# Puzzle Hi-C: an accurate scaffolding software

Guoliang Lin[1]¶, Zhiru Huang[2]¶, Tingsong Yue[1]¶, Jing Chai[3]¶, Yan Li[1]¶, Huimin Yang[1], Wanting Qin[1], Guobing Yang[1], Robert W. Murphy[3,5], Ya-ping Zhang[1,3,6]*, Zijie Zhang[1,6]*, Wei Zhou[4]*, Jing Luo[1,6]*

[1] State Key Laboratory for Conservation and Utilization of Bio-resource, School of Ecology and Environment, School of Life Sciences and School of Medicine, Yunnan University, Kunming, 650091 Yunnan, China.

[2] Department of Biology, University of Waterloo, Waterloo, ON N2L 3G1, Canada.

[3] State Key Laboratory of Genetic Resources and Evolution, Kunming Institute of Zoology, Chinese Academy of Sciences, Kunming, 650223 Yunnan, China.

[4] National Pilot School of Software, Yunnan University, Kunming, 650091 Yunnan, China.

[5] Reptilia Zoo and Education Centre, 2501 Rutherford Rd., Vaughan, ON L4K 2N6, Canada.

[6] Southwest United Graduate School, Yunnan University, Kunming, 650091 Yunnan, China.

¶These authors contributed equally to this work.

*Corresponding author. Email: jingluo@ynu.edu.cn (J.L.); zwei@ynu.edu.cn (W.Z.); zijiezhang@ynu.edu.cn (Z.Z.); zhangyp@mail.kiz.ac.cn (Y.Z.)

## Abstract

High-quality, chromosome-scale genomes are essential for genomic analyses. Analyses, including 3D genomics, epigenetics, and comparative genomics rely on a high-quality genome assembly, which is often accomplished with the assistance of Hi-C data. Current Hi-C-assisted assembling algorithms either generate ordering and orientation errors or fail to assemble high-quality chromosome-level scaffolds. Here, we offer the software Puzzle Hi-C, which uses Hi-C reads to accurately assign contigs or scaffolds to chromosomes. Puzzle Hi-C uses the triangle region instead of the square region to count interactions in a Hi-C heatmap. This strategy dramatically diminishes scaffolding interference caused by long-range interactions. This software also introduces a dynamic, triangle window strategy during assembly. Initially small, the window expands with interactions to produce more effective clustering. Puzzle Hi-C outperforms available scaffolding tools.

## Introduction

Analysis of genomic data rely on high-quality chromosome-level genomes. Accuracy is essential for downstream genomic analyses, and especially for 3D, comparative, and functional genomic analyses. For example, due to significant interactions being calculated on the linear distance of the genome, 3D analyses can create an incorrect

36  assembly, which leads to false positive interactions, resulting in unreliable results(1,2).
37  Chromosome evolution and estimation of recombination rely on the contiguity of the
38  genome(3,4). Only chromosome-scale genomes can drive an understand the complexity
39  of regulatory architecture and how cis-regulatory elements influence genes, because a
40  cis-regulatory element is likely more than 1 Mb apart from the target gene(5–7). A
41  contiguous genome can significantly improve the interpretation of genome-wide
42  association studies (GWAS)(8–11) because usually regions of linkage are not on
43  different contigs. Thus, high-quality chromosome-level genomes are requisite for
44  multiple downstream genomic analyses(12).

45  Long-read sequencing technologies have facilitated genome assembly because they
46  yield large, overlapping repetitive regions. Notwithstanding, such contigs do not always
47  stretch into a complete chromosome or even one arm of a chromosome. To obtain
48  chromosome-scale scaffolds, various strategies have been explored to increase the
49  contiguity of *de novo* genome assemblies. Two primary methods order and orient
50  scaffolds for chromosome-level assembly: genetic mapping and high-throughput
51  chromatin conformation capture (Hi-C). Traditional genetic mapping orders and orients
52  scaffolds based on linkage groups. However, the construction of a genetic map requires
53  a large number of individual offspring to be sequenced. This dramatically limits the
54  application of genetic maps in the genome because many species have K reproductive
55  strategies that cannot satisfy the requirement by having enough offspring(13–15).
56  Further, the sequencing of a large number of individuals costs time, computing
57  resources, storage resources, and other expenses(16). By contrast, recently developed
58  Hi-C provides a powerful tool for chromosome-level assembly. Hi-C merely requires a
59  small number of tissue samples to mount more than 85% of sequences to
60  chromosomes(17,18). Consequently, Hi-C is the most commonly used method for
61  scaffolding at the chromosome level.

62  Three invariant features of Hi-C interactions are used for genome assembly(19): intra-
63  chromosomal interaction enrichment, the random positioning of chromosomes in the
64  nucleus, and the local smoothness of interactions reflected in the Hi-C heat map. In the
65  Hi-C interaction matrix, a locus tends to interact more frequently with another locus
66  within the same chromosome (cis-interactions) than with loci on a different
67  chromosome (trans-interactions). Two phenomena of 3D chromatin may contribute to
68  this feature. In the first phenomenon, chromosomes occupy distinct volumes
69  throughout the cell cycle, leading to physical separation between chromosomes(20,21).
70  The second phenomenon relies on the random positioning of chromosomes in the
71  nucleus(22), which may largely reduce chromosomal interactions, and the probability
72  of intra-chromosomal interaction decreases with increasing linear distance. Thus, in Hi-
73  C interaction maps, the frequency of interaction tends to decrease with genomic
74  distance, i.e., a locus interacts more frequently with other loci that are nearby in the
75  genomic space than with distant loci. When the distance is greater than 100kb, the
76  probability of interaction is about $1/x$, where x is the distance between two
77  points(17,23). Finally, local smoothness of interactions as reflected in the Hi-C heat

78 map interaction of adjacent points tends to be consistent(24). Available software
79 commonly use features one and two for scaffolding.

80 Accompanying the sharply decreasing price of genomic sequencing, Hi-C data for
81 scaffolding at chromosome-level is now accessible and popular. The first Hi-C
82 scaffolding software, LACHESIS(17), developed in 2013, has seen increasingly
83 employment for constructing chromosome-level in genomes. Several software options
84 exist for Hi-C scaffolding(17,25–28), however none can eliminate errors including
85 artificial relocations, translocations, and inversions(26,28), which result in false
86 assemblies and erroneous genomes. Chromosomes may fold into various structures,
87 such as loops, topological associated domains (TADs), and compartments, which lead
88 to many long-range interactions. Such interactions violate the assumption that the
89 probability of the intra-chromosomal interaction decreases with the linear distance,
90 thus causing an incorrect assembly. Some methods use a contig-end solution(29,30),
91 but the employment of limited information results in disappointing performance. To
92 address these issues, we offer an easy-to-use Hi-C scaffolding software, Puzzle Hi-C,
93 which uses a triangle window and iterative assembly strategy to reduce long-range
94 interactions thereby improving performance. This software achieves outstanding
95 scaffolding results in simulated data and real data. The source code and
96 documentation of Puzzle Hi-C are available at GitHub
97 (https://github.com/linguoliang/puzzle-hi-c.git).

## Methods

98

### Datasets

99
100 We used human's Hi-C data from the GM12878 cell line as benchmark data. We used
101 Hi-C data from a gayal (*Bos frontalis*) and puffer fish (*Takifugu bimaculatus*) as
102 examples. We also used the Hi-C data from a broomcorn millet (*Panicum miliaceum*),
103 Indian cobra (*Naja naja*), Peking duck (*Anas platyrhynchos*), water buffalo (*Bubalus*
104 *bubalis*), yellow croaker (*Larimichthys crocea*), and fighting fish (*Betta splendens*) (S1
105 Table).

106

107 **Puzzle Hi-C pipeline.** The Puzzle Hi-C pipeline contains three steps: mapping,
108 scaffolding, and building. Briefly, Puzzle Hi-C uses juicer software for mapping step.
109 Next, scaffolding takes each chromosome group as input and iteratively merges all
110 chromosome groups' contigs into one chromosome. Finally, building reconstructs each
111 chromosome by concatenating the contigs, adding gaps between the contigs, and
112 generating the genome in the FASTA format (Fig 1).

113

114 **Puzzle Hi-C mapping.** Puzzle Hi-C uses juicer software for mapping. Juicer uses bwa
115 mem default parameters for mapping, after mapping juicer program will filter the
116 duplicate matched sequences, in juicer it is considered that if two pairs of double-end
117 sequence matching results, their position information only differ by 4 pb, then it is
118 considered a duplicate sequence, juicer will remove this part duplicate sequences and

119   keep only the result on one pairwise comparison.
120

121   **Puzzle Hi-C scaffolding.** The iterative algorithm used for scaffolding solves two
122   problems: "ordering" assigns a relative position to each scaffold on each chromosome
123   with respect to the other scaffolds assigned to the same chromosome; and "orienting"
124   determines which of the two ends of each scaffold is adjacent to the preceding scaffold,
125   and which end is adjacent to the next scaffold. In each step, subsets of the input
126   scaffolds are ordered and oriented with respect to one another to create a new, longer
127   set of scaffolds, which are then used as inputs for the next step until all the chromosome
128   contigs are scaffolded in one scaffold. The software uses a weight matrix to build a
129   graph. The graph nodes are the scaffolds in a chromosome-group. Weight is defined as
130   follows:

131   Each end of each scaffold is labeled using B (begin) and E (end). Given two scaffolds
132   i and j, there are four possible connections, BB, BE, EB, and EE. We defined a length
133   cutoff of l and considered the read pairs mapped in the region of length l at both ends
134   (B and E) of the scaffolds.

135   The number of links was determined using:

136   $$N_{i,j} = max\{N_{iB,jB}, N_{iB,jE}, N_{iE,jB}, N_{iE,jE}\}$$

137   For each scaffold i, we only considered the top 5 linked edges. Next, we obtained a
138   link-score for each pair as follows:

139   $$W_{topk} = \frac{N_{topk}}{\sum_{j=1}^{5} N_{top\,j}}$$

140

141   When ordering the scaffolds, each node could only have two edges, so we retained the
142   two edges with the largest and second-largest weight. We set a cutoff for the weight; if
143   the weight was less than the cutoff, then the connection was considered unreliable and
144   removed. The graph found the path with all nodes>1, and constructed new scaffolds
145   according to the path and direction of connection for the next iteration, and increased
146   the length of l at the same time. The iteration stopped only if the number of scaffolds
147   equaled to desired number of chromosomes.

148   **Puzzle Hi-C building.** Once scaffolding is completed, Puzzle Hi-C builds a
149   chromosome-level genome. Scaffolds link with 100 bp N gaps (N can be configured
150   with Puzzle Hi-C parameters). Puzzle Hi-C also generates an agp file to record how
151   scaffolds are assembled with the position and direction information.
152

153 **Genomic collinearity.** The genomic collinearity analysis between genomes were
154 completed using NUCMER from the MUMmer package v3.23 with default
155 parameters(33). After alignment, we sorted the scaffolds according to collinearity to
156 draw the final collinearity figures.

157

158 **Scaffolding error statistics.** To compare the performance of each software, we aligned
159 the genome assembly to their respective reference genomes using the program nucmer
160 with default parameters. Alignment quality was assessed using dnadiff(33), a MUMmer
161 utility that provides detailed information on the differences between two genomes. To
162 get a reliable result, we sampled scaffolding 25 times, where each sample contained 5
163 chromosomes.

164

165 **Hi-C scaffolding in LACHESIS, SALSA2, 3D-DNA, and ALLHiC.**
166 Hi-C reads were mapped using bwa with default parameters. The SAM file, which was
167 generate by bwa, was filtered using PreprocessSAMs.pl.

168

169 For LACHESIS, we used default parameters except CLUSTER_N, which depended on
170 how many chromosomes should be clustered.

171

172 For SALSA2, the minimum input files were provided with the following command
173 line:
174 python run_pipeline.py -a seq.fasta -l seq.fasta.fai -b alignment.bed -e GATC -o
175 scaffolds.

176

177 For ALLHiC, we used the default parameters, except -k parameter. The k parameter
178 was set according to how many chromosomes were clustered.

179

180 For 3D-DNA, Hi-C reads were aligned by juicer software using default parameters.
181 Scaffolds >15 kb were retained, and the haploid model was selected in the 3D-DNA
182 pipeline.

# Results

## Overview of the Puzzle Hi-C algorithm

185 Comparative analysis of existing software found that LACHESIS and ALLHiC
186 used of all interaction information between scaffolds for clustering, and achieved high
187 performance on clustering(17,27). However, SALSA2 only uses partial information at
188 both ends of the scaffold, which advances ordering and orientation(28). Taking
189 accuracy of scaffold ordering and clustering into account, Puzzle Hi-C dynamically
190 changes the size of the statistical window at both ends of the scaffold, and increases the
191 window size as the number of iterations increases(Fig 1). Therefore, our software uses
192 local information at both ends of the scaffold for ordering and orientation at the initial
193 assembly, and as the statistical window increases, it uses global interactions for better
194 clustering (Fig 2). Puzzle Hi-C contains three steps: mapping, scaffolding, and building.

195  Mapping uses the Juicer software(31), which filters out duplicate, abnormal alignment,
196  and restriction site information. scaffolding adopt an iterative method to obtain accurate
197  assemblies via multiple iterations. Finally, the genome is assembled according to the
198  ordering and orientation results and output in final fasta and apg format files.

199

200  **Fig 1. Puzzle Hi-C Pipeline.** The Puzzle Hi-C pipeline contains three steps: mapping, scaffolding,
201  and building. Ordering and orientation adopt an iterative method to obtain accurate assemblies via
202  multiple iterations. Puzzle Hi-C introduces a dynamic, triangle window strategy during assembling.
203  The triangle window is initially small and expands with interactions to produce more effective
204  clustering. Finally, the genome is assembled according to the scaffolding results and output in final
205  fasta and apg format files.

206

207  **Fig 2. Contact probability and strategies to evaluate distance between scaffolds adopted by**
208  **different software. a**, ideally distribution of Hi-C contact, Hi-C distribution in line with 1/x. Heat
209  map shows the diagonal position of interaction density is very high, the farther away from the
210  diagonal, the lower the interaction density is. b, Heat map shows the chr2 [0-35MB] assembled by
211  LACHESIS, where c1, c2, c3 and c4 represent scaffolds and the rectangle represents the number of
212  Hi-C reads links with two scaffolds. Due to Compartment and TADs, there are many long-range
213  interactions, which make long range interaction densities is higher than adjacent interaction
214  densities. **c**, strategy to evaluate distance adopted by LACHESIS and ALLHiC. **d**, strategy to
215  evaluate distance adopted by 3D DNA. **e**, strategy to evaluate distance adopted by SALSA. **f**,
216  strategy to evaluate distance adopted by Puzzle Hi-C. **g**, the CV of interaction density with triangle
217  region or square region. **h**, errors of the distance between two scaffolds with different gap size in
218  Puzzle Hi-C. **i**, errors while different strategies to evaluate distance between two scaffolds with
219  1000 samplings.

220

221  **Evaluation using simulated and real data**
222  To compare the performance of Puzzle Hi-C and other software, we used the
223  human genome hg38. We assessed the autosomes using lengths of 200 kb, 600 kb and
224  1 Mb contigs.
225  First, we obtained statistics on the assembly results of different Hi-C scaffolding
226  software in 200 kb, 600 kb, 1 Mb scaffolds (S2-6 Table). For example, LACHESIS
227  assembled a genome size of 2.77 Gb and contained 102, 37 and 33 scaffolds,
228  respectively. Scaffold N50s tended to be stable at about 135 Mb, while scaffold N90s
229  were 79.8 Mb, 68.8 Mb, and 77.5 Mb, respectively (Supplementary Table 2). SALSA2
230  assembled scaffolds of 1193, 593, and 538, respectively. Scaffold N50s were 8.6 Mb,
231  9.3 Mb, and 10.0 Mb, respectively (S3 Table). The assembled scaffolds were relatively
232  short and the clustering effect was not ideal, with scaffold N90s of only 1.2 Mb, 2.5
233  Mb, and 2.8 Mb, respectively. The genome assembled by Puzzle Hi-C was quite similar
234  to that of LACHESIS. The assembled scaffolds were 703, 191, and 109, respectively,
235  and scaffold N50s were 128.5 Mb, 130.9 Mb, and 154.4 Mb, respectively (S6 Table).
236  Scaffold N90s were 44.8 Mb, 55.6 Mb, and 56.6 Mb, respectively. Genome assembly
237  size, scaffold N50, and scaffold N90 can only reflect the clustering effect of Hi-C-based
238  scaffolding software. Because LACHESIS clusters before assembly, the Scaffold N50

239　result was excellent.

240　　　Second, to evaluate and compare the performance of existing scaffolding software
241　in scaffolding and orientation, we used dnadiff to compare the quality of the genomes.
242　We assessed three features: the number of relocations as determined by the number of
243　breaks in the alignment of scaffolds belonging to the same chromosome, but not
244　consistently ordered; the number of translocations, that being the number of breaks in
245　the alignment of scaffolds belonging to different chromosomes; and the number of
246　inversions, or breaks in the alignment by scaffolds inverted with respect to one another.
247　As the size of scaffolds got smaller, the proportion of assembly errors in LACHESIS
248　increased (Table 1). It produced 69 assembly errors in the 1 Mb scaffold size, including
249　9 translocations, 31 orientation assembly errors, and 29 relocations; at the 600kb
250　scaffold size, assembly errors increased to 132, and the 200kb size had 380 errors,
251　which showed an inverse relationship between scaffold size and assembly errors. Other
252　software showed the same pattern. Comparatively, Puzzle Hi-C consistently achieved
253　the greatest assembly accuracy under different sizes of scaffolds (Fig 3a-c, Table 1),
254　for example having 67 assembly errors at the 1 Mb scaffold size (relocations 11,
255　translocations 8, inversions 48).

256

257　　　　　**Table. 1 Statistical errors generated by different software with different**
258　　　　　　　　　　　　　　　　　　　**contig length**

| Contig Length | Error type | LACHESIS | SALSA2 | 3D-DNA | ALLHiC | Puzzle Hi-C |
|---|---|---|---|---|---|---|
| **200k** | Relocations | 152 | 188 | 107 | 226 | 11 |
| | Translocations | 46 | 48 | 48 | 69 | 8 |
| | Inversions | 182 | 387 | 379 | 351 | 48 |
| **600k** | Relocations | 43 | 170 | 24 | 67 | 3 |
| | Translocations | 24 | 14 | 27 | 23 | 9 |
| | Inversions | 65 | 264 | 124 | 59 | 25 |
| **1M** | Relocations | 29 | 129 | 22 | 31 | 2 |
| | Translocations | 9 | 16 | 22 | 17 | 8 |
| | Inversions | 31 | 236 | 69 | 47 | 10 |

259

260　**Fig 3. Statics different errors generated by different software with different**
261　**scaffold size. a-c,** inversions, relocations and translocations generated by LACHESIS, SALSA,
262　3D DNA, ALLHiC and Puzzle Hi-C under different length of Scaffolds. **d-f,** inversions, relocations
263　and translocations generated by LACHESIS, SALSA, 3D DNA, ALLHiC and Puzzle Hi-C with 25
264　sampling data under different length of Scaffolds.

265

266　　　Third, we resampled the human genome 25 times with five different scaffold sizes.
267　Puzzle Hi-C outperformed the other software packages (Fig 3d-f). For example, it
268　produced the best chromatin assembly, showing much less assembly errors. Puzzle Hi-
269　C was not affected by the size of the scaffolds, and it was more robust (Fig 3d-f).

270　　　To test the assembly performance of Puzzle Hi-C using real data, we also
271　employed the scaffold version of the human genome assembly (version:

272   GCA_001013985.1). This analysis used LACHESIS. We compared the assemblies to
273   the human genome GRCh38 using MuMmer software. LACHESIS produced 999 errors
274   in its ordering and orientation of large fragments in assembly, and Puzzle Hi-C gave
275   647 errors, except for chromosome 1, which was composed of three scaffolds. In
276   addition to assembly errors of large fragments, LACHESIS also had more problems in
277   assembling small scaffolds, such as chromosomes 17, 19, 20, and 22. Puzzle Hi-C did
278   not exhibit this problem (Fig 4, Table 2). Other methods also showed more errors than
279   Puzzle Hi-C (Table 2). We also assembled other species genomes across a range of
280   taxonomic groups, genome sizes and initial assembly quality. Puzzle Hi-C consistently
281   generated assemblies with higher contiguity (Supplementary Figure 1).
282
283   **Fig 4. The synteny of chromosomes assembled by Puzzle Hi-C and LACHESIS**
284   **compared with GRCh38.**
285
286   **Table. 2 Statistical errors generated by different software with**
287   **GCA_001013985.1**

| Error type | LACHESIS | SALSA2 | 3D-DNA | ALLHiC | Puzzle Hi-C |
|---|---|---|---|---|---|
| **Relocations** | 529 | 1076 | 4632 | 4188 | 243 |
| **Translocations** | 116 | 420 | 408 | 2526 | 149 |
| **Inversions** | 354 | 714 | 8572 | 2908 | 255 |

288   **Assembling genomes enriched with long-range interactions**
289   To further test the robustness of assembly by Puzzle Hi-C, we employed
290   chromosome 2 of gayal, which contains a Robertsonian translocation. This
291   chromosome has more repetitive sequences and long-range interactions than its
292   relatives. Long-range interaction will obstruct ordering prediction. The Hi-C interaction
293   matrix revealed a very strong internal interaction of compartments on chromosome 2,
294   which indicated a long-distance interaction. Other software assemblies also detected
295   the rearrangement of large fragments, but Puzzle Hi-C scaffolding software obtained
296   relatively fewer chromatin orientation assembly errors. Therefore, Puzzle Hi-C
297   appeared to best assemble chromosomes when chromatin interactions occurred (Fig 5).
298
299   **Fig 5. The synteny of chromosomes assembled by Puzzle Hi-C and other software**
300   **compared with gayal chromosome 2.**
301
302   **Genome's quality is curial for 3D analysis**
303   A high-quality genome is essential for downstream analysis. However, due to the
304   absence of reliable tools, chromosome-level assemblies may contain some error. These
305   errors may directly affect the main results of the analysis. To estimate the effect of
306   chromosome errors on downstream analysis, we downloaded the genome of *Takifugu*
307   *bimaculatus*(32), which was assembled by LACHESIS. Compared to the genome of *T.*
308   *rubripes*, the genome of *T. bimaculatus* has 809 inversions and 2618 relocations. We
309   reassembled this genome using Puzzle Hi-C and obtained 519 inversions and 1791
310   relocations. We performed comparable analysis on both old and new genomes. The

311 results showed that different genome assembles will affect comparative analysis (Fig
312 6).

313

314 **Fig 6. The scaffolding result of LACHESIS and Puzzle Hi-C on *T. bimaculatus*. a**, the syteny
315 between *T. bimaculatus* and *T. rubripes*; **b**, the syteny between Puzzle Hi-C corrected *T.*
316 *bimaculatus* and *T. rubripes*; **c**, i *T. bimaculatus* genome chr1 Hi-C heat map, the black box is the
317 Hi-C heat map suggesting assembly error; ii *T. bimaculatus* genome chr1 Compartment, red is
318 Compartment A and blue is Compartment B; iii is the rearrangement of *T. bimaculatus* genome chr1
319 Compartment according to the corrected chr1; iv Puzzle Hi-C corrected chr1 Compartment after
320 Puzzle Hi-C correction; v Puzzle Hi-C corrected chr1 Hi-C heat map.

321

## Discussion and Conclusion

323 Hi-C data facilitate the assembly of chromosome-level genomes by locating and
324 avoiding long-range interactions. Puzzle Hi-C uses a dynamics triangle window to
325 calculate interaction densities. It dynamically changes the size of the triangle at both
326 ends of the scaffold. Windows start small in initial iterations, which facilitates the
327 assembling of smaller scaffolds. Puzzle Hi-C excludes long-range interactions when
328 the window is small. While LACHHESIS(17), 3D DNA(25), and ALLHiC(27) use all
329 interaction information between two scaffolds, such can result in errors in ordering due
330 to long-range interactions. As iterations increase in Puzzle Hi-C, the window increases
331 in size, therefore obtaining better chromosome clustering by selectively using all
332 interaction information. Such avoids the problem of failing to cluster scaffolds into
333 chromosomes when using only local interaction information, which SALSA(26,28)
334 does. Puzzle Hi-C evaluation on human genome outperforms ALLHiC(27),
335 LACHESIS(17), 3D DNA(25), and SALSA(26,28) in ordering and orientation in both
336 simulated and real data, and with robust performance. The same result occurs upon
337 applying Puzzle Hi-C to all other tested genomes. Further, Puzzle Hi-C outperforms
338 other software when assembling the complex gayal genome, which has many long-
339 range interactions. Similarly, the reassembled genome of the puffer fish reveals
340 improvements when compared with the original assembly(32). Finally, the results
341 suggest that the genome-quality greatly impacts 3D genome analysis. Thus, accurate
342 3D genome analysis requires accurate chromosome-level genomes.

343

## Data and Code Availability

345 All Hi-C data were downloaded from NCBI (S1Table). The Puzzle Hi-C software
346 package with a detailed user tutorial and sample input and output files can be found at
347 https://github.com/linguoliang/puzzle-hi-c.git.

# Acknowledgments

# References

1. Kaul A, Bhattacharyya S, Ay F. Identifying statistically significant chromatin contacts from Hi-C data with FitHiC2. Nat Protoc. 2020 Mar;15(3):991–1012.

2. Wang XT, Cui W, Peng C. HiTAD: Detecting the structural and functional hierarchies of topologically associating domains from chromatin interactions. Nucleic Acids Res. 2017;45(19).

3. Murray GGR, Soares AER, Novak BJ, Schaefer NK, Cahill JA, Baker AJ, et al. Natural selection shaped the rise and fall of passenger pigeon genomic diversity. Science [Internet]. 2017 Nov 17 [cited 2022 Jan 31]; Available from: https://www.science.org/doi/abs/10.1126/science.aao0960

4. O'Connor RE, Romanov MN, Kiazim LG, Barrett PM, Farré M, Damas J, et al. Reconstruction of the diapsid ancestral genome permits chromosome evolution tracing in avian and non-avian dinosaurs. Nat Commun. 2018;9(1):1883.

5. Sagai T, Hosoya M, Mizushina Y, Tamura M, Shiroishi T. Elimination of a long-range cis-regulatory module causes complete loss of limb-specific Shh expression and truncation of the mouse limb. Development. 2005 Feb 15;132(4):797–803.

6. Pombo A, Dillon N. Three-dimensional genome architecture: players and mechanisms. Nat Rev Mol Cell Biol. 2015 Apr;16(4):245–57.

7. Mishra A, Hawkins RD. Three-dimensional genome architecture and emerging technologies: Looping in disease. Genome Med. 2017;9(1):1–14.

8. Choy MK, Javierre BM, Williams SG, Baross SL, Liu Y, Wingett SW, et al. Promoter interactome of human embryonic stem cell-derived cardiomyocytes connects GWAS regions to cardiac gene networks. Nat Commun. 2018 Jun 28;9(1):1–10.

9. Pan DZ, Garske KM, Alvarez M, Bhagat YV, Boocock J, Nikkola E, et al. Integration of human adipocyte chromosomal interactions with adipose gene expression prioritizes obesity-related genes from GWAS. Nat Commun. 2018 Apr 17;9(1):1512.

10. Xu Z, Zhang G, Duan Q, Chai S, Zhang B, Wu C, et al. HiView: an integrative genome browser to leverage Hi-C results for the interpretation of GWAS variants. BMC Res Notes. 2016 Mar 11;9(1):159.

11. Lu L, Liu X, Huang WK, Giusti-Rodríguez P, Cui J, Zhang S, et al. Robust Hi-C Maps of Enhancer-Promoter Interactions Reveal the Function of Non-coding

388    Genome in Neural Development and Diseases. Mol Cell. 2020 Aug 6;79(3):521-
389    534.e15.

390    12. Lu B, Jiang J, Wu H, Chen X, Song X, Liao W, et al. A large genome with
391    chromosome-scale assembly sheds light on the evolutionary success of a true toad
392    (Bufo gargarizans). Mol Ecol Resour. 2021;21(4):1256–73.

393    13. Heesch S, Cho GY, Peters AF, Corguillé GL, Falentin C, Boutet G, et al. A
394    sequence-tagged genetic map for the brown alga Ectocarpus siliculosus provides
395    large-scale assembly of the genome sequence. New Phytol. 2010;188(1):42–51.

396    14. Yu Q, Tong E, Skelton RL, Bowers JE, Jones MR, Murray JE, et al. A physical
397    map of the papaya genome with integrated genetic map and genome sequence.
398    BMC Genomics. 2009 Aug 7;10(1):371.

399    15. Wu P, Zhou C, Cheng S, Wu Z, Lu W, Han J, et al. Integrated genome sequence
400    and linkage map of physic nut (Jatropha curcas L.), a biodiesel plant. Plant J.
401    2015;81(5):810–21.

402    16. Rice ES, Green RE. New Approaches for Genome Assembly and Scaffolding.
403    Annu Rev Anim Biosci. 2019;7(1):17–40.

404    17. Burton JN, Adey A, Patwardhan RP, Qiu R, Kitzman JO, Shendure J.
405    Chromosome-scale scaffolding of de novo genome assemblies based on chromatin
406    interactions. Nat Biotechnol. 2013 Dec;31(12):1119–25.

407    18. Mascher M, Gundlach H, Himmelbach A, Beier S, Twardziok SO, Wicker T, et al.
408    A chromosome conformation capture ordered sequence of the barley genome.
409    Nature. 2017 Apr;544(7651):427–33.

410    19. Oddes S, Zelig A, Kaplan N. Three invariant Hi-C interaction patterns:
411    Applications to genome assembly. Methods. 2018 Jun 1;142:89–99.

412    20. Cremer T, Cremer M. Chromosome Territories. Cold Spring Harb Perspect Biol.
413    2010 Jan 3;2(3):a003889.

414    21. Meaburn KJ, Misteli T. Chromosome territories. Nature. 2007 Jan;445(7126):379–
415    81.

416    22. Sun HB, Shen J, Yokota H. Size-Dependent Positioning of Human Chromosomes
417    in Interphase Nuclei. Biophys J. 2000 Jul 1;79(1):184–90.

418    23. Lieberman-Aiden E, Van Berkum NL, Williams L, Imakaev M, Ragoczy T, Telling
419    A, et al. Comprehensive mapping of long-range interactions reveals folding
420    principles of the human genome. Science. 2009;326(5950):289–93.

421    24. Lajoie BR, Dekker J, Kaplan N. The Hitchhiker's guide to Hi-C analysis: Practical
422    guidelines. Methods. 2015 Jan 15;72:65–75.

423    25. Dudchenko O, Batra SS, Omer AD, Nyquist SK, Hoeger M, Durand NC, et al. De
424    novo assembly of the Aedes aegypti genome using Hi-C yields chromosome-length
425    scaffolds. Science. 2017 Apr 7;356(6333):92–5.

426    26. Ghurye J, Pop M, Koren S, Bickhart D, Chin CS. Scaffolding of long read
427    assemblies using long range contact information. BMC Genomics. 2017 Jul
428    12;18(1):527.

429    27. Zeng W, Wu M, Jiang R. Prediction of enhancer-promoter interactions via natural
430    language processing. BMC Genomics. 2018;19(Suppl 2).

431    28. Ghurye J, Rhie A, Walenz BP, Schmitt A, Selvaraj S, Pop M, et al. Integrating Hi-C

links with assembly graphs for chromosome-scale assembly. PLOS Comput Biol. 2019 Aug 21;15(8):e1007273.

29. Wang S, Wang H, Jiang F, Wang A, Liu H, Zhao H, et al. EndHiC: assemble large contigs into chromosome-level scaffolds using the Hi-C links from contig ends. BMC Bioinformatics. 2022 Dec 8;23(1):528.

30. Zhou C, McCarthy SA, Durbin R. YaHS: yet another Hi-C scaffolding tool [Internet]. bioRxiv; 2022 [cited 2022 Jul 29]. p. 2022.06.09.495093. Available from: https://www.biorxiv.org/content/10.1101/2022.06.09.495093v2

31. Durand NC, Shamim MS, Machol I, Rao SSP, Huntley MH, Lander ES, et al. Juicer Provides a One-Click System for Analyzing Loop-Resolution Hi-C Experiments. Cell Syst. 2016;3(1):95–8.

32. Zhou Z, Liu B, Chen B, Shi Y, Pu F, Bai H, et al. The sequence and de novo assembly of Takifugu bimaculatus genome using PacBio and Hi-C technologies. Sci Data. 2019 Sep 30;6(1):187.

33. Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M, Antonescu C, et al. Versatile and open software for comparing large genomes. Genome Biol. 2004 Jan 30;5(2):R12.

# Supporting information

S1 Text. Puzzle Hi-C pipeline.
S1 Fig. Hi-C contact maps of assemblies constructed from puzzle hic.
S1 Table. SRA data used in this study.
S2 Table. LACHESIS scaffolding results.
S3 Table. SALSA2 scaffolding results.
S4 Table. 3D DNA scaffolding results.
S5 Table. ALLHiC scaffolding results.
S6 Table. Puzzle Hi-C scaffolding results.
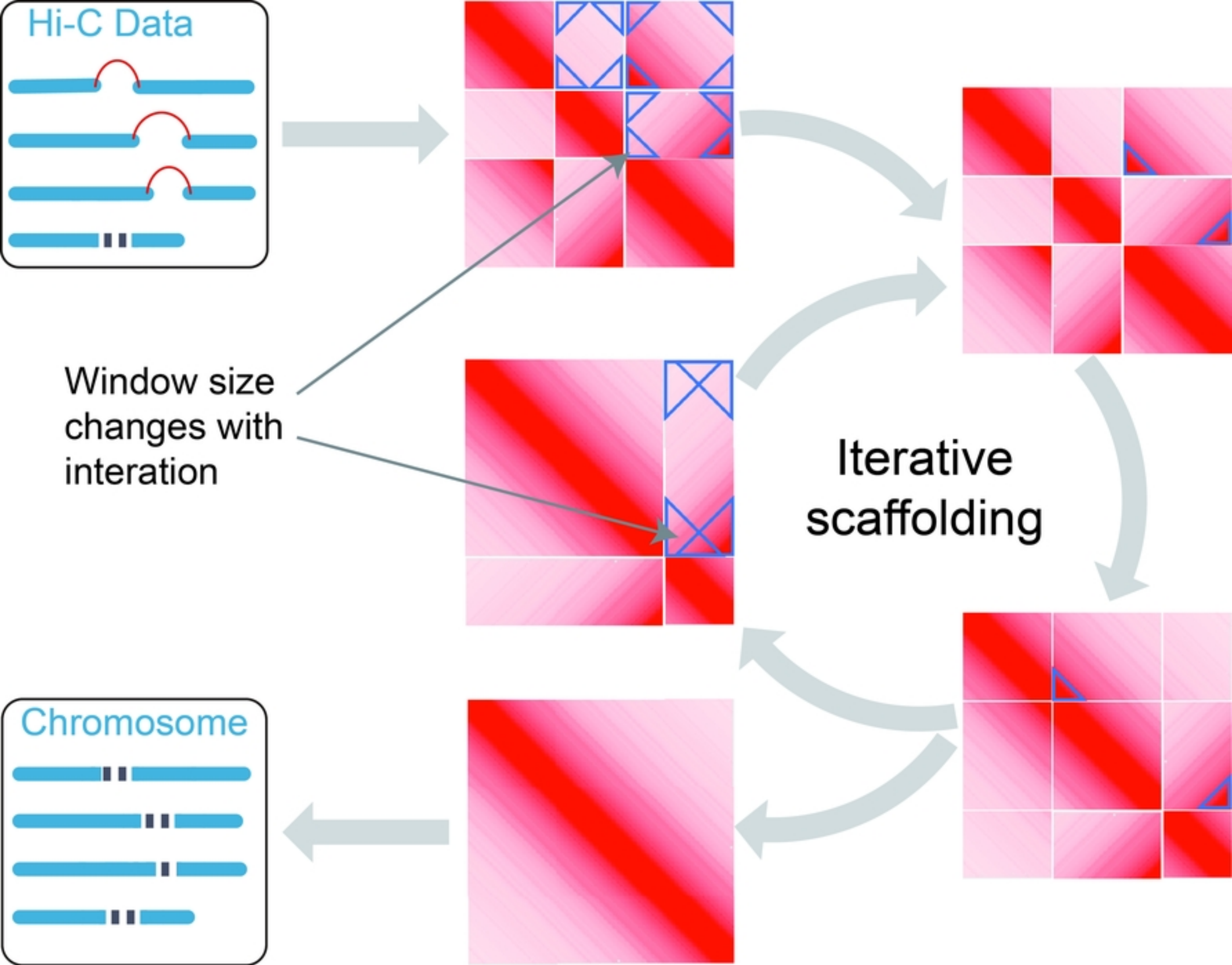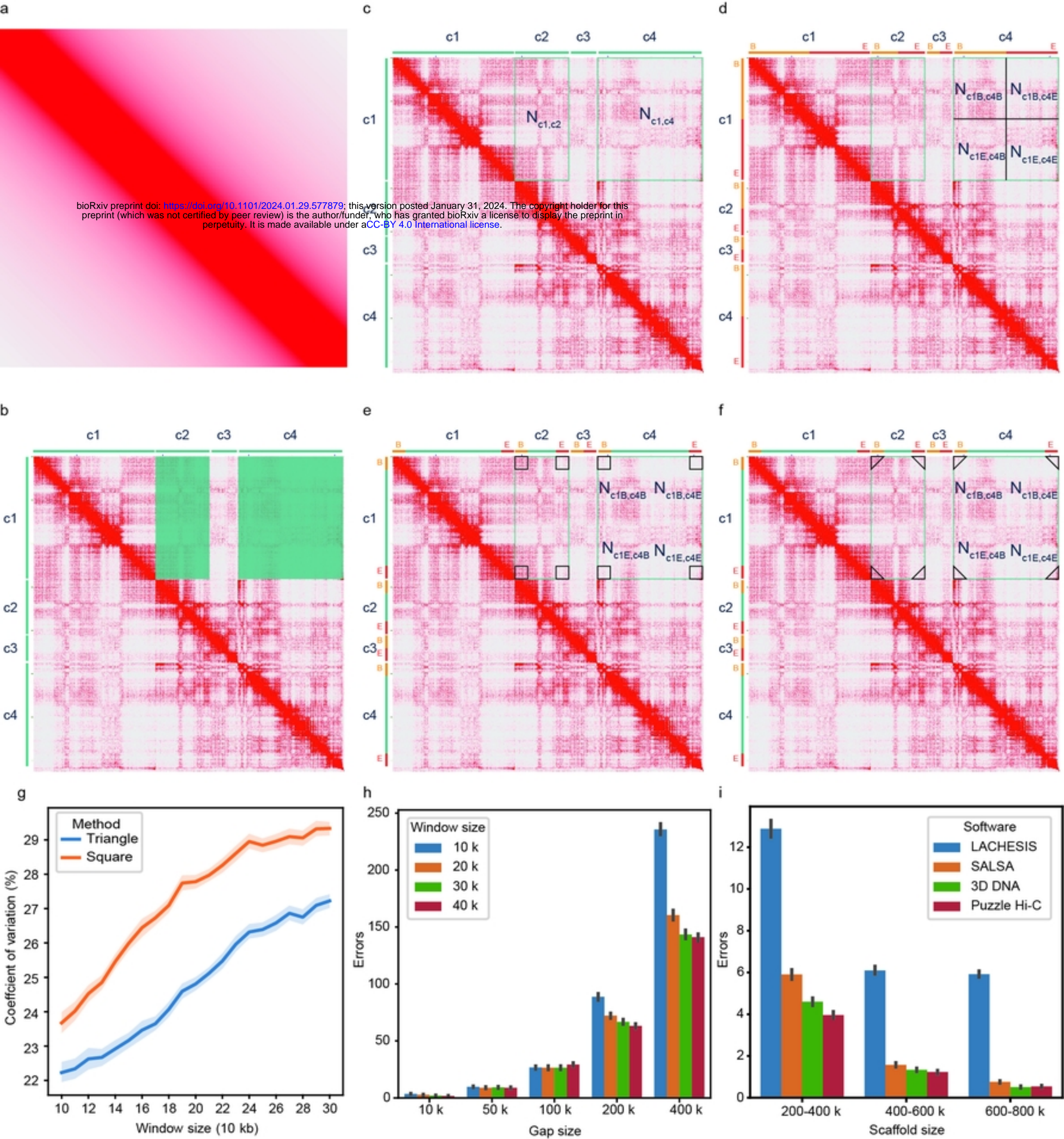
**Hi-C Data**

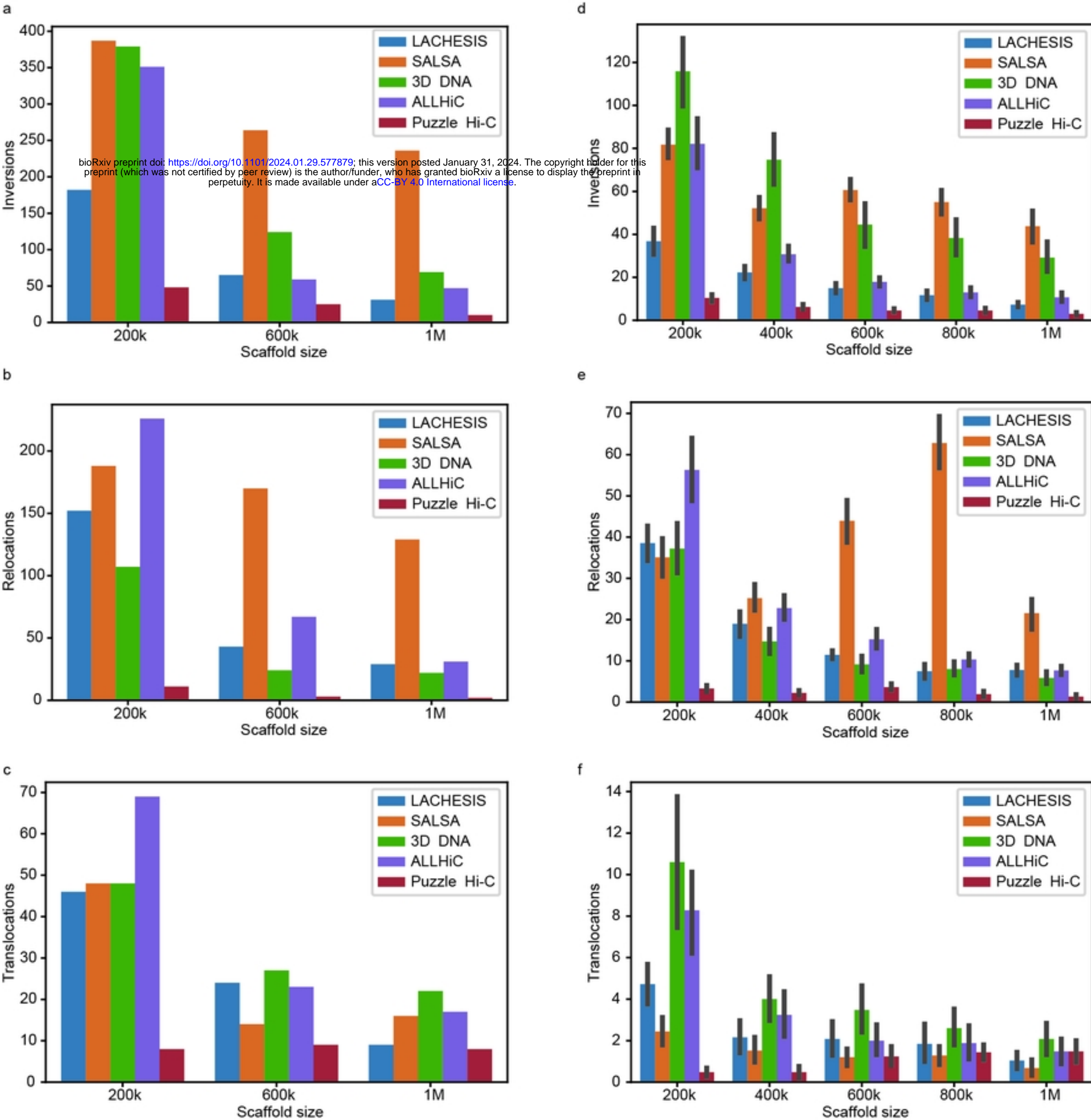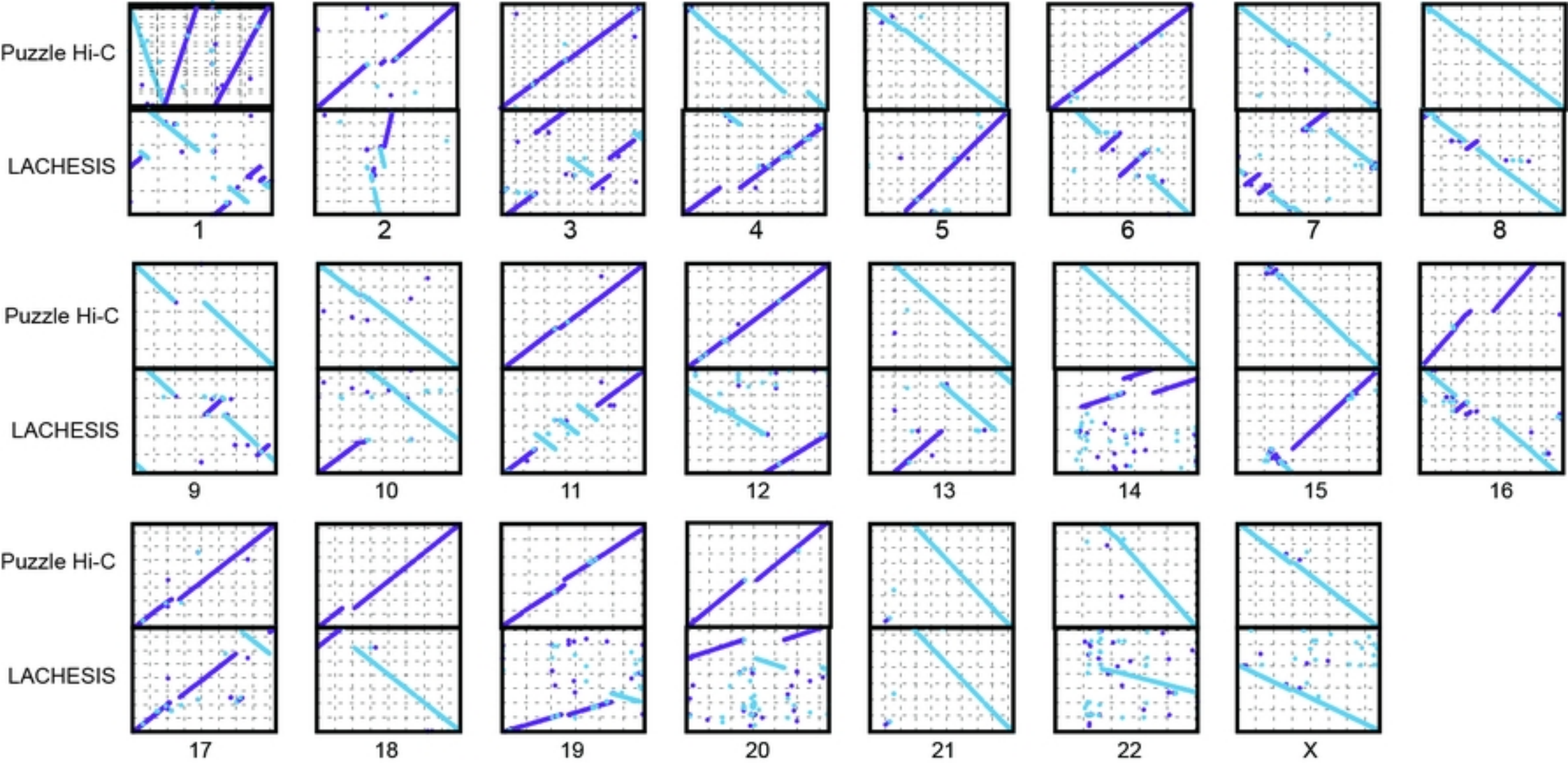Window size changes with interation

**Chromosome**

Iterative scaffolding

Fig. 1

Fig. 2

Fig. 3

Fig. 4

Fig. 5

a

T. rubripes

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 1819 20 21 22

T. bimaculatus assembled by LACHESIS

b

T. rubripes

1 2 3 4 5 6 7 8 9 1011 1213 14 15 1617 1819 20 21 22

T. bimaculatus assembled by Puzzle Hi-C

c

i          ii   iii    iv  v

T. bimaculatus chr1 assembled by LACHESIS

T. bimaculatus chr1 assembled by Puzzle Hi-C

# Fig. 6