

UC Riverside

UC Riverside Previously Published Works

Title

OMGS: Optical Map-Based Genome Scaffolding

Permalink

<https://escholarship.org/uc/item/2mr9b7tn>

Journal

Journal of Computational Biology, 27(4)

ISSN

1066-5277

Authors

Pan, Weihua
Jiang, Tao
Lonardi, Stefano

Publication Date

2020-04-01

DOI

10.1089/cmb.2019.0310

Peer reviewed

OMGS: Optical Map-Based Genome Scaffolding

WEIHUA PAN, TAO JIANG, and STEFANO LONARDI

ABSTRACT

Due to the current limitations of sequencing technologies, *de novo* genome assembly is typically carried out in two stages, namely contig (sequence) assembly and scaffolding. While scaffolding is computationally easier than sequence assembly, the scaffolding problem can be challenging due to the high repetitive content of eukaryotic genomes, possible mis-joins in assembled contigs, and inaccuracies in the linkage information. Genome scaffolding tools either use paired-end/mate-pair/linked/Hi-C reads or genome-wide maps (optical, physical, or genetic) as linkage information. Optical maps (in particular Bionano Genomics maps) have been extensively used in many recent large-scale genome assembly projects (e.g., goat, apple, barley, maize, quinoa, sea bass, among others). However, the most commonly used scaffolding tools have a serious limitation: they can only deal with one optical map at a time, forcing users to alternate or iterate over multiple maps. In this article, we introduce a novel scaffolding algorithm called OMGS (Optical Map-based Genome Scaffolding) that for the first time can take advantages of multiple optical maps. OMGS solves several optimization problems to generate scaffolds with optimal contiguity and correctness. Extensive experimental results demonstrate that our tool outperforms existing methods when multiple optical maps are available and produces comparable scaffolds using a single optical map.

Keywords: combinatorial optimization, *de novo* genome assembly, optical maps, scaffolding.

1. INTRODUCTION

GENOME ASSEMBLY IS A FUNDAMENTAL PROBLEM in genomics and computational biology. Due to the current limitations of sequencing technologies, the assembly is typically carried out in two stages, namely contig (sequence) assembly and scaffolding. Scaffolds are arrangements of oriented contigs with gaps representing the estimated distance separating them. The scaffolding process can vastly improve the assembly contiguity and can produce chromosome-level assemblies. Despite significant algorithmic progress, the scaffolding problem can be challenging due to the high repetitive content of eukaryotic genomes, possible mis-joins in assembled contigs, and the inaccuracies of the linkage information.

Genome scaffolding tools either use paired-end/mate-pair/linked/Hi-C reads or genome-wide maps. The first group includes scaffolding tools for second generation sequencing data, such as Bambus (Pop et al., 2004; Koren et al., 2011), GRASS (Gritsenko et al., 2012), MIP (Salmela et al., 2011), Opera (Gao et al., 2011), SCARPA (Donmez and Brudno, 2012), SOPRA (Dayarian et al., 2010), and SSPACE (Boetzer et al., 2010) and the scaffolding modules from assemblers ABySS (Simpson et al., 2009), SGA (Simpson

and Durbin, 2012), and SOAPdenovo2 (Luo et al., 2012). Since the relative orientation and approximate distance between paired-end/mate-pair/linked/Hi-C reads are known, the consistent alignment of a sufficient number of reads to two contigs can indicate their relative order, their orientation, and the distance between them. An extensive comparison of scaffolding methods in this first group of tools can be found in Hunt et al. (2014).

The second group uses genome-wide maps such as genetic maps (Tang et al., 2015), physical maps, or optical maps. According to the markers provided by these maps, contigs can be anchored to specific positions so that their order and orientations can be determined. The distance between contigs can also be estimated with varying degree of accuracy depending on the density of the map.

The optical mapping technologies currently on the market (e.g., BioNano Genomics Irys systems and OpGen Argus) allow computational biologists to produce genome-wide maps by fingerprinting long DNA molecules (up to 1 Mb), using nicking restriction enzymes (Samad et al., 1995). Linear DNA fragments are stretched on a glass surface or in a nanochannel array, and then the locations of restriction sites are identified with the help of dyes or fluorescent labels. The results are imaged and aligned to each other to map the locations of the restriction sites relative to each other. While the assembly process for optical molecules is highly reliable, there is clear evidence that a small fraction of the optical molecules is chimeric (Jiao et al., 2017).

A few scaffolding algorithms that use optical maps are available. SOMA appears to be the first published tool that can take advantage of optical maps, but it can only deal with a nonfragmented optical map (Nagarajan et al., 2008). The scaffolding tool proposed in Saha and Rajasekaran (2014) was used for two bacterial genomes *Yersinia pestis* and *Yersinia enterocolitica*, but the software is no longer publicly available. In the last few years, Bionano optical maps have become very popular and have been used to improve the assembly contiguity in many large-scale de novo genome assembly projects (e.g., goat, apple, barley, maize, quinoa, and sea bass) (Pendleton et al., 2015; Bickhart et al., 2017; Daccord et al., 2017; Mascher et al., 2017). To the best of our knowledge, the main tools used to generate scaffolds using Bionano optical maps are SEWINGMACHINE from KSU (Shelton et al., 2015) and HYBRIDSCAFFOLD from Bionano Genomics (unpublished, 2016). SEWINGMACHINE seems to be favored by practitioners over HYBRIDSCAFFOLD.

Both HYBRIDSCAFFOLD and SEWINGMACHINE have, however, a serious limitation: they can only deal with one optical map at a time, forcing users to alternate or iterate over optical maps when multiple maps are available. In this article, we introduce a novel scaffolding algorithm called OMGS (Optical Map-based Genome Scaffolding) that for the first time can take advantage of any number of optical maps. OMGS solves several optimization problems to generate scaffolds with optimal contiguity and correctness.

2. PROBLEM DEFINITION

The input to the problem is the genome assembly to be scaffolded (represented by a set of assembled contigs) and one or more optical maps (represented by a set of sets of genomic distances). We use $C = \{c_i | i = 1, \dots, l\}$ to denote the set of contigs in the genome assembly, where each c_i is a string over the alphabet $\{A, C, G, T\}$. Henceforth, we assume that the contigs in C are chimera free.

An optical map is composed by a set of optical molecules, each of which is represented by an ordered set of positions for the restriction enzyme sites. As said, optical molecules are obtained by an assembly process similar to sequence assembly, but we will reserve the term *contig* for sequenced contigs. We use $M = \{m_i | i = 1, \dots, n\}$ to denote the optical map, where each optical molecule m_i is an ordered set of integers, corresponding to the distances in base pairs between two adjacent restriction enzyme sites on molecule m_i . By digesting in silico the contigs in C using the same restriction enzyme used to produce the optical map and matching the sequence of adjacent distances between sites, one can align the contigs in C to the optical map M . If one is given multiple optical maps obtained using different restriction enzymes, M will be the union of the molecules from all optical maps. In this case, each genomic location is expected to be covered by multiple molecules in M . As said, high quality alignments allow one to anchor and orient contigs to specific coordinates on the optical map. When multiple contigs align to the same optical map molecule, one can order them and estimate the distance between them. By filling these gaps with a number of N 's equal to the estimated distance, longer DNA sequences called *scaffolds* can be obtained.

A series of practical factors make the problem of scaffolding nontrivial. These factors include imprecisions in optical maps (e.g., mis-joins introduced during the assembly of the optical map) (Jiao et al., 2017), unreliable alignments between contigs and optical molecules, and multiple inconsistent anchoring positions for the same contigs. As a consequence, it is appropriate to frame this scaffolding problem as an optimization problem.

We are now ready to define the problem. We are given an assembly represented by a set of contigs C , a set of optical map molecules M , and a set of alignments $A = \{a_{1,1}, a_{1,2}, \dots, a_{l,n}\}$ of C to M , where $a_{i,j}$ is the alignment of contig c_i to optical map molecule o_j . The problem is to obtain a set of scaffolds $S = \{s_1, s_2, \dots, s_k\}$ where each s_i is a string over the alphabet $\{A, C, G, T, N\}$, such that (i) each contig c_i is contained/assigned to exactly one scaffold, (ii) the *contiguity* of S is maximized, and (iii) the conflicts of S with respect to A are minimized. This optimization problem is not rigorously defined unless one defines precisely the concepts of *contiguity* and *conflict*, but this description captures the spirit of what we want to accomplish. In genome assembly, the assembly contiguity is usually captured by statistical measures like the N50/L50 or the NG50/LG50. The notion of conflict is not easily quantified, and even if it was made precise, this multi-objective optimization problem would be hard to solve. We decompose this problem into two separate steps, namely (1) scaffold detection and (2) gap estimation, as explained below.

3. METHODS

As said, our proposed method is composed of two phases: scaffold detection and gap estimation. In the first phase, contigs are grouped into scaffolds, and the order of contigs in each scaffold is determined. In the second phase, distances between neighboring contigs assigned to scaffolds are estimated. The pipeline of the proposed algorithm is illustrated in Figure 1.

3.1. Phase 1: detecting scaffolds

Phase 1 has three major steps. In Step 1, we align in silico-digested chimeric-free contigs to the optical maps (e.g., for a Bionano optical map, we use REFALIGNER), but not all alignments are used in Step 2. We only consider alignments that (i) exceed a minimum confidence level (e.g., confidence 15 in the case of REFALIGNER); (ii) do not overlap each other more than a given genomic distance (e.g., 20 kbp); and (iii) do not create conflict with each other. The method we use here to select conflict-free alignments was introduced in our previous work (Pan et al., 2018). In Step 2, we compute candidate scaffolds by building the *order graph* and formulating an optimization problem on it. In Step 3, either the exhaustive algorithm or a log n -approximation algorithm is used to solve the optimization problem (depending on the size of the graph) and produce the final scaffolds.

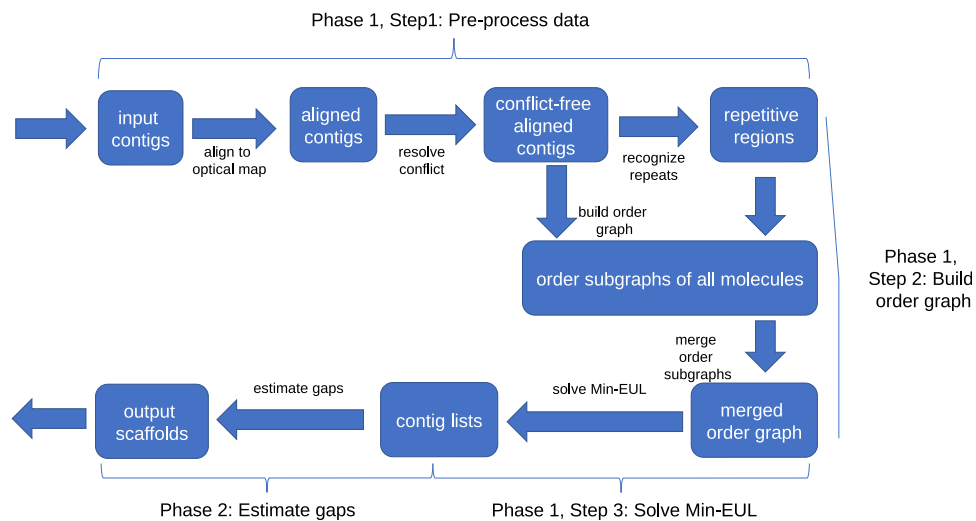


FIG. 1. Pipeline of the proposed algorithm.

3.1.1. Building the order graph. The order graph O is a directed weighted graph, in which each vertex represents a contig. Given two contigs c_i and c_j aligned to an optical molecule o with alignments a_i and a_j , we create a directed edge (c_i, c_j) in O if (i) the starting coordinate of alignment a_i (that we call a_i .start henceforth) is smaller than the starting coordinate of alignment a_j (that we call a_j .start henceforth), (ii) there is no other alignment a_k such that a_k .start is between a_i .start and a_j .start, and (iii) there are no conflict sites between a_i .end and a_j .start on the optical molecule, as defined below. For each alignment a between optical molecule o and contig c , we compute the left overhang l_o and right overhang r_o from o and the left overhang l_c and right overhang r_c from c . The left end of alignment a is declared a *conflict site* if (i) both l_o and l_c are longer than some minimum length (e.g., 50 kbp) and (ii) at least one restriction enzyme site appears in both l_o and l_c . A symmetric argument applies to the right end of the alignment, which determines the values for r_o and r_c .

Directed edge (c_i, c_j) is assigned a weight equal to $\text{qual}(o, a_i.\text{end}, a_j.\text{start}) * (\text{conf}(a_i) + \text{conf}(a_j))$, where (i) $\text{qual}(o, a_i.\text{end}, a_j.\text{start})$ is the *quality* of the region between a_i .end and a_j .start on molecule o (higher is better, defined next), and (ii) $\text{conf}(a)$ is the confidence score provided by REFALIGNER alignment a (higher is better). The quantity $\text{qual}(o, s, t)$ is defined based on the length of a repetitive region between coordinates (s, t) . Based on our experience, assembly mis-joins on optical molecule almost always happen in repetitive regions (Jiao et al., 2017). Given the length of repetitive region $\text{len_rep}(o, s, t)$ in base pairs (defined below), we define the quality of o in the interval (s, t) as $\text{qual}(o, s, t) = e^{-\text{len_rep}(o, s, t)/100,000}$. When a_i and a_j have a small overlap (e.g., shorter than 20 kbp), we set $\text{len_rep}(o, s, t) = 0$.

We recognize repetitive regions in optical molecules based on the distribution of restriction enzyme sites. For a molecule o with n sites, let m_i be the coordinate of the i th site for $i = 1, \dots, n$. As said, molecule o can be represented as a list of positions $\{m_i | i = 1, \dots, n\}$. To determine the repetitive regions in o , we slide a window that covers k sites (e.g., $k = 10$ sites). At each position $j = 1, \dots, n - k + 1$, we select window $w_j = \{m_j, \dots, m_{j+k-1}\}$. While repetitive regions in genome can be highly complex (Zheng and Lonardi, 2005), we observed only two types of repetitive regions in optical molecules, namely single-site repetitive region (Fig. 2A) and two-site repetitive region (Fig. 2B). It is entirely possible that more complex repetitive regions exist: if they do, they seem rare. Based on this observation, to decide whether window w_j is repetitive, we first compute two lists of pairwise distances between sites, namely $D_{j,1} = \{m_{j+l} - m_{j+l-1} | l = 1, \dots, k-1\}$ and $D_{j,2} = \{m_{j+l+1} - m_{j+l-1} | l = 1, \dots, k-2\}$ that we call *distance lists*, then we apply the statistical test described next.

In our statistical test we assume that the values in the distance lists that belong to repetitive regions are independent and identically distributed as a Gaussian. We further assume that each specific distance list ($D_{j,1}$ or $D_{j,2}$) is associated with a Gaussian with a specific mean $\mu_{j,q}$ ($q \in \{1, 2\}$). Finally, we assume that the variance σ^2 is globally shared by all molecules. An estimator of the mean $\mu_{j,q}$ is $\hat{\mu}_{j,q} = \sum_{i=1}^{k-q} d_i / (k-q)$, where $d_i \in D_{j,q}$ and k are the window sizes. To estimate σ^2 , we first get an initial (rough) estimate of the repetitive regions on all molecules. Given a particular $D_{j,q}$, let d_{\max} and d_{\min} be the maximum and minimum distance in $D_{j,q}$. We declare a distance list $D_{j,q}$ to be *estimated repetitive* if $d_{\max} - d_{\min}$ is smaller than a given distance (e.g., 1.5 kbp). We collect all estimated repetitive lists in set $R = \{D_p \text{ is estimated repetitive} | p = 1, \dots, P\}$ and the estimated mean $\hat{\mu}_p$ for each distance list D_p in the set R , where P is the total number of estimated repetitive lists. Then, we define the log likelihood function L as follows (additional details can be found in Section 1.2 of Appendix 1)

$$\log L(\sigma^2) = -\frac{\log \sigma^2}{2} \sum_{p=1}^P |D_p| - \frac{1}{2\sigma^2} \sum_{p=1}^P \sum_{d_i \in D_p} (d_i - \hat{\mu}_p)^2.$$



FIG. 2. Examples of single-site repetitive region (A) and two-site repetitive region (B) in optical maps. Observe the small variations in the repetitive patterns in (B).

By maximizing $\log L(\sigma^2)$, the estimator for the variance becomes

$$\hat{\sigma}^2 = \sum_{p=1}^P \sum_{d_i \in D_p} (d_i - \hat{\mu}_p)^2 / \sum_{p=1}^P |D_p|.$$

Then, we carry out the test on the statistic $d_{\max} - d_{\min}$ for each $D_{j,q}$. The joint density function of (d_{\max}, d_{\min}) is

$$f_{d_{\max}, d_{\min}}(u, v) = n(n-1)f_{d_i}(u)f_{d_i}(v)[F_{d_i}(v) - F_{d_i}(u)]^{n-2}$$

for $-\infty < u < v < +\infty$, where F_{d_i} and f_{d_i} are the distribution function and density function of $d_i \sim N(\hat{\mu}_{j,q}, \hat{\sigma}^2)$, respectively. The density function of $d_{\max} - d_{\min}$ is

$$f_{d_{\max} - d_{\min}}(x) = \int_{-\infty}^{+\infty} n(n-1)f_{d_i}(y)f_{d_i}(x+y)[F_{d_i}(x+y) - F_{d_i}(y)]^{n-2} dy,$$

defined when $x \geq 0$ (additional details can be found in Section 1.3 of Appendix 1). Let now X be a random variable associated with the distribution $f_{d_{\max} - d_{\min}}$. If the p -value $p(X > d_{\max} - d_{\min})$ is greater than a pre-defined threshold (e.g., 0.001), we accept the null hypothesis and declare that window w_j is repetitive. The repetitive regions for the entire molecule o are the union of all the windows w_j 's recognized as repetitive according to the test above.

Once the order graph of each optical molecule is built, we connect all the order graphs, which share the same contigs using the association graph introduced in Pan et al. (2018). The association graph is an undirected graph in which each vertex represents an optical molecule, and an edge indicates that the two molecules share at least one contig aligned to both of them. We use depth first search (DFS) to first build a spanning forest of the association graph. Then, we traverse each spanning tree and connect the corresponding order subgraph to the final order graph. Every time we add a new graph, new vertices and new edges might be added. If an edge already exist, the weights of the new edges are added to the weights of existing edges.

3.1.2. Generating scaffolds. Once the order graph O is finalized, we generate the ordered sequence of contigs in each scaffold. In the ideal case, each connected component O_i of O is a directed acyclic graph (DAG) because the genome is one-dimensional and the order of any pair of contigs is unique. In practice however, O_i may contain cycles caused by the inaccuracy of the alignments and mis-joins in optical molecules. To convert each cyclic component O_i into a DAG, we solve the Minimum Feedback Arc Set problem on O_i . In this problem, the objective is to find the minimum subset of edges (called *feedback arc set*) containing at least one edge of every cycle in the input graph. Since the minimum feedback edge set problem is APX hard, we use the greedy local heuristics introduced in Baharev et al. (2015) to solve it.

We then break each DAG G_i of connected component O_i into subgraphs as follows. In each subgraph, we require the order of every pair of vertices to be uniquely determined by the directed edges. This allows us to uniquely determine the order of the contigs for each scaffold. The formal definition of this optimization problem is as follows.

Definition 1 (Minimum Edge Unique Linearization problem). Input: A weighted DAG $G=(V, E)$. Output: A subset of edges $E' \subseteq E$ such that (i) in each connected component G'_i of the graph $G'=(V, E-E')$ obtained after removing E' , the order of all vertices can be uniquely determined, and (ii) the total weights of the edges in E' are the minimum among all the subset of edges satisfying (i).

In Theorem 1 below, we show that the Minimum Edge Unique Linearization (MIN-EUL) problem is NP-hard by proving that it is equivalent to the Minimum Edge Clique Partition (MIN-ECP) problem, which is known to be NP-hard (Dessmark et al., 2007). In MIN-ECP, we are given a general undirected graph, and we need to partition its vertices into disjoint clusters such that each cluster forms a clique and the total weight of the edges between clusters is minimized.

Theorem 1 MIN-EUL is equivalent to MIN-ECP.

Proof. First, we show that MIN-EUL polynomially reduces to MIN-ECP. Given an instance $G=(V, E)$ of MIN-EUL, we build an instance $G'=(V', E')$ of MIN-ECP as follows. Let $V'=V$. For each pair of vertices

$u, v \in V'$ where v is reachable from u , we define an undirected edge between u and v in E' . For each directed edge $(u, v) \in E$, set the weight of the corresponding undirected edge $(u, v) \in E'$ as 1. Set the weights of the other edges in E' as 0. Then it is easy to see that a MIN-EUL solution to G' is equivalent to a MIN-ECP solution to G and vice versa.

Now we show that MIN-ECP polynomially reduces to MIN-EUL. Given an instance $G' = (V', E')$ (assuming G' is connected) of MIN-ECP, we build an instance $G = (V, E)$ of MIN-EUL as follows. Let $V = V'$. Pick any total linear order O of all vertices in V' . For each undirected edge $(u, v) \in E'$ where $\text{rank}(u) < \text{rank}(v)$ in O , we define a directed edge from u to v in E and set its weight to be the same as its corresponding undirected edge in E' . For any two vertices $u, v \in V$, where $\text{rank}(u) < \text{rank}(v)$ and $(u, v) \notin E^{\text{prime}}$, add a new vertex $x_{uv} \in V$ with $\text{rank}(x_{uv}) = \text{rank}(v)$ and a directed edge u to x_{uv} of weight 1 in E . Now for each pair of vertices $u, v \in V$ where $\text{rank}(u) < \text{rank}(v)$ and $(u, v) \notin E$, add a directed edge u to v with weight zero in E . Then it is easy to see that a MIN-EUL solution to G corresponds to a MIN-ECP solution to G' and vice versa.

Given the complexity of MIN-EUL, we propose an exponential time exact algorithm and a polynomial time $\log n$ -approximation algorithm for solving it. To describe the exact algorithm, we need to introduce some notations. A *conjunction* vertex in a DAG is a vertex which has more than one incoming edge or outgoing edge. A *candidate* edge is an edge which connects at least one conjunction vertex. In Theorem 2 below, we prove that the optimal solution E' of MIN-EUL must only contain candidate edges. Let E_c be the set of all candidate edges in the DAG G ; for each subset E'_j of E_c , we check whether the graph $G' = (V, E - E'_j)$ satisfies requirement (i) in Definition 1 after removing E'_j from G . Among all the feasible E'_j , we produce the set of edges with minimum total weights. To check whether E'_j is feasible, we use a variant of topological sorting, which requires one to produce a unique topological ordering. To do so, we require that in every iteration of topological sorting, the candidate node to be added to sorted graph is always unique. Details of this algorithm are shown as Algorithm 1 in Section 1.1 of Appendix 1.

Theorem 2 *The optimal solution E' of MIN-EUL only contains candidate edges.*

Proof. For sake of contradiction, we assume that E' contains noncandidate edges (u, v) . Since E' is optimal, $G' = (V, E - E')$ satisfies condition (i) in Definition 1. Since both u and v are conjunction vertices, u has only one incoming edge and v has only one outgoing edge. Therefore, by adding (u, v) to $G' = (V, E - E')$, we still satisfy condition (i) in Definition 1. Since the weight of (u, v) is positive, the total weight of $E - E' + \{(u, v)\}$ is larger than $E - E'$. Therefore $E' - \{(u, v)\}$ is optimal, contradicting the optimality of E' .

As said, MIN-EUL is equivalent to MIN-ECP (Theorem 1). In addition, the authors of Dessmark et al. (2007) showed that for any instance of MIN-ECP one can find an equivalent instance of the MINIMUM DISAGREEMENT CORRELATION CLUSTERING problem. As a consequence, any algorithm for the Minimum Disagreement Correlation Clustering problem could be used to solve MIN-EUL. In our tool OMGS, we implemented a $O(\log n)$ -approximation algorithm based on linear programming, originally proposed in Demaine and Immorlica (2003). Standard linear programming packages (e.g., GLPK or CPLEX) are used to solve the linear program. We use the exact algorithm for DAGs with no more than 20 candidate edges and the approximation algorithm for larger DAGs. ■

3.2. Phase 2: estimating gaps

Let $s = \{c_i | i = 1, \dots, h\}$ be one of the scaffolds generated in Phase 1 where each c_i is a contig. In Phase 2, we estimate the length l_i of the gap between each pair c_i and c_{i+1} of adjacent contigs. We estimate all gap lengths $L = \{l_i | i = 1, \dots, h-1\}$ at the same time using the distances between the contigs provided by the alignments and the corresponding order subgraphs. We assume that each l_i is chi-square distributed with α_i degrees of freedom. The choice of chi-square distribution is due to its additive properties, namely the sum of independent chi-squared variables is also chi-squared distributed. Recall that each order subgraph O_k provides a unique ordering $x_k = \{c_j | j = 1, \dots, r\}$ of the contigs aligned to molecule o_k , while the coordinates of the alignment provide the distances between all pairs of adjacent contigs c_j and c_{j+1} as $y_k = \{d_j | j = 1, \dots, r-1\}$. We use the distances d_j as samples to estimate gap lengths l_i . If edge (c_j, c_{j+1}) in O_k is removed in the order graph O when solving MIN-EUL in Phase 1, d_j will be considered not reliable and removed from y_k .

In the ideal case, d_j should be a sample of a single l_i (i.e., $c_j c_{j+1}$ in x_k corresponds to $c_p c_{p+1}$ in s). In practice, however, $c_j c_{j+1}$ in x_k will correspond to a different pair $c_p c_q$ in s where $q > p+1$ (i.e., $c_{p+1} \dots c_{q-1}$ are missing from the order subgraph because some alignments with low confidence were removed in Step 1 of Phase 1). In this situation, after subtracting the length of missing contigs from d_j , $d_j - \sum_{c=c_{p+1}}^{c_{q-1}} |c|$ is a sample of $\sum_{i=p}^{q-1} l_i$ where $|c|$ represents the length of contig c . Since l_p, \dots, l_{q-1} are independent chi-square random variables, $\sum_{i=p}^{q-1} l_i$ is chi-square distributed with degree of freedom $\sum_{i=p}^{q-1} \alpha_i$, so that the log likelihood of this sample is

$$\log l = (\beta - 1) \log \gamma - \frac{\gamma}{2} - \beta \log 2 - \log \Gamma(\beta).$$

Where $\beta = \sum_{i=p}^{q-1} \frac{\alpha_i}{2}$, $\gamma = d_j - \sum_{c=c_{p+1}}^{c_{q-1}} |c|$ and Γ is the gamma function (additional details can be found in Section 1.4 of Appendix 1). The total log likelihood is the sum of the log likelihoods across all samples. To find the α_i maximizing the total log likelihood, we use the Broyden-Fletcher-Goldfarb-Shanno (BFGS) algorithm (Avriel, 2003). Since the mean of a chi-square distribution equals its degree of freedom, we obtain the estimated gaps $\hat{l}_i = \hat{\alpha}_i$. For the case in which the l_i is pre-estimated as negative in the first step, the second and third steps are ignored and the pre-estimated distances are used as final estimates.

Finally, we add $\lceil \hat{l}_i \rceil$ nucleotides (represented by N s) between each pair of contigs c_i and c_{i+1} . When $\hat{l}_i < 0$, we add exactly 100 N s between c_i and c_{i+1} , which is the convention for a gap of unknown length.

4. EXPERIMENTAL RESULTS

We compared OMGS against KSU SEWINGMACHINE (version 1.0.6, released in 2015) and Bionano HYBRIDSCAFFOLD (version 4741, released in 2016) which, to the best of our knowledge, are the only available scaffolding tools for Bionano Genomics optical maps. All tools were run with default parameters, unless otherwise specified. We collected experimental results on scaffolds of (i) cowpea (*Vigna unguiculata*) and (ii) fruit fly (*Drosophila melanogaster*).

4.1. Experimental results on cowpea

Cowpea is a diploid with a chromosome number $2n=22$ and an estimated genome size of 620 Mb. We sequenced the cowpea genome using single-molecule real-time sequencing (Pacific Biosciences RSII). A total of 87 single molecule real time (SMRT) cells yielded about 6M reads for a total of 56.84 Gbp ($91.7 \times$ genome equivalent). We tested the three scaffolding tools on a high-quality assembly produced by CANU (Berlin et al., 2015; Koren et al., 2017) with parameters `corMhapSensitivity=high` and `corOutCoverage=100`, then polished it with QUIVER. We used CHIMERICOGNIZER to detect and break chimeric contigs, using seven other assemblies generated by CANU, FALCON (Chin et al., 2016), and ABRUIJN (Lin et al., 2016) as explained in Pan and Lonardi (2019).

In addition to standard contiguity statistics ($N50^*$, $L50^\dagger$), total assembled size, and scaffold length distribution, we determined incorrect/chimeric scaffolds by comparing them against the high-density genetic map available from Muñoz-Amatriaín et al. (2017). We BLASTed 121-bp long design sequence for the 51,128 genome-wide single nucleotide polymorphisms (SNPs) described in Muñoz-Amatriaín et al. (2017) against each assembly, then we identified which contigs had SNPs mapped to them, and what linkage group (chromosome) of the genetic map those mapped SNPs belonged to. Chimeric contigs were revealed when their mapped SNPs belonged to more than one linkage group. The last line of Tables 1 and 2 reports the total size of contigs in each assembly for which (i) they have at least one SNP mapped to it and (ii) all SNPs belong to the same linkage group (i.e., likely to be nonchimeric).

As said, the three scaffolding tools were run on a chimera-free assembly of cowpea described above using two available Bionano Genomics optical maps (the first obtained using the *BspQI* nicking enzyme, and the second obtained with the *BssSI* nicking enzyme). Since SEWINGMACHINE can only use a single optical map, we alternated the optical maps in input (*BspQI* map first, then *BssSI* and vice versa). SEWINGMACHINE provides two outputs depending on the minimum allowed alignment confidence, namely “default” and

*Length for which the set of contigs/scaffolds of that length or longer accounts for at least half of the assembly size.

†Minimum number of contigs/scaffolds accounting for at least half of the assembly.

TABLE 1. COMPARING OPTICAL MAP-BASED GENOME SCAFFOLDING, SEWINGMACHINE, AND HYBRIDSCAFFOLD ON A COWPEA ASSEMBLY USING ONE OR TWO OPTICAL MAPS

	One optical map						Two optical maps								
	BspQI			BssSI			BspQI+			BssSI+			BspQI & BssSI		
	SM (default)	SM (relax)	HS	SM (default)	SM (relax)	HS	SM (default)	SM (relax)	HS	SM (default)	SM (relax)	HS	SM (default)	SM (relax)	OMGS
Contig/scaffold N50, bp	5,633,882	13,154,336	13,154,336	12,211,658	14,339,314	10,620,326	10,886,079	N/A	11,536,649	14,892,230	14,892,230	14,892,230	13,527,997	14,892,235	16,364,046
Contig/scaffold L50	28	15	17	14	18	17	N/A	N/A	15	13	13	14	14	13	12
Total assembled, bp	511,101,122	521,209,608	521,210,640	516,455,893	518,987,660	518,945,404	N/A	518,252,638	525,577,823	525,198,231	525,827,900	525,105,345	525,827,900	525,105,345	521,324,385
No. of contigs/scaffolds	948	863	877	847	849	846	N/A	832	822	823	816	814	816	814	802
No. of contigs/scaffolds ≥ 100 kbp	269	185	198	170	177	174	N/A	165	149	150	145	143	145	143	137
No. of contigs/scaffolds ≥ 1 Mbp	94	59	63	56	63	62	N/A	59	46	46	48	46	48	46	44
No. of contigs/scaffolds ≥ 10 Mbp	10	20	21	20	18	18	N/A	17	21	21	22	22	22	22	21
Contigs/scaffolds with consistent LG, bp	425,812,490	404,408,642	404,409,674	381,974,417	410,552,582	425,572,265	425,530,009	N/A	424,143,108	385,449,577	385,069,985	425,678,421	403,637,207	432,639,234	

Numbers in boldface highlight the best N50 and scaffold consistency with the genetic map for one map (BspQI and BssSI) or two maps (“A+B” refers to the use of map A followed by map B, “A&B” refers to the use of both maps at the same time).

HS, HYBRIDSCAFFOLD; OMGS, Optical Map-based Genome Scaffolding; SM, SEWINGMACHINE.

TABLE 2. COMPARING OPTICAL MAP-BASED GENOME SCAFFOLDING, SEWINGMACHINE, AND HYBRIDSCAFFOLD ON A COWPEA ASSEMBLY USING OPTICAL MAPS CORRECTED BY CHIMERICOGNIZER

	<i>one optical map</i>						<i>Two optical maps</i>								
	BspQI			BssSI			BspQ/+			BssSI/+			BspQI& BssSI		
	SM (default)	SM (relax)	HS	SM (default)	SM (relax)	HS	OMGS	SM (default)	SM (relax)	HS	OMGS	SM (default)	SM (relax)	HS	OMGS
Contig/scaffold N50, bp	5,633,882	12,487,373	12,495,655	13,505,314	9,420,899	10,886,079	N/A	11,256,770	14,354,752	14,354,752	13,527,997	14,892,235	14,892,235	16,364,046	
Contig/scaffold L50	28	16	15	14	19	17	N/A	16	14	14	14	13	13	12	12
Total assembled, bp	511,101,122	519,785,777	519,785,777	518,405,022	517,678,278	517,636,022	N/A	517,318,151	523,520,329	523,139,705	521,540,185	525,105,345	520,697,623	520,697,623	805
No. of contigs/ scaffolds	948	863	871	849	854	851	N/A	837	823	824	817	814	814	805	
No. of contigs/ scaffolds ≥ 100 kbp	269	185	192	172	182	179	N/A	169	150	151	146	143	143	139	
No. of contigs/ scaffolds ≥ 1 Mbp	94	60	60	58	66	65	N/A	62	48	48	48	46	46	46	
No. of contigs/ scaffolds ≥ 10 Mbp	10	19	19	19	17	17	N/A	17	21	21	21	22	22	21	
Contigs/scaffolds with consistent LG, bp	425,812,490	413,819,557	402,840,302	421,466,164	424,262,883	424,220,627	N/A	423,117,331	402,344,751	401,964,127	420,269,616	403,637,207	431,921,182		

Numbers in boldface highlight the best N50 and scaffold consistency with the genetic map for one map (BspQI and BssSI) or two maps (“A+B” refers to the use of map A followed by map B, “A&B” refers to the use of both maps at the same time).

TABLE 3. COMPARING OPTICAL MAP-BASED GENOME SCAFFOLDING, SEWINGMACHINE, AND HYBRIDSCAFFOLD ON THREE *DROSOPHILA MELANOGASTER* ASSEMBLIES (PRODUCED BY MINIASM, CANU, AND DBG2OLC) USING THE *Bst*QI OPTICAL MAP

	MiniASM assembly						CANU assembly						DBG2OLC assembly							
	Input	SM (default)	SM (relax)	HS	OMGS	Input	SM (default)	SM (relax)	HS	OMGS	Input	SM (default)	SM (relax)	HS	OMGS	Input	SM (default)	SM (relax)	HS	OMGS
Contig/scaffold N50, bp	3,866,686	4,494,241	4,906,224	3,866,686	4,906,224	3,004,953	3,004,953	3,004,953	3,918,649	5,336,340	10,113,899	11,223,142	11,223,142	12,785,467	12,928,771					
contig/scaffold L50	9	8	8	9	8	11	11	11	10	7	6	5	5	5	4					
Total assembled, bp	131,856,353	132,480,826	133,233,999	132,138,056	132,838,677	140,720,404	140,923,974	140,923,974	140,867,226	140,960,395	134,109,164	134,164,629	134,164,629	134,162,857	134,208,377					
No. of contigs/scaffolds	208	205	203	206	206	295	291	291	286	280	339	337	337	331	327					
No. of contigs/scaffolds ≥ 100 kbp	85	82	80	83	83	111	107	107	102	96	78	76	76	70	66					
No. of contigs/scaffolds ≥ 1 Mbp	26	26	25	26	25	31	31	31	29	27	22	22	22	17	16					
No. of contigs/scaffolds ≥ 10 Mbp	2	2	2	2	2	1	1	1	1	5	6	6	6	5	7					
Nonchimeric contigs/scaffolds, bp	131,317,873	125,305,638	132,695,519	131,174,201	132,300,197	140,720,404	140,923,974	140,923,974	140,867,226	140,960,395	134,109,164	134,164,629	134,164,629	134,162,857	134,208,377					

Numbers in boldface highlight the best N50 and the best scaffold consistency with the reference genome.

“relax.” Mode “relax” considers more alignments than “default,” but it has a higher chance of introducing mis-joins. HYBRIDSCAFFOLD failed on the *BssSI* map, so we could not test it on alternating maps.

Table 1 shows that when using a single optical map, OMGS can generate comparable or better scaffolds than SEWINGMACHINE and HYBRIDSCAFFOLD. With two optical maps, OMGS’ correctness (“contigs/scaffolds with 100% consistent LG”) and contiguity (N50) are significantly better than other two tools. Observe that OMGS’ correctness (“contigs/scaffolds with 100% consistent LG”) is even better than the input assembly. This can happen when contigs with SNPs belonging to same linkage group are scaffolded with contigs that have no SNP.

We also compared the performance of OMGS, SEWINGMACHINE, and HYBRIDSCAFFOLD when using optical maps corrected by CHIMERICOGNIZER (on the same cowpea assembly). Observe in Table 2 that OMGS, SEWINGMACHINE, and HYBRIDSCAFFOLD increased the correctness but decreased the contiguity when the corrected *BspQI* optical map was used. The results on the corrected *BssSI* optical map or both corrected optical maps did not change significantly. But again, OMGS produced better scaffolds than SEWINGMACHINE and HYBRIDSCAFFOLD.

4.2. Experimental results on *D. melanogaster*

D. melanogaster has four pairs of chromosomes: three autosomes and one pair of sex chromosomes. The fruit fly’s genome is about 139.5 Mb. We downloaded three *D. melanogaster* assemblies generated in Solares et al. (2018) (https://github.com/danrdanny/Nanopore_ISO1). The first assembly (295 contigs, total size 141 Mb, N50=3 Mb) was generated using CANU (Berlin et al., 2015; Koren et al., 2017) on Oxford Nanopore (ONT) reads longer than 1 kb. The second assembly (208 contigs, total size 132 Mb, N50=3.9 Mb) was generated using MINIMAP and MINIASM (Li, 2016) using only ONT reads. The third assembly (339 contigs, total size 134 Mb, N50=10 Mb) was generated by PLATANUS (Kajitani et al., 2014) and DBG2OLC (Ye et al., 2016) using $67.4\times$ of Illumina paired-end reads and the longest $30\times$ ONT reads. The first and third assemblies were polished using NANOPOLISH (Loman et al., 2015) and PILON (Walker et al., 2014). The Bionano optical for *D. melanogaster* map was provided by the authors of Solares et al. (2018). This *BspQI* optical map (363 molecules, total size=246 Mb, N50=841 kb) was created using IRYSOLVE 2.1 from 78,397 raw Bionano molecules (19.9 Gb of data with a mean read length of 253 kb).

As said, all tools were run with default parameters, with the exception of OMGS’ minimum confidence, which was set at 20 (default is 15). To evaluate the performance of OMGS, HYBRIDSCAFFOLD, and SEWINGMACHINE, we compared their output scaffolds to the high-quality reference genome of *D. melanogaster* (release 6.21, downloaded from FlyBase). We reported the total length of correct/non-chimeric scaffolds as a measure of the overall correctness. To determine which scaffolds were incorrect/chimeric we first selected BLAST alignments of the scaffolds against the reference genome which had an e-value lower than $1e-50$ and an alignment length higher than 30 kbp. We defined a scaffold *S* to be *chimeric* if *S* had at least two high-quality alignments, which satisfied one or more of the following conditions: (i) *S* aligned to different chromosomes; (ii) the orientation of *S*’s alignments was different; or (iii) the difference between the distance of alignments on the scaffold and the distance of alignments on the reference sequence was larger than 100 kbp.

Table 3 reports the main statistics for the three *D. melanogaster* scaffolded assemblies. Even with one map, OMGS’ scaffolds are better than SEWINGMACHINE and HYBRIDSCAFFOLD.

5. CONCLUSIONS

We presented a scaffolding tool called OMGS for improving the contiguity of *de novo* genome assembly using one or multiple optical maps. OMGS solves several optimization problems to generate scaffolds with optimal contiguity and correctness. Experimental results on *V. unguiculata* and *D. melanogaster* clearly demonstrate that OMGS outperforms SEWINGMACHINE and HYBRIDSCAFFOLD both in contiguity and correctness using multiple optical maps.

AUTHOR DISCLOSURE STATEMENT

The authors declare they have no competing financial interests.

FUNDING INFORMATION

This work was supported, in part, by National Science Foundation grants IIS-1814359, IOS-1543963, IIS-1526742, and IIS-1646333, the Natural Science Foundation of China grant 61772197, and the National Key Research and Development Program of China grant 2018YFC0910404.

REFERENCES

- Avriel, M. 2003. *Nonlinear Programming: Analysis and Methods*. Courier Corporation.
- Baharev, A., Schichl, H., Neumaier, A., et al. 2015. An exact method for the minimum feedback arc set problem. *Univ. Vienna* 10, 35–60.
- Berlin, K., Koren, S., Chin, C.-S., et al. 2015. Assembling large genomes with single-molecule sequencing and locality-sensitive hashing. *Nat. Biotechnol.* 33, 623.
- Bickhart, D.M., Rosen, B.D., Koren, S., et al. 2017. Single-molecule sequencing and chromatin conformation capture enable de novo reference assembly of the domestic goat genome. *Nat. Genet.* 49, 643.
- Boetzer, M., Henkel, C.V., Jansen, H.J., et al. 2010. Scaffolding pre-assembled contigs using sspace. *Bioinformatics* 27, 578–579.
- Chin, C.-S., Peluso, P., Sedlazeck, F.J., et al. 2016. Phased diploid genome assembly with single-molecule real-time sequencing. *Nat. Methods* 13, 1050.
- Daccord, N., Celton, J.-M., Linsmith, G., et al. 2017. High-quality de novo assembly of the apple genome and methylome dynamics of early fruit development. *Nat. Genet.* 49, 1099.
- Dayarian, A., Michael, T.P., and Sengupta, A.M. 2010. Sopra: Scaffolding algorithm for paired reads via statistical optimization. *BMC Bioinformatics* 11, 345.
- Demaine, E.D., and Immorlica, N. 2003. Correlation clustering with partial information, 1–13. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*. Springer.
- Dessmark, A., Jansson, J., Lingas, A., et al. 2007. On the approximability of maximum and minimum edge clique partition problems. *Int. J. Found. Comput. Sci.* 18, 217–226.
- Donmez, N., and Brudno, M. 2012. Scarpa: Scaffolding reads with practical algorithms. *Bioinformatics* 29, 428–434.
- Gao, S., Nagarajan, N., and Sung, W.-K. 2011. Opera: Reconstructing optimal genomic scaffolds with high-throughput paired-end sequences, 437–451. In *International Conference on Research in Computational Molecular Biology*. Springer.
- Gritsenko, A.A., Nijkamp, J.F., Reinders, M.J., et al. 2012. Grass: A generic algorithm for scaffolding next-generation sequencing assemblies. *Bioinformatics* 28, 1429–1437.
- Hunt, M., Newbold, C., Berriman, M., et al. 2014. A comprehensive evaluation of assembly scaffolding tools. *Genome Biol.* 15, R42.
- Jiao, W.-B., Accinelli, G.G., Hartwig, B., et al. 2017. Improving and correcting the contiguity of long-read genome assemblies of three plant species using optical mapping and chromosome conformation capture data. *Genome Res.* 27, 778–786.
- Kajitani, R., Toshimoto, K., Noguchi, H., et al. 2014. Efficient de novo assembly of highly heterozygous genomes from whole-genome shotgun short reads. *Genome Res.* 24, 1384–1395.
- Koren, S., Treangen, T.J., and Pop, M. 2011. Bambus 2: Scaffolding metagenomes. *Bioinformatics* 27, 2964–2971.
- Koren, S., Walenz, B.P., Berlin, K., et al. 2017. Canu: Scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Res.* 27, 722–736.
- Li, H. 2016. Minimap and miniasm: Fast mapping and de novo assembly for noisy long sequences. *Bioinformatics* 32, 2103–2110.
- Lin, Y., Yuan, J., Kolmogorov, M., et al. 2016. Assembly of long error-prone reads using de bruijn graphs. *Proc. Natl. Acad. Sci. U.S.A.* 113, E8396–E8405.
- Loman, N.J., Quick, J., and Simpson, J.T. 2015. A complete bacterial genome assembled de novo using only nanopore sequencing data. *Nat. Methods* 12, 733.
- Luo, R., Liu, B., Xie, Y., et al. 2012. Soapdenovo2: An empirically improved memory-efficient short-read de novo assembler. *Gigascience* 1, 18.
- Mascher, M., Gundlach, H., Himmelbach, A., et al. 2017. A chromosome conformation capture ordered sequence of the barley genome. *Nature* 544, 427.
- Muñoz-Amatriaín, M., Mirebrahim, H., Xu, P., et al. 2017. Genome resources for climate-resilient cowpea, an essential crop for food security. *Plant J.* 89, 1042–1054.
- Nagarajan, N., Read, T.D., and Pop, M. 2008. Scaffolding and validation of bacterial genome assemblies using optical restriction maps. *Bioinformatics* 24, 1229–1235.
- Pan, W., and Lonardi, S. 2019. Accurate detection of chimeric contigs via bionano optical maps. *Bioinformatics* 35, 1760–1762.

- Pan, W., Wanamaker, S.I., Ah-Fong, A.M., et al. 2018. Novo&stitch: Accurate reconciliation of genome assemblies via optical maps. *Bioinformatics* 34, i43–i51.
- Pendleton, M., Sebra, R., Pang, A.W.C., et al. 2015. Assembly and diploid architecture of an individual human genome via single-molecule technologies. *Nat. Methods* 12, 780.
- Pop, M., Kosack, D.S., and Salzberg, S.L. 2004. Hierarchical scaffolding with bambus. *Genome Res.* 14, 149–159.
- Saha, S., and Rajasekaran, S. 2014. Efficient and scalable scaffolding using optical restriction maps. *BMC Genomics* 15, S5.
- Salmela, L., Mäkinen, V., Välimäki, N., et al. 2011. Fast scaffolding with small independent mixed integer programs. *Bioinformatics* 27, 3259–3265.
- Samad, A., Huff, E., Cai, W., et al. 1995. Optical mapping: A novel, single-molecule approach to genomic analysis. *Genome Res.* 5, 1–4.
- Shelton, J.M., Coleman, M.C., Herndon, N., et al. 2015. Tools and pipelines for bionano data: Molecule assembly pipeline and fasta super scaffolding tool. *BMC Genomics* 16, 734.
- Simpson, J.T., and Durbin, R. 2012. Efficient de novo assembly of large genomes using compressed data structures. *Genome Res.* 22, 549–556.
- Simpson, J.T., Wong, K., Jackman, S.D., et al. 2009. Abyss: A parallel assembler for short read sequence data. *Genome Res.* 19, 1117–1123.
- Solares, E.A., Chakraborty, M., Miller, D.E., et al. 2018. Rapid low-cost assembly of the *Drosophila melanogaster* reference genome using low-coverage, long-read sequencing. *G3 (Bethesda)* 8, 3143–3154.
- Tang, H., Zhang, X., Miao, C., et al. 2015. Allmaps: Robust scaffold ordering based on multiple maps. *Genome Biol.* 16, 3.
- Walker, B.J., Abeel, T., Shea, T., et al. 2014. Pilon: An integrated tool for comprehensive microbial variant detection and genome assembly improvement. *PLoS one.* 9, e112963.
- Ye, C., Hill, C.M., Wu, S., et al. 2016. Dbg2olc: Efficient assembly of large genomes using long erroneous reads of the third generation sequencing technologies. *Sci. Rep.* 6, 31900.
- Zheng, J., and Lonardi, S. 2005. Discovery of repetitive patterns in dna with accurate boundaries, 105–112. In *Fifth IEEE Symposium on Bioinformatics and Bioengineering (BIBE'05)*. IEEE, Piscataway, NJ.

Address correspondence to:

Prof. Stefano Lonardi
 Department of Computer Science and Engineering
 University of California
 Riverside, CA 92521

E-mail: stelo@cs.ucr.edu

Appendix 1

1.1. DIRECTED ACYCLIC GRAPH UNIQUE ORDERING

Algorithm 1 Sketch of the algorithm for checking whether a directed acyclic graph provides an unique ordering

```

1: procedure ORDER_UNIQUENESS_CHECK( $G=(V, E)$ )
2:    $S$  = nodes with no incoming edges
3:   while  $S \neq \emptyset$  do
4:     if  $|S| > 1$  then
5:       return False
6:     remove a node  $n$  from  $S$ 
7:     for each node  $m$  with an edge  $e=(n, m)$  do
8:       remove edge  $e$  from the  $E$ 
9:       if  $m$  has no other incoming edges then
10:        insert  $m$  into  $S$ 
11:  return True

```

1.2. STATISTICAL TEST FOR REPETITIVE REGIONS

Here we provide additional details for the estimation of σ^2 during the analysis of repetitive regions. Recall that we collect all estimated repetitive lists in set $R = \{D_p \text{ is estimated repetitive} | p = 1, \dots, P\}$ and the estimated mean $\hat{\mu}_p$ for each distance list D_p in the set R , where P is the total number of estimated repetitive lists. For each D_p , the distances d_i 's are distributed as a Gaussian with mean $\hat{\mu}_p$ and variance σ^2 . According to the density function of Gaussian distribution, the log likelihood of one D_p is

$$-\frac{|D_p|}{2} \log(2\pi) - \frac{|D_p|}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{d_i \in D_p} (d_i - \hat{\mu}_p)^2.$$

The total log likelihood is the sum of the log likelihoods across all D_p 's in R , which is

$$\log L(\sigma^2) = -\frac{\sum_{p=1}^P |D_p|}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{p=1}^P \sum_{d_i \in D_p} (d_i - \hat{\mu}_p)^2,$$

after ignoring all terms not related to σ^2 . To maximize $\log L(\sigma^2)$, we require that the derivative of total log likelihood

$$\frac{\partial \log L(\sigma^2)}{\partial \sigma^2} = 0,$$

that is,

$$-\frac{\sum_{p=1}^P |D_p|}{2\sigma^2} + \frac{1}{2(\sigma^2)^2} \sum_{p=1}^P \sum_{d_i \in D_p} (d_i - \hat{\mu}_p)^2 = 0.$$

After some simplification, the estimator for variance becomes

$$\hat{\sigma}^2 = \frac{\sum_{p=1}^P \sum_{d_i \in D_p} (d_i - \hat{\mu}_p)^2}{\sum_{p=1}^P |D_p|}.$$

1.3. DENSITY FUNCTION OF $d_{\max} - d_{\min}$

Here we provide additional details for calculating the density function of $d_{\max} - d_{\min}$. It is well known that the joint density function of order statistics is

$$f_{X(i), X(j)}(u, v) = \frac{n!}{(i-1)!(j-1-i)!(n-j)!} f_x(u) f_x(v) [F_x(u)]^{i-1} [F_x(v) - F_x(u)]^{j-1-i} [1 - F_x(v)]^{n-j} \quad (1)$$

for $-\infty < u < v < +\infty$, where $X(i)$ and $X(j)$ are the i th and j th order statistics in X_1, \dots, X_n , and F_x and f_x are the distribution function and density function of each X_i , respectively. Using Equation (1), the joint density function of (d_{\max}, d_{\min}) can be expressed as

$$f_{d_{\max}, d_{\min}}(u, v) = n(n-1) f_{d_i}(u) f_{d_i}(v) [F_{d_i}(v) - F_{d_i}(u)]^{n-2}$$

for $-\infty < u < v < +\infty$, where F_{d_i} and f_{d_i} are the distribution function and density function of $d_i \sim N(\hat{\mu}_{i,q}, \hat{\sigma}^2)$, respectively.

Now, let $X = d_{\max} - d_{\min}$ and $Y = d_{\min}$. Then $d_{\max} = X + Y$ and $d_{\min} = Y$, and the corresponding Jacobian determinant is

$$J = \begin{vmatrix} \partial d_{\max} / \partial X & \partial d_{\max} / \partial Y \\ \partial d_{\min} / \partial X & \partial d_{\min} / \partial Y \end{vmatrix} = \begin{vmatrix} 1 & 1 \\ 0 & 1 \end{vmatrix} = 1.$$

Thus, the joint density function of (X, Y) is given by

$$f_{X, Y}(x, y) = f_{d_{\max}, d_{\min}}(x + y, y) |J| = n(n-1) f_{d_i}(y) f_{d_i}(x+y) [F_{d_i}(x+y) - F_{d_i}(y)]^{n-2},$$

where $x \geq 0$ and $-\infty < y < +\infty$. By integrating over Y , the density function of $X = d_{\max} - d_{\min}$ becomes

$$f_{d_{\max} - d_{\min}}(x) = \int_{-\infty}^{+\infty} n(n-1) f_{d_i}(y) f_{d_i}(x+y) [F_{d_i}(x+y) - F_{d_i}(y)]^{n-2} dy, \quad x \geq 0.$$

1.4. GAP ESTIMATION

Here we provide additional details for calculating the log likelihood function when estimating gaps. Recall that l_p, \dots, l_{q-1} are independent chi-square random variables, and $\sum_{i=p}^{q-1} l_i$ is chi-square distributed with degree of freedom $\sum_{i=p}^{q-1} \alpha_i$. Since the density function of a chi-square random variable X with degree of freedom k is

$$f_X(x) = \frac{1}{2^{k/2} \Gamma(k/2)} x^{k/2-1} e^{-x/2}$$

where Γ is the gamma function, the likelihood of $\sum_{i=p}^{q-1} l_i$ with observation

$$\gamma = d_j - \sum_{c=c_{p+1}}^{c_{q-1}} |c|$$

is

$$\frac{1}{2^\beta \Gamma(\beta)} \gamma^{\beta-1} e^{-\gamma/2},$$

where $\beta = \sum_{i=p}^{q-1} \frac{\alpha_i}{2}$. Therefore, the log likelihood function for one sample is

$$\log l = (\beta - 1) \log \gamma - \frac{\gamma}{2} - \beta \log 2 - \log \Gamma(\beta).$$

The total log likelihood is the sum of the log likelihoods across all samples.