

1 **Supplementary Information for Marie-Nelly et al.**

2

3 **Supplementary Data 1. Most likely genome structure for the Malaysian yeast strain after**
4 **47,880 iterations**

5 This file generated by GRAAL recapitulates the correspondence between the genome used for the
6 initialization of GRAAL and the most likely genomic structure recovered at the end of the
7 process. Superscaffolds generated by GRAAL are indicated in the first column, with the
8 corresponding bin from the original genome indicated below each superscaffold under the
9 “init_published_scaffold” label. The index, orientation, and initial coordinates of each bin within
10 the initial genome sequence are also indicated.

11

12 **Supplementary Data 2. Most likely genome structure for the *T. reesei* strain QM6A after**
13 **31,920 iterations**

14 This file generated by GRAAL recapitulates the correspondence between the genome used for
15 initializing the algorithm and the most likely genomic structure recovered. Superscaffolds
16 generated by GRAAL are indicated in the first column, with the corresponding “bin” from the
17 original genome indicated below each superscaffold under the “init_published_scaffold” label.
18 The index, orientation, and initial coordinates within the initial genome sequence are also
19 indicated.

20

21 **Supplementary Data 3. Fasta file of the most likely genome structure of the UWOPS03-**
22 **461.4 Malaysian yeast strain after 47,880 iterations**

23

24 **Supplementary Data 4. Fasta file of the most likely genome structure of the *T. reesei* strain**
25 **QM6A after 31,920 iterations**

26

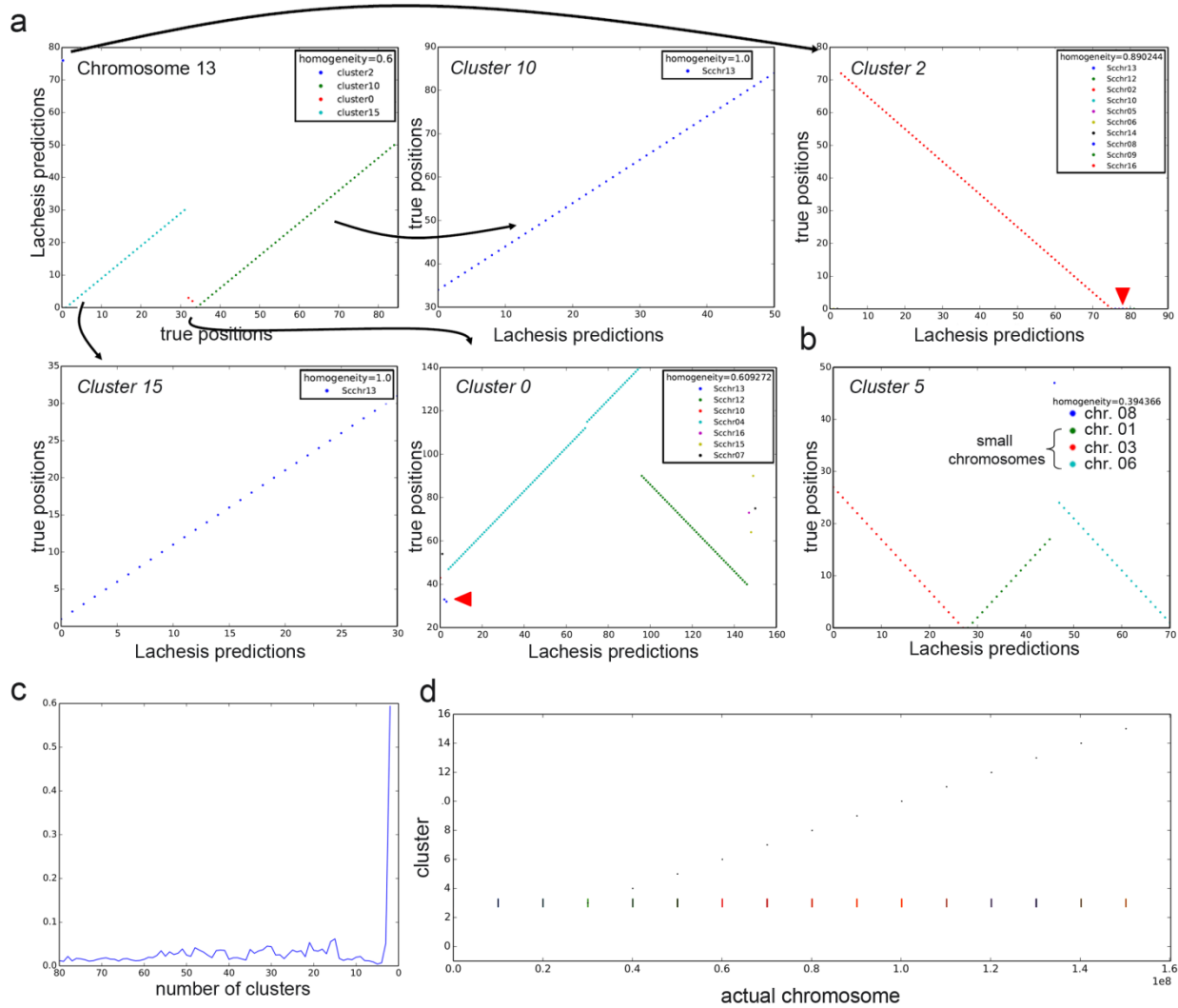
27 **Supplementary Data 5.** List of the 2,917 *de novo* contigs of chromosome 14 from sequencing
28 libraries downloaded from the GAGE competition website used for initializing GRAAL.

29

30 **Supplementary Data 6.** List of the 8,382 bins generated from these 2,917 contigs from
31 Supplementary Data 5.

32

33



34

35 **Supplementary Figure 1:** Assembly of virtual *S. cerevisiae* contigs using Lachesis¹ and dnaTri².

36 (a) Example of inaccurate clustering by Lachesis, starting with the set of bins assembled by

37 GRAAL in Fig. 2. Two large chromosomal segments of chromosome 13 were attributed to

38 clusters 10 and 15, whereas two small regions of the same chromosome were incorporated in

39 clusters 0 and 2 (red arrowheads). (b) Example of inaccurate clustering by Lachesis of small

40 chromosomes 1, 3, and 6 into a single cluster. Note that although the 95% of the bins are

41 correctly aligned with respect to their neighbors, this measure does not reflect the overall quality

42 of the assembly. (c) dnaTri also fails to retrieve the correct number of yeast chromosomes when

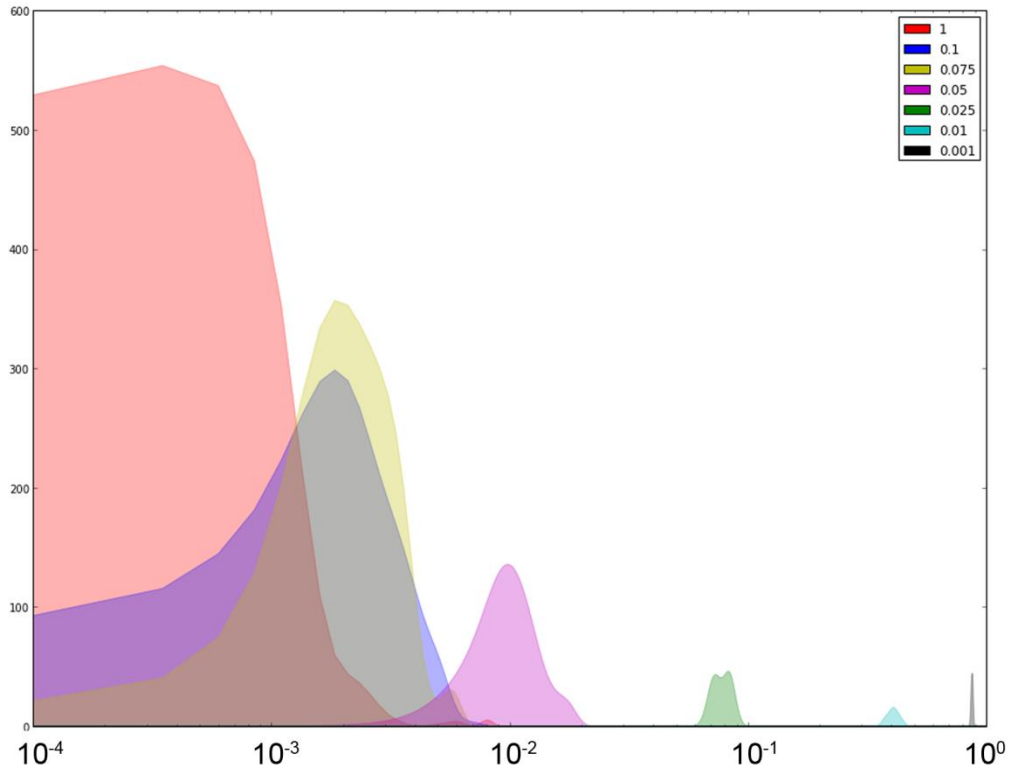
43 applied to yeast contact data. Both plots were generated by dnaTri. The left plot shows the

44 average clustering step length as function of the number of clusters tested (see Figure 3a of the

45 dnaTri paper²). The number of clusters chosen by dnaTri corresponds to the maximum of this

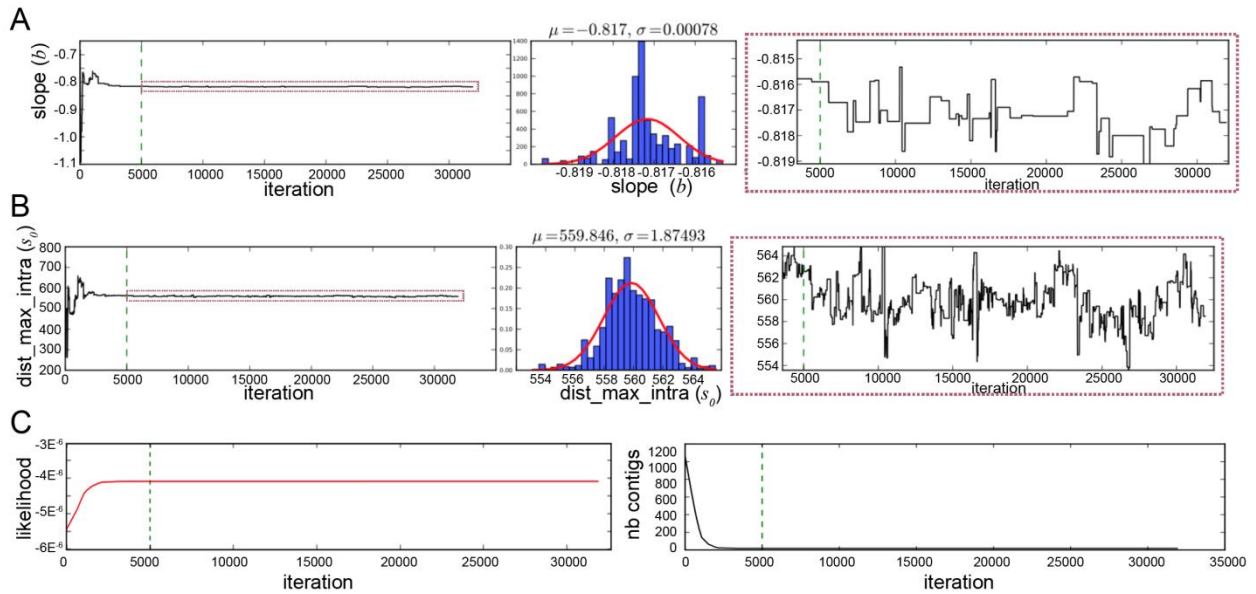
46 graph and in this case equals 2 instead of the expected number of 16 chromosomes. The right plot
47 shows the assignment of contigs from the 16 chromosomes to clusters, revealing that the vast
48 majority was incorrectly grouped into a single cluster.

49



50
 51 **Supplementary Figure 2:** Distribution of the error rate for sets of randomly down-sampled 3C
 52 dataset (from 1X to 0.001X for the *S. cerevisiae* matrix containing 21,457,486 contacts). The x-
 53 axis represents the error rate in log scale. For each down-sampled dataset, 15,000 iterations were
 54 performed.

55
 56



57
58

59 **Supplementary Figure 3: Evolution of the parameters** of the model (a) The slope reflects the
 60 intrachromosomal contact frequencies as a function of genomic separation (i.e. nuisance
 61 parameter b in the model; Material and Methods), repeatedly revisited over the 50,000 iterations
 62 in light of contact data. (b) Dist_max_intra represent the threshold, in kb, allowing discrimination
 63 between intra- and inter-chromosomal contacts, with inter-chromosomal frequencies assumed
 64 constant, corresponding to nuisance parameter s_0 . Both the slope and the Dist_max_intra values
 65 are repeatedly reassessed based on the 3D data, fluctuating around an average value (μ), as
 66 illustrated by the close-ups (red (red dotted squares) on each curve. (c) Evolution of the
 67 likelihood and the number of contigs as function of the number of iterations.

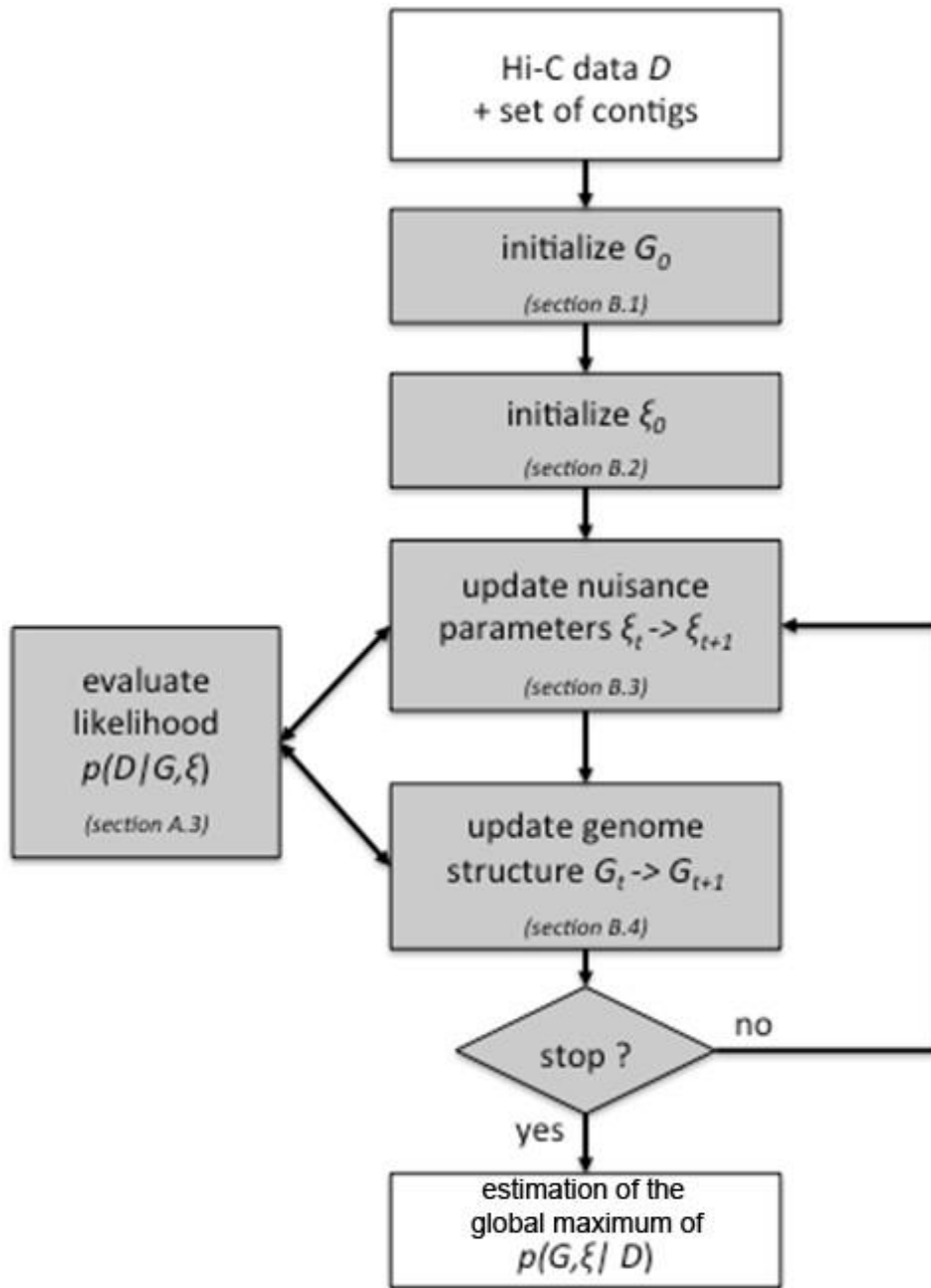
68

69

70

71

72



73

74

75 **Supplementary Figure 4: Overview of the GRAAL algorithm.**

76

77

78 **Supplementary Table 1. Features of three genome assembly algorithms based on 3D contact**
 79 **data**

80

Feature	GRAAL	Lachesis ¹	dnaTri ²
Predicts number of chromosomes	Yes	No	Yes
Corrects automatically initial misassemblies	Yes	No	No*
Orients contigs	Yes	Yes	No
Identifies repeated regions	Yes	No	No
Estimates assembly uniqueness	Yes	No	No**

81

82 * not directly: the user can still cut the initializing contigs before the clustering step

83

84 ** the probabilistic framework of the dnaTri algorithm is very elegant and allows it to estimate
 85 the likelihood of the structure, but, as acknowledged by the authors, there is no guarantee on the
 86 global optimum of the solution.

86

87

88 **Supplementary Table 2: Sequencing adapters used in this study**

oligos	sequence	library
MM70	GTANNNNNNAGATCGGAAGAGCGGTTTCAGCAGGAATGCCGAG	Malaysian yeast strain
MM71	ACACTCTTCCCTACACGACGCTCTTCCGATCTNNNNNNNTACT	
MM182	TCTNNNNNNAGATCGGAAGAGCGGTTTCAGCAGGAATGCCGAG	YKF1246 yeast strain
MM183	ACACTCTTCCCTACACGACGCTCTTCCGATCTNNNNNNNAGAT	
MM106	TGGNNNNNNAGATCGGAAGAGCGGTTTCAGCAGGAATGCCGAG	QM6a <i>T. reesei</i> strain
MM107	ACACTCTTCCCTACACGACGCTCTTCCGATCTNNNNNNCCAT	
MM108	CCANNNNNNAGATCGGAAGAGCGGTTTCAGCAGGAATGCCGAG	<i>E. histolytica</i> strain
MM109	ACACTCTTCCCTACACGACGCTCTTCCGATCTNNNNNNNTGGT	

89

90

91

92

93 **Supplementary Table 3. Summary of the initialization parameters for the different analysis**

Dataset	enzyme	nb of bins	n contacts	Mean nb of RFs per bin	mean bin size (kb)
<i>S. cerevisiae</i>	<i>DnpII</i>	1086	21457086	27	11
<i>Trichoderma reesei</i> QM6a	<i>DnpII</i>	1193	15014468	27	27,9
YKF1246 <i>S.c.</i> strain	<i>DnpII</i>	1295	21830579	27	9,5
Malaysian <i>S.c.</i> strain	<i>DnpII</i>	3136	8353283	9	3,8
Human chr7/17/19/22	<i>HindIII</i>	3607	19672219	9	95,9
Human <i>de novo</i> chr14	<i>HindIII</i>	8382	1156115	3	8,8

94

95

96

97

98 **Supplementary Method**

99

100 The following provides a more detailed description of the algorithm implemented in GRAAL.

101

102 The description of GRAAL can be divided into two main components: (i) the probabilistic model
103 that assigns a likelihood to a given linear (one-dimensional) genome structure given a specific
104 contact/Hi-C data set, and (ii) the sampling algorithm used to explore the space of linear genome
105 structures (and nuisance parameters).

106

107 **A. Probabilistic model**

108

109 **A.1. Bayesian inference approach:**

110 We consider the genome assembly problem as a Bayesian inference problem, taking inspiration
111 from previous work in protein structure determination³. In its simplest form, the Bayes rule reads:

$$p(G|D) \propto p(D|G)p(G)$$

112 where G denotes the linear genome structure to be determined, D is the Hi-C data set (both will
113 be defined more precisely below), $p(A|B)$ is the probability density of A conditioned on B , and \propto
114 indicates proportionality. Our goal is to determine, or at least approximate, the posterior
115 probability $p(G|D)$. The above formula provides a means to compute this probability density (up
116 to a normalizing factor) given a probabilistic data generation model, $p(D|G)$ (called likelihood)
117 and data-independent assumptions about the structure, encapsulated by the prior probability
118 $p(G)$.

119

120 In practice, our data generation model involves several parameters (called nuisance parameters)
121 that are not known *a priori* (see below). Therefore, we include these parameters, collectively
122 noted as ξ , in the Bayesian formulation, yielding:

$$p(G, \xi|D) \propto p(D|G, \xi)p(G, \xi) = p(D|G, \xi)p(G)p(\xi)$$

123 where for the latter identity we assumed statistical independence of the genome structures G and
124 the nuisance parameters ξ .

125

126 We next assume that in absence of data, all possible genome structures and nuisance parameters
 127 are equally probable, i.e. that $p(G)$ and $p(\xi)$ are constants (flat priors). With these assumptions
 128 the Bayes rule reduces to:

$$p(G, \xi|D) \propto p(D|G, \xi)$$

129
 130 To compute the likelihood $p(D|G, \xi)$ we need a data generation model that relates the contact
 131 frequencies measured by the Hi-C experiment to an assumed linear genome structure and the
 132 nuisance parameters.

133
 134 A.2. Notations and definitions for the genome structure G and the Hi-C data D :
 135 Before describing our model for $p(D|G, \xi)$, we need a more formal definition of the variables G ,
 136 and D . The parameters ξ will be defined in section A.1.3.

137
 138 Genome structure:

139 First, we define G as an unordered set of N contigs \mathcal{C}_i :

$$G = \{\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_{N_c}\}$$

140 If the genome is perfectly assembled, each contig corresponds exactly to a single chromosome.
 141 Hi-C reads are mapped to restriction fragments μ_i defined by the restriction enzyme cutting sites.
 142 We therefore consider the restriction fragment μ_i as the elementary units of a genome assembly.
 143 However, many operations performed by GRAAL are not applied to individual restriction
 144 fragments, but to ordered sets of p consecutive fragments, which we call 'bins' and note :

$$f_k = (\mu_{k,1}, \dots, \mu_{k,p})$$

145 Whenever possible, we choose $p = p_m$ where p_m is a single user-defined constant typically set to
 146 3. However, if the number of restriction fragments in a contig is not a multiple of p_m , then some
 147 bins will consist of $p < p_m$ fragments.

148 We define a contig as an ordered sequence of bins, noted:

$$\mathcal{C}_k = (f_{\varphi_k^1}, f_{\varphi_k^2}, \dots, f_{\varphi_k^{L_k}})$$

149 where $\mathfrak{F}_G = \{f_1, f_2, \dots, f_n\}$ is the set of all bins, L_k is the number of bins in contig \mathcal{C}_k and φ_k is
 150 an indexing function with $\varphi_k^i \in [1, \dots, n]$. The subscript in \mathfrak{F}_G is used to indicate that the bins

151 rely on an initial assumed set of contigs G_0 . We also introduce the two functions $\psi_1(i)$ and $\psi_2(i)$
 152 such that:

$$\begin{cases} \psi_1(\varphi_k^i) = i \\ \psi_2(\varphi_k^i) = k \end{cases}$$

153 Next, we define $s_G(f_i, f_j)$ as the genomic distance (in units of base pairs) between two bins. This
 154 distance is obviously only defined for bins belonging to the same contig, i.e. for $\psi_2(i) = \psi_2(j)$.
 155 For $\psi_2(i) \neq \psi_2(j)$ we consider that $s_G(f_i, f_j) = \infty$.

156

157 Hi-C/3C-seq contact data:

158 The chromosome contact data used by GRAAL are obtained after mapping the Hi-C/3C-seq
 159 reads to an initial set of fragments $\mu_i, i = 1..N_f$. We define D as the matrix whose entries $D_{i,j}$,
 160 $\left((i, j) \in [1, \dots, N_f]^2 \right)$ are the number of Hi-C3C-seq reads pair mapped to each pair of fragments
 161 (μ_i, μ_j) . Please note that, although the sampling algorithm of GRAAL (described below)
 162 manipulates the genome structure at the level of super-contigs, the likelihood will always be
 163 evaluated by considering the contact data at the resolution of individual fragments.

164

165

166 A.3. Likelihood and nuisance parameters:

167 We now need a means to relate the probability of the matrix D to the assumed linear genome
 168 structure G . Our first step is to relate the probability of each entry $D_{i,j}$ to the contact probability
 169 between fragments μ_i and μ_j , which we note $q_{i,j}$. Since $D_{i,j}$ results from a counting process, its
 170 probability can be modeled as a Poisson distribution:

$$\begin{cases} P(D_{i,j} = n) = \text{Poisson}(n; \lambda_{i,j}) \stackrel{\text{def}}{=} \begin{cases} \frac{\lambda_{i,j}^n}{n!} e^{-\lambda_{i,j}} & \text{for } n \in \mathbb{N} \\ 0 & \text{otherwise} \end{cases} \\ \lambda_{i,j} = N_D q_{i,j} \end{cases}$$

171 where N_D is the total number of independent counts. Although N_D is strictly speaking also a
 172 random number (which depends chiefly on sequencing depth and criteria used to validate the read
 173 pairs), for simplicity we treat it as a constant and simply set: $N_D = \sum_{i=1}^{N_f} \sum_{j=i}^{N_f} D_{i,j}$. We further
 174 assume that contacts between distinct pairs of fragments are independent from each other, such
 175 that :

$$p(D) = \prod_{i=1}^n \prod_{j=i}^n p(D_{i,j})$$

176 which implies:

$$p(D) = \prod_{i=1}^n \prod_{j=i}^n \text{Poisson}(D_{i,j}; N_D q_{i,j})$$

177 In order to be able to calculate $p(D|G, \xi)$, we now need to relate the contact probabilities $q_{i,j}$ to
178 G and ξ .

179 Contact probabilities are intimately dependent on how chromosomes are folded and positioned
180 relative to each other. The detailed relationship between the 3D and 1D architecture of the
181 genome is in general complex, depends on the organism and cell state, and is subject of much
182 current research (e.g. ⁴). Nevertheless, some important features are present in all Hi-C/3C-seq
183 data sets obtained so far, and are also in good agreement with predictions from polymer physics.
184 Specifically, the contact probability P_{cis} between loci on the same chromosome (*cis* contacts)
185 decays as a power law with increasing genomic distances s (as expressed in bp). This relation
186 holds up to a genomic distance s_0 above which contact probabilities are approximately constant,
187 i.e.:

$$\begin{cases} P_{cis}(s) = P_t \left(\frac{s}{s_0}\right)^b & \text{for } s \leq s_0 \\ P_{cis}(s) = P_t & \text{for } s \geq s_0 \end{cases}$$

188 These relationships have been approximately verified in a number of different organisms, with
189 variable values for b , s_0 , and P_t , and can be recapitulated by computational simulations of
190 polymer dynamics ⁵⁻⁹.

191 The contact probabilities between loci on distinct chromosomes (*trans* contacts) are less
192 amenable to simple theoretical predictions and arguably more sensitive to biological specificities
193 such as organism and cell type. They are also on average much weaker than *cis* contact
194 probabilities. For simplicity and generality, we therefore simply assume that *trans* contacts have
195 the same probability as long-range *cis* contacts:

$$P_{trans} = P_t$$

196 Our probabilistic model is therefore characterized by only 3 parameters, collectively noted as ξ :

197
$$\xi = (P_t, s_0, b).$$

198 Because these parameters cannot reliably be predicted *a priori*, they will be sampled together
 199 with G as will be detailed below.

200

201 With the equations above, we now have all ingredients to calculate $p(D|G, \xi)$:

202

$$203 \quad p(D|G, \xi) = \prod_{i=1}^n \prod_{j=i}^n \text{Poisson}(D_{i,j}; N_D q_{i,j}) \quad (\text{Eq 1})$$

204

205 with:

$$206 \quad p(q_{i,j}|G, \xi) = g(s_G(i,j); \mathbf{P}_t, \mathbf{s}_0, \mathbf{b}) \stackrel{\text{def}}{=} \begin{cases} \mathbf{P}_t \left(\frac{s_G(i,j)}{s_0} \right)^b & \text{if } s_G(i,j) \leq \mathbf{s}_0 \\ \mathbf{P}_t & \text{otherwise} \end{cases} \quad (\text{Eq 2})$$

207

208 **B. Sampling algorithm**

209 The above formulas allow us to compute the posterior probability of any assumed linear genome
 210 structure G (and the nuisance parameters ξ) given D , a contact data set obtained by mapping 3C-
 211 seq/Hi-C reads to an initial set of contigs. In order to explore the entire probability density
 212 $p(G, \xi)$, we need a method to sample the extremely large (or infinite) space of possible linear
 213 genome structures. For this purpose, we implemented an algorithm inspired by the Markov-Chain
 214 Monte-Carlo (MCMC) Gibbs sampler. Starting from an initialization (G_0, ξ_0) the algorithm
 215 makes a large number of random moves across the space (G, ξ) to be sampled, and uses a
 216 probabilistic rule to accept or reject individual moves $(G_t, \xi_t) \rightarrow (G_{t+1}, \xi_{t+1})$. After a sufficient
 217 number of steps, once the chain has reached equilibrium, a subset of the accepted samples can be
 218 used to approximate the global maximum of the probability density $p(G, \xi)$. An overview of the
 219 algorithm's main modules is provided in .

220 At each iteration, GRAAL updates first nuisance parameters ξ_t and then the genome structure G_t .
 221 Below, we separately describe first the initialization of $G_{t=0}$ and $\xi_{t=0}$ and then the update rules
 222 for G_t and ξ_t .

223

224

225

226

227

228 B.1. Initialization of the genome, G_0

229 Different initializations can be considered for the initial set of contigs G_0 depending on the
230 availability of a preliminary assembly of the organism under study, or of a related genome. The
231 initial set of contigs does not need to be perfect, since GRAAL can split and rearrange incorrectly
232 assembled contigs. However, in our current implementation, the restriction fragments μ_i and the
233 bins f_i , whose definition depends on G_0 , cannot be broken. In our paper, we considered the
234 following different types of initializations:

- 235 • The reference budding yeast genome (16 chromosomes; GCF_000146045.1) was used as
236 validation data, since a high quality assembly of this genome is already available. In order
237 to simulate an incomplete assembly of this genome, we split the genome into $N_c = 1,086$
238 bins (of approximately 11Kb) to initialize GRAAL.
- 239 • For YFK1246, the structural mutant of budding yeast¹⁰, $N_c = 3,171$ bins of 9 RFs (*DpnII*
240 restriction enzyme) of the reference budding yeast genome were used to initialize
241 GRAAL.
- 242 • We also used this initialization to assemble the Malaysian budding yeast isolate
243 (UWOPS03-461.4).
- 244 • The *Trichoderma* genome (*ATCC* 13631)¹¹ of strain QM6a was only partly assembled
245 (including using long-insert paired-end data), yielding 77 scaffolds. Rather than
246 initializing G_0 with these scaffolds, those were split into bins of 81 RF, which led to
247 $N_c = 1193$ bins that were used to initialize GRAAL.
- 248 • For the human chromosome 14, we downloaded the 4,722 contigs obtained from the
249 ALLPATHS-LG *de novo* assembly (average size 20kb)¹². A filter was applied to
250 identify RFs (from the *HindIII* restriction enzyme used in the Hi-C experiment¹³)
251 presenting little or no read coverage. If reads appears sparse along a RF compared to
252 the distribution of read coverage over the entire population of RFs, the RF was discarded.
253 If the entire contig appeared undercovered, it was therefore discarded. A similar filtering
254 step is used by dnaTri2. We then split the remaining 2,917 contigs into bins of 3 RFs.

255 As a general strategy to complete the assembly of an imperfectly assembled genome, we
256 recommend starting from the existing contigs and splitting them into bins as illustrated here for
257 *Trichoderma*. The user of GRAAL has the option to choose whether to split these contigs or not
258 (see section B).

259

260 B.2. Initialization of the nuisance parameters, ξ_0

261 The initial values of the parameters $\xi_0 = (P_t, s_0, b)$ are obtained based on the Hi-C data D and
262 the initial genome structure G_0 as follows:

263 First, the initial value of P_t is set to the contact probability averaged over all pairs of bins
264 belonging to different contigs, i.e.:

$$P_t = \frac{\sum_{i<j} (1 - \delta_{\psi_2(i), \psi_2(j)}) D_{i,j}}{\sum_{i<j} (1 - \delta_{\psi_2(i), \psi_2(j)})}$$

265 where $\delta_{i,j} = 1$ if $i=j$ and $\delta_{i,j} = 0$ otherwise.

266 Next, we construct a histogram of *cis*-contact frequencies with genomic intervals $[d_0, \dots, d_M]$,
267 ranging from $d_0 = 0$ to $d_M = \max (L_k)_{k=1..N_c}$, the length of the longest contig in G_0 . For each
268 genomic bin $[d_l, d_{l+1}]$, the histogram reports the mean contact frequency, among all N_c contigs,
269 between bins sharing a contig and separated by genomic distances $s \in [d_l, d_{l+1}]$:

$$F_l = \frac{\sum_{i<j} \delta_{\psi_2(i), \psi_2(j)} H(s_{G_0}(f_i, f_j) - d_l) H(d_{l+1} - s_{G_0}(f_i, f_j)) D_{i,j}}{\sum_{i<j} \delta_{\psi_2(i), \psi_2(j)} H(s_{G_0}(f_i, f_j) - d_{li}) H(d_{l+1} - s_{G_0}(f_i, f_j))}$$

270 where H is the Heaviside function ($H(x) = 1$ for $x \geq 0$ and $H(x) = 0$ otherwise) . We then
271 estimate the initial values of s_0 and b by least squares fitting of (Eq1) to F_l , i.e.:

$$272 (s_0, b) = \arg \min \sum_{l=1}^M \left(F_l - g\left(\frac{1}{2}(d_l + d_{l+1}); P_t, s_0, b\right) \right)^2$$

273 This minimization is performed using a quasi-Newton method¹².

274

275 B.3. Monte Carlo modifications of the genome $G_t \rightarrow G_{t+1}$:

276

277 B.3.1. Virtual mutations:

278 We will call 'virtual mutations' the random changes applied to the genome structure. GRAAL
279 considers 5 different types of elementary virtual mutations and 4 composite mutations as detailed
280 below.

281

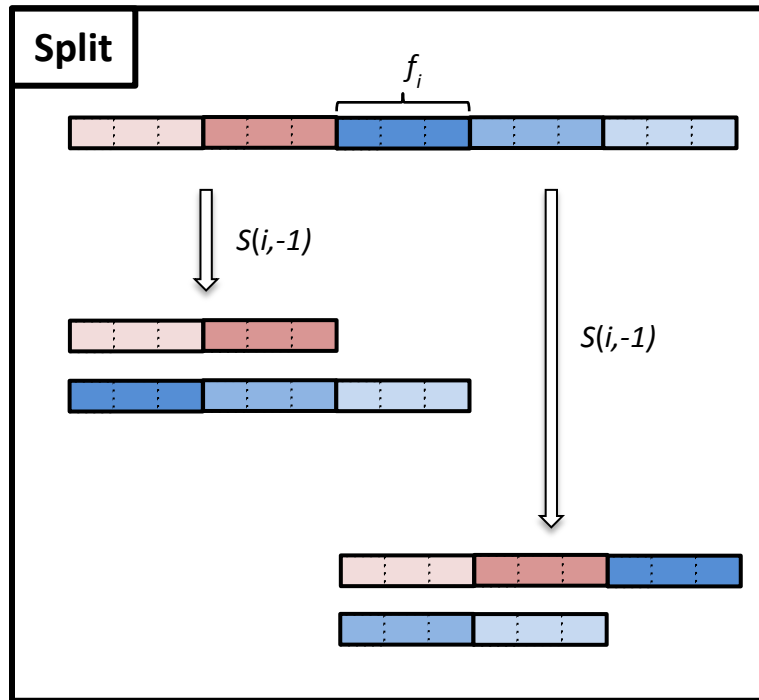
282 Elementary mutations:

283 The 5 elementary mutations are defined as follows:

284 • Split: this mutation splits a contig at a bin and is formally noted as $\mathcal{S}(i, \varepsilon)$, where i is the
 285 index of the bin f_i , and $\varepsilon = \pm 1$ indicates whether the split occurs to the left or right of
 286 the bin. As a result of this operation, contig $\mathcal{C}_{\psi_2(i)}$ is replaced by two new contigs:

$$\begin{cases} \mathcal{C}_{\text{new,left}} = (f_{\varphi_k^1}, \dots, f_{\varphi_k^{l-1}}) & \text{if } \varepsilon = -1 \\ \mathcal{C}_{\text{new,left}} = (f_{\varphi_k^1}, \dots, f_{\varphi_k^l}) & \text{if } \varepsilon = +1 \\ \mathcal{C}_{\text{new,right}} = (f_{\varphi_k^l}, \dots, f_{\varphi_k^{L_k}}) & \text{if } \varepsilon = -1 \\ \mathcal{C}_{\text{new,right}} = (f_{\varphi_k^{l+1}}, \dots, f_{\varphi_k^{L_k}}) & \text{if } \varepsilon = +1 \\ G_{\text{new}} = G \setminus \mathcal{C}_{\psi_2(i)} \cup \mathcal{C}_{\text{new,left}} \cup \mathcal{C}_{\text{new,right}} \end{cases}$$

287 where $(l, k) = (\psi_1(i), \psi_2(i))$.

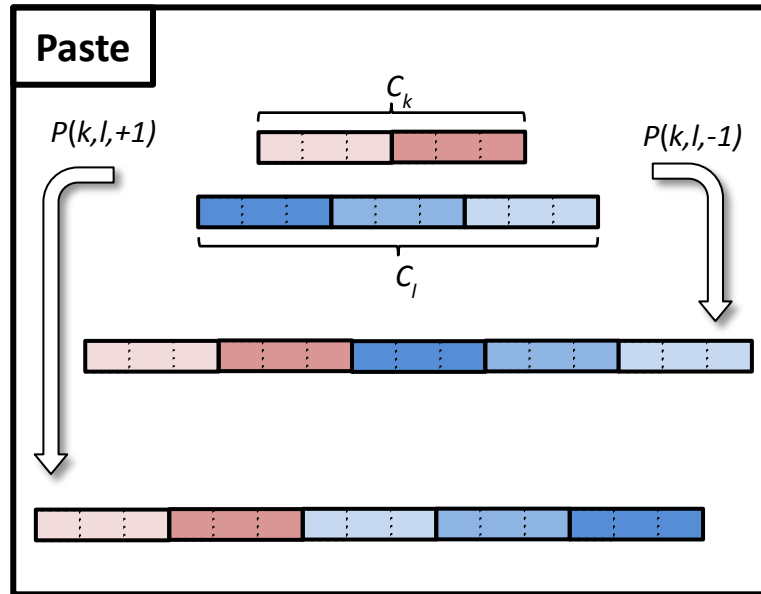


288
 289
 290 • Paste: this mutation concatenates two contigs and is formally noted as $\mathcal{P}(k, l, \varepsilon)$, where k
 291 and l are the indices of the two contigs to be pasted, and $\varepsilon = \pm 1$ indicates whether or not
 292 contig \mathcal{C}_l is flipped before pasting. As a result of this mutation, the two contigs \mathcal{C}_k and \mathcal{C}_l
 293 are replaced by a single new contig obtained by concatenating \mathcal{C}_l (or its flipped version) to
 294 the right of \mathcal{C}_k . Note that the orientation of the bins inside each contig are preserved.

$$\begin{cases} \mathcal{C}_{\text{new}} = \left(f_{\varphi_k^1}, f_{\varphi_k^2}, \dots, f_{\varphi_k^{L_k}}, f_{\varphi_l^1}, f_{\varphi_l^2}, \dots, f_{\varphi_l^{L_l}} \right) & \text{if } \varepsilon = +1 \\ \mathcal{C}_{\text{new}} = \left(f_{\varphi_k^1}, f_{\varphi_k^2}, \dots, f_{\varphi_k^{L_k}}, \mathcal{F} \left(f_{\varphi_l^{L_l}} \right), \mathcal{F} \left(f_{\varphi_l^{L_l-1}} \right), \dots, \mathcal{F} \left(f_{\varphi_l^1} \right) \right) & \text{if } \varepsilon = -1 \end{cases}$$

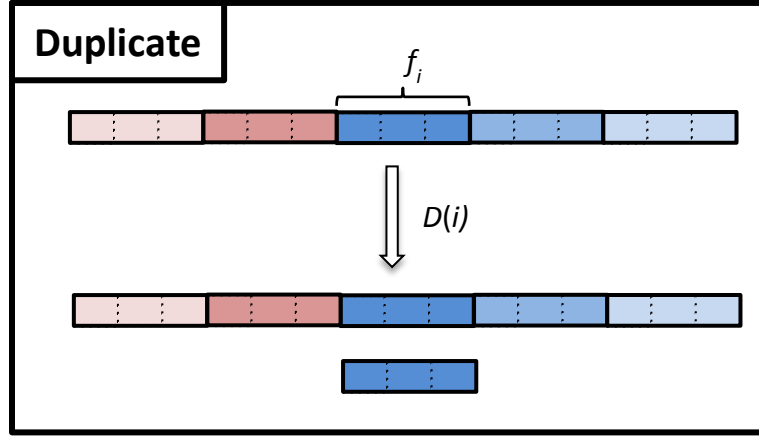
$$G_{\text{new}} = G \setminus \{ \mathcal{C}_k, \mathcal{C}_l \} \cup \mathcal{C}_{\text{new}}$$

295 where the operator \mathcal{F} (flipping a bin) is defined below.
 296 Split and paste are reciprocal operations, i.e. one mutation can reverse the effect of the
 297 other, such that: $\mathcal{P}(k_i^-, k_i^+, -1) \circ \mathcal{S}(i, \varepsilon_1) = \mathcal{S}(i_l, \varepsilon_1) \circ \mathcal{P}(k, l, +1) = \mathcal{N}$, where k_i^- is the
 298 index of the contig resulting from the split operation that was originally to the left of bin
 299 f_i , $k_i^+ = \psi_2(i)$ is the index of the contig still containing f_i after the split, $i_{k,l}$ is the
 300 leftmost bin of contig l , and \mathcal{N} is the "null" mutation, which leaves the genome
 301 unchanged.



302
 303 • Duplicate: This mutation duplicates a bin f_i and is formally noted as $\mathcal{D}(i)$. As a result of
 304 this mutation, a copy of f_i is added to the current set of bins and a new contig consisting of
 305 this single bin is added to the current contig set :

$$\begin{cases} \mathcal{F}_{\text{new}} = \mathcal{F} \cup \{ f_i \} \\ G_{\text{new}} = G \cup \{ f_i \} \end{cases}$$



306

307

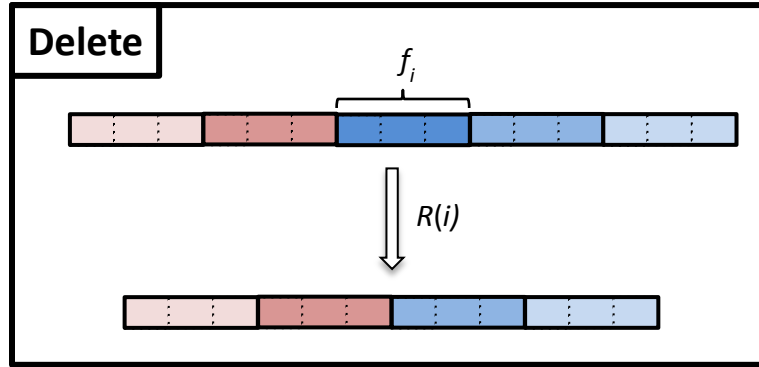
308

- Delete: formally noted as $\mathcal{R}(i)$, this mutation leads to the removal of bin f_i from the current set of bins \mathfrak{F}_t and from the contig that contained it:

$$\begin{cases} \mathfrak{F}_{\text{new}} = \mathfrak{F} \setminus \{f_i\} \\ \mathcal{C}_{k,\text{new}} = (f_{\varphi_k^1}, \dots, f_{\varphi_k^{l-1}}, f_{\varphi_k^{l+1}}, \dots, f_{\varphi_k^{L_k}}) \\ G_{\text{new}} = G \setminus \{\mathcal{C}_k\} \cup \mathcal{C}_{k,\text{new}} \end{cases}$$

309

where $(l, k) = (\psi_1(i), \psi_2(i))$.



310

311

312

313

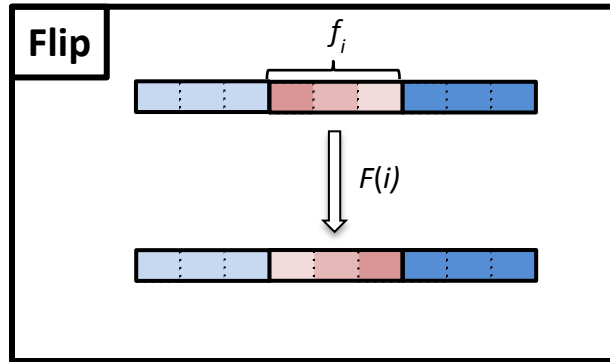
- Flip: formally noted as $\mathcal{F}(i)$, this mutation flips the orientation of bin f_i in its containing contig.

$$\begin{cases} f_{\varphi_k^l, \text{new}} = (\mu_{k,p}, \mu_{k,p-1}, \dots, \mu_{k,1}) \\ \mathfrak{F}_{\text{new}} = \mathfrak{F} \setminus \{f_{\varphi_k^l}\} \cup \{f_{\varphi_k^l}\} \\ \mathcal{C}_{k,\text{new}} = (f_{\varphi_k^1}, \dots, f_{\varphi_k^{l-1}}, f_{\varphi_k^l, \text{new}}, f_{\varphi_k^{l+1}}, \dots, f_{\varphi_k^{L_k}}) \\ G_{\text{new}} = G \setminus \{\mathcal{C}_k\} \cup \mathcal{C}_{k,\text{new}} \end{cases}$$

314

315

where $(l, k) = (\psi_1(i), \psi_2(i))$. The reciprocal operation of a flip is itself: $\mathcal{F}(i) \circ \mathcal{F}(i) = \mathcal{N}$.



316

317

318 Any complex alteration of the genome (defined at the resolution of bins) can be decomposed into
 319 a sequence of these five mutations \mathcal{S} , \mathcal{P} , \mathcal{D} , \mathcal{R} , and \mathcal{F} . However, for complex structural changes
 320 such as translocations, the required sequence may be very long, and it might take unreasonable
 321 time for the sampler to achieve them using Monte Carlo moves. Therefore, we introduce the
 322 following composite mutations:

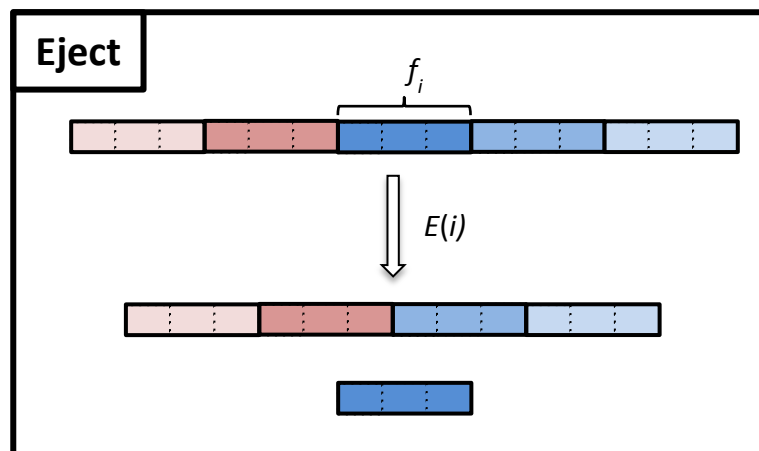
323

324 Composite mutations:

- 325 • Eject: this mutation, noted as $\mathcal{E}(i)$, pops out bin f_i from its contig, and pastes together the
 326 two extremities flanking the bin, leaving f_i as a new contig. It is therefore a composite of
 327 two split and one paste mutations:

$$\mathcal{E}(i) = \mathcal{P}(k_i^-, k_i^+, +1) \circ \mathcal{S}(i, +1) \circ \mathcal{S}(i, -1)$$

328 where k_i^- and k_i^+ are the indices of the contigs resulting from the two splits and originally
 329 located to the left and the right of bin f_i .



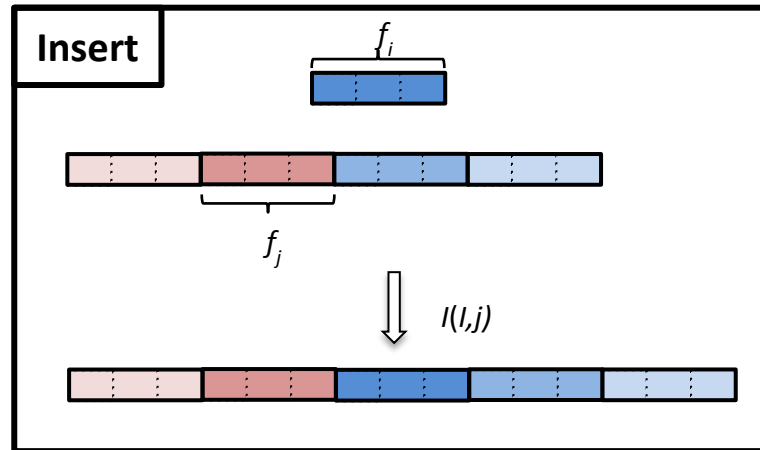
330

331

- 332 • Insert: this mutation, noted as $\mathcal{I}(i, j)$ inserts an isolated bin f_i (i.e. a contig consisting of a
 333 single bin) to the right of bin f_j into its contig $\mathcal{C}_{\psi_2(j)}$. It is a composite of one split and
 334 two paste mutations:

$$\mathcal{I}(i, j) = \mathcal{P}(i, k_j) \circ \mathcal{P}(j, i, +1) \circ \mathcal{S}(j, +1)$$

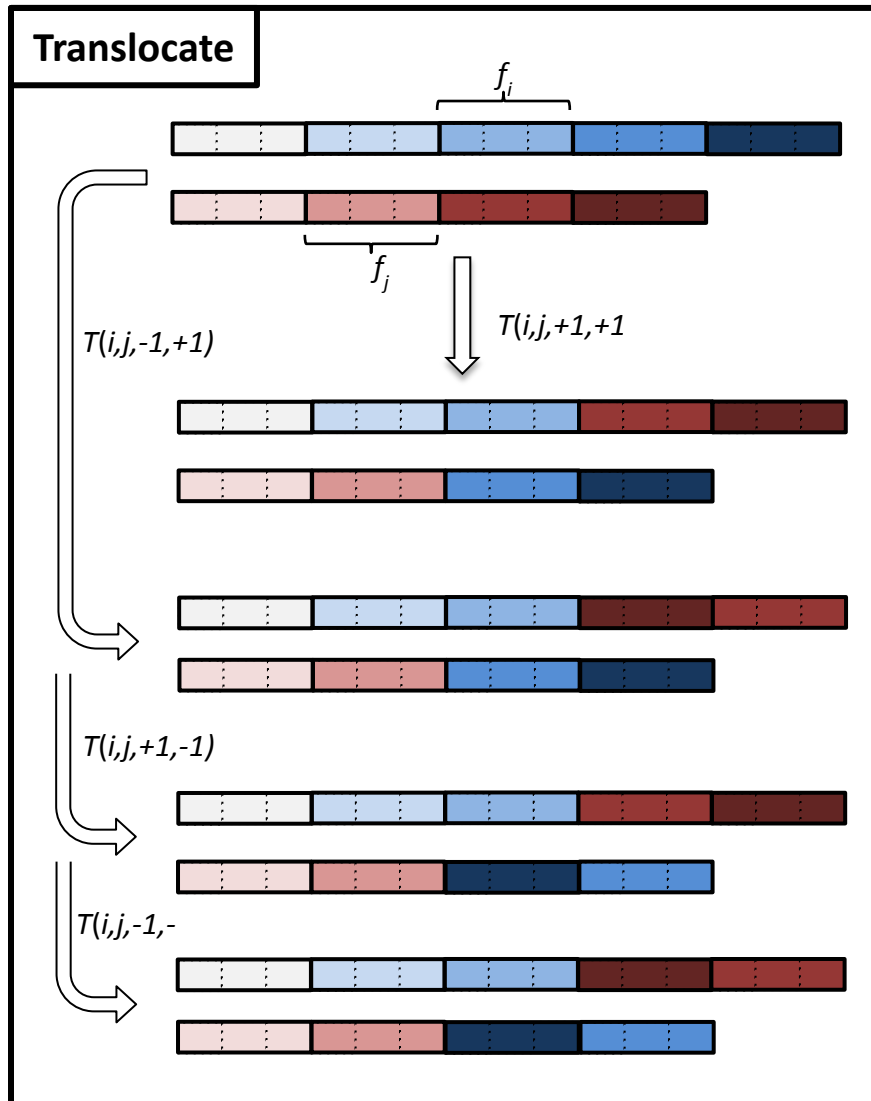
- 335 where $k_j = \psi_1(j) + 1$ designates the bin immediately to the right of f_j within contig
 336 $\mathcal{C}_{\psi_2(j)}$. Ejection and insertions are reciprocal operations, i.e. : $\mathcal{I}(i, j) \circ \mathcal{E}(i) = \mathcal{E}(i) \circ$
 337 $\mathcal{I}(i, j) = \mathcal{N}$.



- 338
 339
 340 • Translocate: this mutation mimics a biological translocation which swaps two parts of
 341 distinct chromosomes and is denoted as $\mathcal{T}(i, j, \varepsilon_1, \varepsilon_2)$, where i and j designate the bin on
 342 the two contigs $\mathcal{C}_{\psi_2(i)}$ and $\mathcal{C}_{\psi_2(j)}$, to the right of which the translocation events take
 343 place, and where $\varepsilon_1 = \pm 1$ and $\varepsilon_2 = \pm 1$ indicate whether the two swapped regions are
 344 flipped or not. This operation is a composite of two split and two paste mutations:

$$\mathcal{T}(i, j, \varepsilon_1, \varepsilon_2) = \mathcal{P}(j, k_j^+, \varepsilon_2) \circ \mathcal{P}(i, k_i^+, \varepsilon_1) \circ \mathcal{S}(j, +1) \circ \mathcal{S}(i, +1)$$

- 345 where $k_i = \psi_1(i) + 1$ and $k_j = \psi_1(j) + 1$ are the indices of the bin immediately to
 346 the right of f_i and f_j , respectively, on contigs $\mathcal{C}_{\psi_2(i)}$ and $\mathcal{C}_{\psi_2(j)}$.



347

348

The reciprocal operation of a translocation is itself:

$$\mathcal{T}(i, j, \varepsilon_1, \varepsilon_2) \circ \mathcal{T}(i, j, \varepsilon_1, \varepsilon_2) = \mathcal{N}$$

349

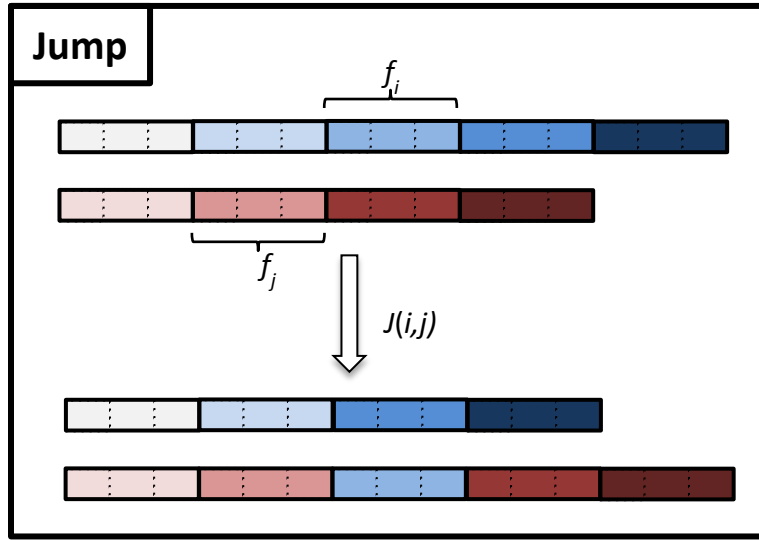
350

- Jump: this mutation, noted $\mathcal{J}(i, j)$, extracts bin f_i from its contig $\mathcal{C}_{\psi_2(i)}$ and inserts it to the right of bin f_j on contig $\mathcal{C}_{\psi_2(j)}$. It can be decomposed into an ejection followed by an insertion:

351

352

$$\mathcal{J}(i, j) = \mathcal{I}(i, j) \circ \mathcal{E}(i)$$



353

354

355 These composite mutations can generate more complex and drastic alterations of genome
 356 structure in a single step, thereby allowing faster exploration of larger regions of structure space
 357 than the elementary mutations.

358

359 We now introduce some notations that will be important for the following section. First, we call
 360 $\mathcal{M} = \{\mathcal{P}, \mathcal{S}, \mathcal{D}, \mathcal{R}, \mathcal{F}, \mathcal{E}, \mathcal{J}, \mathcal{T}, \mathcal{J}\}$ the set of all 9 mutations and use the generic notation $\Theta_i \in \mathcal{M}$
 361 with $i \in \{1, 2, \dots, 9\}$ for individual members of this set (for example, Θ_2 is the paste mutation. We
 362 point out that each of the 9 mutations can be defined by either one (mutations $\mathcal{S}, \mathcal{D}, \mathcal{R}, \mathcal{F}, \mathcal{E}$) or
 363 two indices of bins (mutations $\mathcal{P}, \mathcal{J}, \mathcal{T}, \mathcal{J}$) and, for some mutations, one or two auxiliary binary
 364 parameters $\varepsilon_i = \pm 1$. To formally note the parameters of an arbitrary mutation Θ_k in \mathcal{M} , we can
 365 therefore use the notation: $\Theta_k(i, j, \alpha)$, where it is understood that j is relevant only for mutations
 366 $\mathcal{P}, \mathcal{J}, \mathcal{T}, \mathcal{J}$ and α corresponds to the auxiliary parameter, if relevant (e.g. $\alpha = (\varepsilon_1, \varepsilon_2)$ for $\Theta_8 = \mathcal{T}$,
 367 $\alpha = \{\emptyset\}$ for $\Theta_3 = \mathcal{D}$). We call \mathcal{A}_k the set of all possible values of the auxiliary parameter for
 368 mutation Θ_k . For example, $\mathcal{A}_8 = \{(-1, -1), (-1, +1), (+1, -1), (+1, +1)\}$. Finally, we note
 369 $G^* = \Theta_k(G)$ the structure resulting from application of mutation Θ_k to the genome G .

370

371 B.3.2. Multiple Try Metropolis updates of genome structure

372 Now that we have defined the possible mutations, we explain how they are used to update
 373 genome structures.

374 In devising the sampling algorithm, we initially implemented a basic Metropolis-Hastings
 375 algorithm¹². However, this led to very low acceptance rates of individual moves and excessive
 376 computation time. In order to accelerate the sampling, we therefore implemented a new algorithm
 377 based on a more sophisticated sampling strategy known as Multiple-Try Metropolis that
 378 evaluates several candidate moves at each step and has been shown to allow significantly
 379 improved computation times¹³.

380

381 The canonical MTM method works as follow:

382

- 383 1. Randomly pick one bin f_i by choosing a random integer i between 1 and N (the current
 384 number of bins) with uniform probability.
- 385 2. Next, randomly pick a number K of distinct bins $(f_j)_{j=1..K}$ with $f_j \neq f_i$. In contrast to the
 386 first bin f_i , however, these bins are not drawn with uniform probability, but with a
 387 probability:

$$V_i(j) = \frac{D_{i,j}}{\sum_{k \in [1,N], k \neq i} D_{i,k}}$$

388 As a consequence, the sampled bins f_i tend to have high contact probability with f_i and
 389 are therefore likely to be located in close linear proximity vicinity on the same
 390 chromosome.

- 391 3. Consider the set \mathfrak{G} of all candidate genome structures Γ_l obtained by separately
 392 applying each of the 9 mutations Θ_k to the current genome structure G_t with all possible
 393 values of the auxiliary parameters, i.e.:

$$\mathfrak{G} = (\Gamma_l) = \{\Theta_k(i, j, \alpha)(G_t); k \in [1,9], j = [1, K], \alpha \in \mathcal{A}_k\}$$

394 Among all structures in this set, we pick a random subset of K_1 structures (with uniform
 395 probability):

$$\mathfrak{G}_{K_1} = \{\Gamma_1, \dots, \Gamma_{K_1}\}$$

396 For each of these candidate structure, we evaluate the likelihood $\pi(\Gamma_l) = p(D|\Gamma_l, \xi_t)$
 397 using equations 1 and 2. Note that the nuisance parameters are held constant (they are
 398 updated separately as described in section A.2.4).

- 399 4. For each candidate structure $\Gamma_l = \Theta_k(i, j, \alpha)(G_t)$, we define:

$$w(G_t, \Gamma_l) = \pi(G_t)T(G_t, \Gamma_l)$$

400 where the proposal function T is chosen as:

$$T(G_t, \Gamma_l) = V_i(j)$$

401 5. Among the K_1 proposed candidate structures, we select one, called Γ with probability
402 proportional to

$$w(\Gamma, G_t) = \pi(\Gamma)T(\Gamma, G_t) = \pi(\Gamma)V_j(i)$$

403 6. We note j the index of the bin f_j that led to this structure $\Gamma = \Theta_k(i, j, \alpha)(G_t)$. We then
404 randomly pick another set of K bins f_p , with probability $V_j(p)$ and define a new set of
405 genome structures :

$$\mathfrak{G}^* = (G_l^*) = \{\Theta_k(j, p, \alpha)(\Gamma)\}; k \in [1, 9], p = [1, K], \alpha \in \mathcal{A}_k\}$$

406 Among this set, we randomly pick (with uniform probability) $K_1 - 1$ structures.

$$\mathfrak{G}_{K_1-1}^* = \{G_1^*, \dots, G_{K_1-1}^*\}$$

407 7. Finally, we compute the generalized Metropolis-Hastings acceptance ratio as:

$$r = \min \left\{ 1, \frac{w(\Gamma_1, G_t) + w(\Gamma_2, G_t) + \dots + w(\Gamma_{K_1}, G_t)}{w(G_1^*, \Gamma) + \dots + w(G_{K_1-1}^*, \Gamma) + w(G_{K_1}^*, \Gamma)} \right\}$$

408 With probability r , we accept the new structure Γ and set $G_{t+1} = \Gamma$. In case of rejection,
409 we set: $G_{t+1} = G_t$

410

411 However, in order to lower the computing load of the process we implemented an alternative
412 version of the algorithm. At step 5, we set $G_{t+1} = \Gamma$ and therefore skip steps 6 and 7. The
413 resulting random process is no longer a time homogeneous Markov chain, but the efficiency of
414 this strategy is experimentally verified.

415

416 B.4. Monte Carlo updates of the nuisance parameters $\xi_t \rightarrow \xi_{t+1}$:

417 The nuisance parameters are updated as follows:

418 First, we randomly pick one of the three parameters with equal probability 1/3, i.e. we choose
419 $\theta \in \{P_t, s_0, b\}$. Second, we consider a new candidate value for this parameter by addition of a
420 normally distributed random variable:

$$421 \theta^* = \theta_t + \Delta\theta \quad \text{with} \quad \Delta\theta \sim \mathcal{N}(0, \sigma_\theta)$$

422 We chose to set the variance of the parameter change to a small fraction of the initial value:
423 $\sigma_\theta = 10^{-4} \theta_{t=0}$. This choice was made because of the high sensitivity of the likelihood to small
424 variations of the parameters.

425 We note ξ^* the new candidate set of parameters obtained by replacing parameter θ by θ^* in ξ_t .

426 Next, we accept this candidate with probability:

$$r = \min\left(1, \frac{p(D|G_t, \xi^*)}{p(D|G_t, \xi)}\right)$$

427 where the ratio of likelihoods on the right is computed using Eqs 1 and 2.

428 If this move is accepted, we set: $\xi_{t+1} = \xi^*$, otherwise we keep: $\xi_{t+1} = \xi_t$.

429

430

431 B.5. Sampling from the Markov chain:

432 Starting from the initialization of G_0 and ξ_0 as defined in A.2.1 and A.2.2, we let the Markov
433 moves update G_t and ξ_t under the rules specified in sections A.2.3 and A.2.4 for a total number
434 of iterations N_{\max} . In order to approximate the probability distribution $p(G, \xi|D)$, we discard all
435 samples obtained during an initial burn-in period specified by a number of iterations $N_{\text{burn-in}}$ and
436 use all samples thereafter, i.e. we use (G_t, ξ_t) with $N_{\text{burn-in}} \leq t \leq N_{\max}$. We chose N_{\max} and
437 $N_{\text{burn-in}}$ depending on N_f , the number of restriction fragments in the Hi-C data set D . Typically
438 used values are: $N_{\text{burn-in}} = 3 N_f$ and $N_{\max} = 10 N_f$.

439

440 B.6. Metrics

441 We use different metrics to quantify assembly quality or otherwise characterize the sampled
442 structure probability density.

- 443 • $\text{iqr}(N_{\text{contigs}})$: One simple way to measure the variability among the sampled structures is
444 to measure the variability of contig number. Here, we use the interquartile range (i.e. the
445 difference between the 75% and the 25% percentiles) of the number of contigs in the
446 structure samples $(G_{N_{\text{burn-in}}}, \dots, G_{N_{\max}})$.
- 447 • Error: In order to quantify the quality of assembly on a known genome, we define an error
448 measured as follows: we examine the position of each bin f_i , $i = 1..N_f$ and ask if its
449 immediate flanking neighbors and its orientation are correct. Depending on the answer,

450 we attribute a bin error $E_i \in \{0,1,2,3\}$, where $E_i = 0$ if both neighbors and orientation are
451 correct, and $E_i = 3$ if all are incorrect. We then define the total normalized error as

452 $E = \frac{\sum_{i=1}^{N_f} E_i}{3N_f}$. The normalization ensures that $0 \leq E \leq 1$. A perfect assembly (at the level of

453 bins) yields $E = 0$. Note that this measure is quite sensitive to assembly errors, since any
454 displacement of a bin from its true position (irrespective of the magnitude of this
455 displacement) and any incorrect orientation will increase E .

456

457

458

459 **Supplementary references**

460

- 461 1. Burton, J. N. *et al.* Chromosome-scale scaffolding of *de novo* genome assemblies based on
462 chromatin interactions. *Nat. Biotechnol.* **31**, 1119–1125 (2013).
- 463 2. Kaplan, N. & Dekker, J. High-throughput genome scaffolding from in vivo DNA interaction
464 frequency. *Nat. Biotechnol.* **31**, 1143–1147 (2013).
- 465 3. Rieping, W., Habeck, M. & Nilges, M. Inferential structure determination. *Science* **309**, 303
466 (2005).
- 467 4. Rosa, A. & Zimmer, C. Computational models of large-scale genome architecture. *Int. Rev.*
468 *Cell Mol. Biol.* **307**, 275–349 (2014).
- 469 5. Barbieri, M. *et al.* Complexity of chromatin folding is captured by the strings and binders
470 switch model. *Proc. Natl. Acad. Sci. U. S. A.* **109**, 16173–8 (2012).
- 471 6. Lieberman-Aiden, E. *et al.* Comprehensive mapping of long-range interactions reveals
472 folding principles of the human genome. *Science* **326**, 289–93 (2009).
- 473 7. Duan, Z. *et al.* A three-dimensional model of the yeast genome. *Nature* **465**, 363–367 (2010).
- 474 8. Wong, H. *et al.* A predictive computational model of the dynamic 3D interphase yeast
475 nucleus. *Curr. Biol. CB* **22**, 1881–90 (2012).
- 476 9. Tanizawa, H. *et al.* Mapping of long-range associations throughout the fission yeast genome
477 reveals global genome organization linked to transcriptional regulation. *Nucleic Acids Res.*
478 **38**, 8164–77 (2010).
- 479 10. Koszul, R., Caburet, S., Dujon, B. & Fischer, G. Eucaryotic genome evolution through the
480 spontaneous duplication of large chromosomal segments. *EMBO J.* **23**, 234–243 (2004).
- 481 11. Martinez, D. *et al.* Genome sequencing and analysis of the biomass-degrading fungus
482 *Trichoderma reesei* (syn. *Hypocrea jecorina*). *Nat. Biotechnol.* **26**, 553–560 (2008)
- 483 12. Salzberg, S. L. *et al.* GAGE: A critical evaluation of genome assemblies and assembly
484 algorithms. *Genome Res.* **22**, 557–567 (2011).
- 485 13. Dixon, J. R. *et al.* Topological domains in mammalian genomes identified by analysis of
486 chromatin interactions. *Nature* **485**, 376–380 (2012).
- 487 14. Press, W. H., Teukolsky, S. A., Vetterling, W. T. & Flannery, B. P. *Numerical recipes 3rd*
488 *edition: the art of scientific computing*. 1235 pages (Cambridge University Press, 2007).

- 489 15. Liu, Jun S and Liang, Faming and Wong, W. H. The multiple-try method and local
490 optimization in Metropolis sampling. *J. Am. Stat. Assoc.* **95**, 121–134 (2000).
491