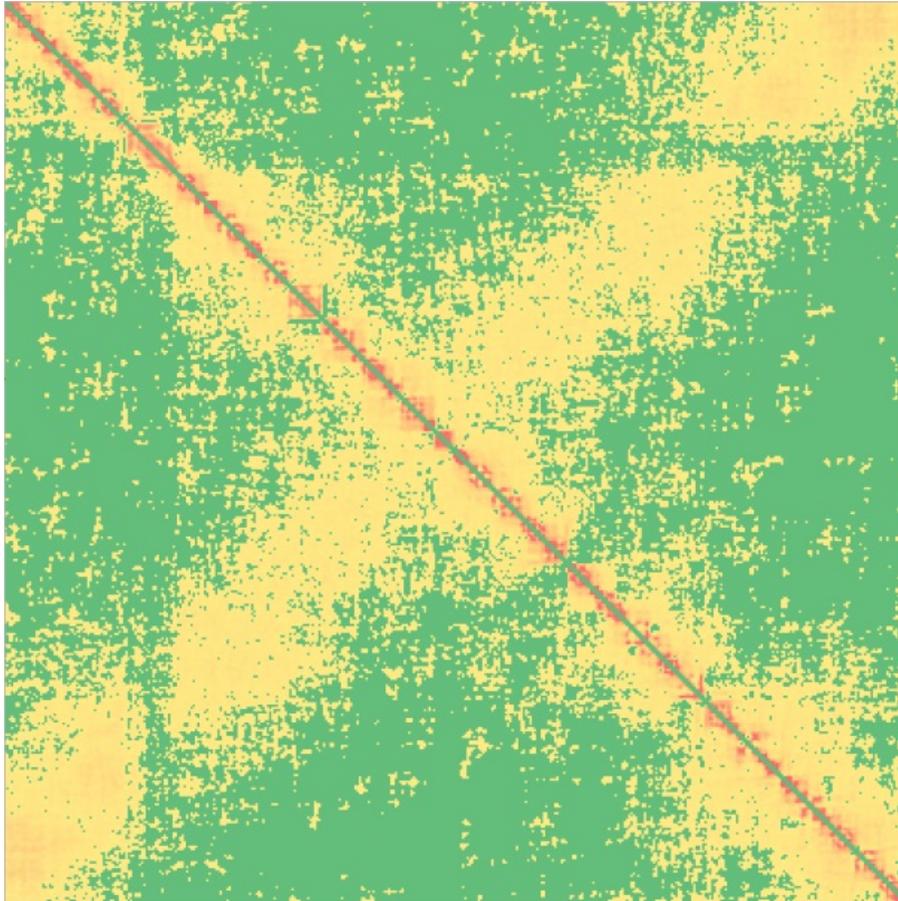


Figure27. Synthetic contact matrix of 100 conformations produced by ChromStruct4 on normalized data from Caulobacter.

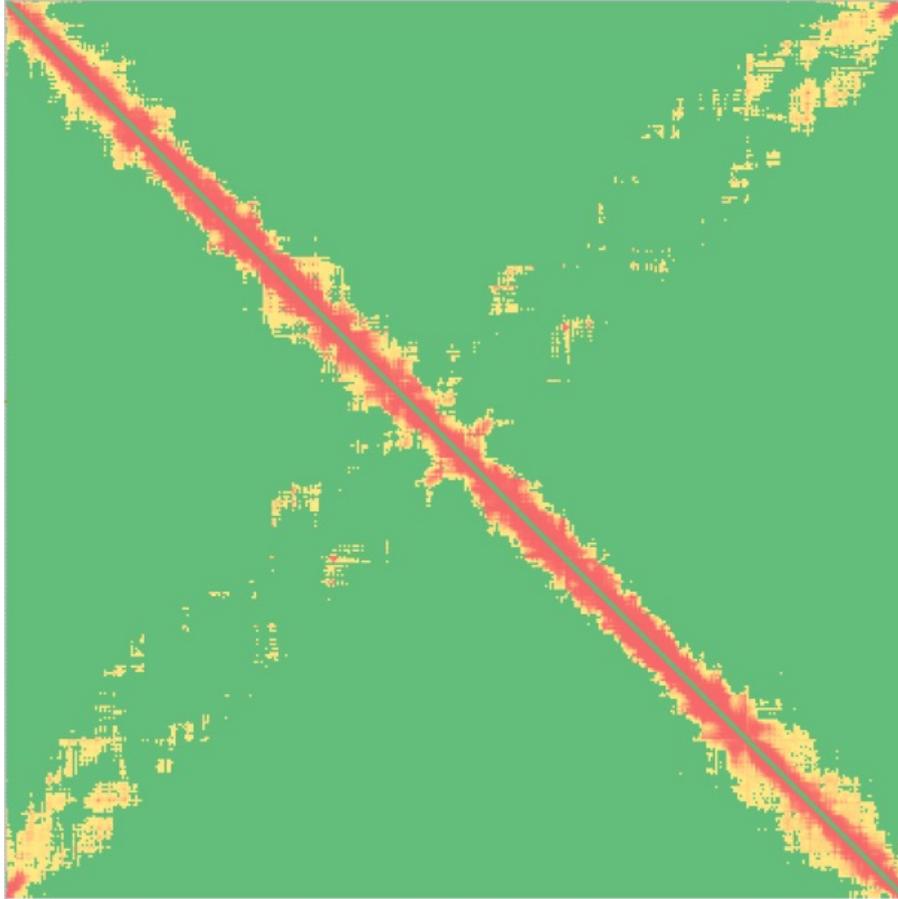


Figure28. Synthetic contact matrix of 100 conformations produced by TADbit on normalized data from Caulobacter.

4 Conclusions

From comparative experiments between ChromStruct4 and TADbit, the following conclusions can be drawn:

- ChromStruct4 is robust against biases, because boxplots Euclidean-distance vs genomic-distance and synthetic matrices, for both Chromosome 1 and *Caulobacter*, are very similar in raw and normalized data (ICE).
- TADbit and ChromStruct4 produce different conformations, both in structure and sizes. This means that the two approaches focus on different strategies in the chromatin's three-dimensional structure reconstruction. Starting from the same data, they produce different results. The TADbit solutions are also much less variable than those produced by ChromStruct4.

References

1. Baú, D. *et.al.* (2012): Genome structure determination via 3C -based data integration by the Integrative Modeling Platform, *Methods*, 58(3): 300306
2. Lieberman-Aiden, E. *et.al.* (2009): Comprehensive Mapping of Long-Range Interactions Reveals Folding Principles of the Human Genome, *Science* 326: 289-293.
3. Imakaev, M. *et.al.* (2012): Iterative correction of Hi-C data reveals hallmarks of chromosome organization, *Nature Methods* 9(10): 999-1003.
4. Dixon, J.R. *et.al.* (2012): Topological domains in mammalian genomes identified by analysis of chromatin interactions, *Nature* 485: 376-380.
5. Cudai, C. *et.al.* (2015): Inferring 3D chromatin structure using a multiscale approach based on quaternions, *BMC Bioinformatics*, 16: 234.
6. Cudai, C. *et.al.* (2015): A Statistical Approach to Infer 3D Chromatin Structure, *Mathematical Models in Biology*, Springer, 161-171.