



INSTITUTO DE COMPUTAÇÃO  
UNIVERSIDADE ESTADUAL DE CAMPINAS

**A program for building contig scaffolds in  
double-barrelled shotgun genome sequencing**

*J. Setubal and R. Werneck*

Technical Report - IC-01-005 - Relatório Técnico

April - 2001 - Abril

The contents of this report are the sole responsibility of the authors.  
O conteúdo do presente relatório é de única responsabilidade dos autores.

# A program for building contig scaffolds in double-barrelled shotgun genome sequencing\*

João Carlos Setubal<sup>†</sup>

Renato F. Werneck<sup>‡</sup>

## Abstract

We describe a program that builds contig scaffolds from contig assemblies, to be used in a whole-genome sequencing project. Our program builds scaffolds based on forward/reverse pair information (both from small clones, such as plasmids, and from large clones, such as cosmids). The program assumes that a DNA assembly, preferably with large repeats masked, is available. A scaffold is a path in a weighted graph, and the main novelty of our approach is a careful weighting scheme for arcs in this graph, such that heavier paths represent more reliable scaffolds. This weighting scheme takes into account the presence of repeats, possible clone duplication, existence of different clone libraries, and hybrid (small clones mixed with large clones) links between contigs. The program provides two different algorithms for scaffold building: one that uses a simple greedy strategy, and one that produces scaffolds that correspond to paths of maximum weight. If  $|N|$  is the number of contigs and  $|L|$  is the number of F/R pairs, the complexities are  $O(|N| + |L| \log |L|)$  (greedy) and  $O(|N|^2 + |L| \log |L|)$  (maximum-weight path). This program has been successfully used in several bacterial genome projects.

## 1 Introduction

There seems to be a general agreement that the most efficient way to sequence the whole genome of a prokaryote is by doing shotgun sequencing. This technique was first effectively demonstrated in the genome of *Haemophilus influenzae* [4], and has been used many times since. We assume knowledge of this technique. The problem with using only shotgun reads is the possible presence of long repeats in the genome. When one uses a “standard” assembly software such as `phrap` [5] to assemble a genome using only shotgun reads, repeats can cause the following problems:

1. Some contigs may be misassembled.

---

\*This is a revised version of TR IC-00-20.

<sup>†</sup>setubal@ic.unicamp.br. Institute of Computing, University of Campinas, 13083-970 Campinas, SP, Brazil. Research supported in part by FAPESP and CNPq.

<sup>‡</sup>rwerneck@inf.puc-rio.br. Department of Informatics, Pontifícia Universidade Católica do Rio de Janeiro, 22453-900 Rio de Janeiro, RJ, Brazil. Research supported in part by CAPES.

2. Two distinct regions that are very similar to one another may be merged into one (this is a *collapsed repeat*).

Example of misassembled contigs:

- result of incorrect assembly:

– A-**rrr**-B ... C-**rrr**-D

where A-**rrr**-B and C-**rrr**-D are contigs, and **rrr** is a repeat.

- correct assembly: A-**rrr**-D ... C-**rrr**-B.

Example of collapsed repeat:

- result of incorrect assembly:

1. A-**rrr**-B-C-**sss**-D

2. **sss**-E-F-**rrr**

- correct assembly: A-**rrr**-B-C-**sss**-E-F-**rrr**-B-C-**sss**-D (note that the actual repeat is **rrr**-B-C-**sss**).

Other examples of problems in assembling in the presence of repeats are presented by Myers et al. [8].

To explain one way to deal with these problems we need the concept of a *forward-reverse pair* (F/R pair for short). A given clone can be sequenced from either end. One end results in the *forward read* and the other results in the *reverse read*. Given that sequencers can read about 800 bp of sequence, then if the clone is larger than, say, 2 kb, the sequences derived from the end reads will not overlap. This information can be used to “virtually” link contigs when read F from a clone is in one contig and read R is in another contig, and each is “pointing” to the other (relying on F/R pairs in an assembly project is also known as *double-barrelled sequencing*, and has been described in [9]). This information can also be used to span repeats, as follows. In addition to obtaining shotgun inserts cloned in plasmids (whose size is no more than 5 kb) the target DNA is cut in larger pieces. These fragments can be cloned in cosmids (where average insert length is about 40 kb) or in BACs (average length 100 kb). This means that these clones are able to span any repeats that might reasonably be expected in prokaryotes. One then sequences only the ends (but both of them) of these large inserts. One then could use an assembly software that would take into account the fact that certain reads belong to clones that are a certain distance apart (the F/R pairs), and therefore deal correctly with repeats.<sup>1</sup>

There are many assembly programs available. Three popular ones are: Staden [1], **phrap** [5], and CAP3 [6]. Staden and **phrap** do not use information from F/R pairs in determining contigs. CAP3 does, but in a limited fashion. This paper describes a program that can be used in conjunction with any assembly program to build a *scaffold* of a genome

---

<sup>1</sup>For an alternative way of handling repeats, which relies mainly on read discrepancies, see [7].

based on F/R pair information. A scaffold in this paper is an ordered sequence of contigs given by links provided by F/R pairs. The use of F/R pair information is done in a much more careful way than either in CAP3 or in `phrap` (which does produce contig scaffolds as part of its output, when F/R information can be deduced from read names). On the other hand, being an add-on to an assembly program our program relies on some manual intervention to be of any practical use in a real genome project. The program has been successfully used in the *Xanthomonas axonopodis* pv *citri* genome project [2] in the following way:

1. All reads are assembled using `phrap`.
2. Repeats are identified by comparing contig sequences to one another (several tools are available for this; one is `cross_match` [5]). Even collapsed repeats can be identified, because usually differing flanking sequences will cause parts of the repeat to separate, as shown in the example above.
3. After the longer ( $\geq 400$  bp, i.e. larger than the average length of a read) repeats have been identified, a new whole assembly should be done, but screening (masking) the reads for those repeats.
4. Our program is applied on the resulting contigs.

In the remainder of this paper we describe the model upon which our program is based and the scaffold construction algorithms. Section 6 provides documentation for the program.

## 2 Overview of Model and Algorithms

In this section we provide an overview that will be detailed in the following sections.

We model the problem of building a scaffold by the problem of finding paths in a weighted directed graph. In this graph, nodes represent contigs and an arc exists between nodes  $u$  and  $v$  if there is at least one F/R link between the corresponding contigs. Given the assembly output, it is relatively straightforward to build such a graph, but one has to be careful with read and contig orientation. One has also to deal with the possible many F/R links between two nodes, and this is a crucial part of our program.

The novelty of our approach is in determining arc weights. The idea is that the weight of an arc  $(u, v)$  represents the degree of confidence that we have that its  $u$  and  $v$  contigs are indeed linked. Therefore, the program has a preprocessing step in which all F/R links for each pair of nodes are carefully analyzed. The result of this analysis is the weight of the arc, and is based on a simple scoring scheme.

Path finding can be done by two different algorithms, at the user's choice. In one of the algorithms (called MWP), we determine maximum weight paths in the graph  $G$ . This algorithm assumes  $G$  is acyclic. In the presence of repeats or errors the kind of graphs we build would not necessarily be acyclic. We use our weighting scheme to throw away arcs that would make  $G$  cyclic. In the other algorithm (called GP) we greedily build vertex-disjoint paths by selecting heavier edges first.

### 3 Arc List Construction

In this section we describe in detail how the list of arcs between nodes is constructed based upon the DNA assembly and F/R link information.

A basic aspect of this construction is that it relies on naming conventions for the reads used in the assembly. This means that it is crucial for our program that most (but not all) read names do correspond to the actual physical clones. It is well known that this practice is subject to a number of errors, but this problem is becoming less of a concern with the increased use of capillary sequencers. We require three pieces of information from a read name: the clone end which it came from (either forward or reverse), the clone library which it came from, and whether it is an `sclone` or an `lclone` (see below). Most genome projects nowadays include such information in read names.

Another basic aspect of our construction is that it makes a distinction between *small clones* and *large clones*. We define an `sclone` as an F/R pair derived from a small clone ( $\leq 5$  kb); and an `lclone` as an F/R pair derived from a large clone (in the range 30–55 kb). With the size ranges mentioned, `sclones` correspond to plasmids and `lclones` correspond to cosmids, but our program is parameterized so that other kinds of clones can be used.

A third basic aspect of our construction has to do with the information we require from the DNA assembly (step 1 in the process described at the end of Section 1). We will base the description below on `phrap`, but other assembly programs should be able to provide the same kind of information. If reads have been named according to the rules above, `phrap` will list all F/R pairs that it has detected as part of its output. In particular it lists F/R pairs whose components are in different contigs. For each such pair  $(r_f, r_r)$ , it will output (among other things) the following information, which is used by our program:

- The orientation of  $r_f$  and of  $r_r$  with respect to the contigs where they have been placed.
- Contig alignment position of the 5' end of each read.
- The *primary* and *secondary alignment scores* of each read. The primary score of a read is the alignment score that the read has with respect to the contig to which it belongs. The secondary score is defined by the `phrap` documentation [5] thus: “the highest score of a match of the read against some other read in a different contig or elsewhere in the same contig (so reads for which this number is non-zero are those which overlap a repeat, or an incorrectly or incompletely assembled region).”

A final basic aspect of arc list construction has to do with the definition of an arc. Contigs have an orientation with respect to each other (and discovering the correct relative orientation is part of the problem). So we consider each contig as having a *left end* and a *right end*. Given a contig  $c_u$  we represent its left end by  $c_u^l$  and its right end by  $c_u^r$ . The paths we will look for in the graph will always have to enter a node (contig) by one of its ends (either left or right) and then leave by the other (right or left). Therefore, an arc is a link between contig *ends* and not between contigs as a whole. Given an F/R pair and their respective contigs, how do we know what ends to link? This depends on the orientation of

Table 1: How links between contig ends are determined.  $c_u$  and  $c_v$  are contigs, and  $r_p$  and  $r_q$  are an F/R pair, such that  $r_p$  is in  $c_u$  and  $r_q$  is in  $c_v$ . When arrows are in the same direction it means that the direct sequence of the read was aligned; when arrows have opposite orientation, it means that the reverse complement of the read was aligned.

	$\overrightarrow{c_u} \overrightarrow{r_p}$	$\overrightarrow{c_v} \overrightarrow{r_q}$	$\overrightarrow{c_u} \overleftarrow{r_p}$	$\overrightarrow{c_v} \overleftarrow{r_q}$
$\overrightarrow{c_u} \overrightarrow{r_p}$	$c_u^r \leftrightarrow c_v^r$	$c_u^r \leftrightarrow c_v^r$	$c_u^r \leftrightarrow c_v^l$	$c_u^r \leftrightarrow c_v^l$
$\overrightarrow{c_u} \overleftarrow{r_p}$	$c_u^l \leftrightarrow c_v^r$	$c_u^l \leftrightarrow c_v^r$	$c_u^l \leftrightarrow c_v^l$	$c_u^l \leftrightarrow c_v^l$

each read in its contig. Assuming contigs  $c_u$  and  $c_v$  and an F/R pair  $(r_p, r_q)$ , so that read  $r_p$  is in contig  $c_u$  and read  $r_q$  is in contig  $c_v$ , we have arcs as given by Table 1.

Now we finally come to actual arc list construction. This is done in two phases. In the first phase all F/R pairs whose components are in different contigs are scanned. From the contig alignment positions we compute  $\mathcal{L}$ , which is the sum of the distance between the 5' end of one member of the pair to the far end (following the 5'  $\rightarrow$  3' direction of the read) of the contig it has been placed in with the analogous distance for the other member of the pair. F/R pairs derived from `lclones` are retained only if their  $\mathcal{L} \leq \text{LMAX}$  (user-defined). F/R pairs derived from `sclones` are retained only if their  $\mathcal{L} \leq \text{SMAX}$  (user-defined). For F/R pairs that are retained,  $\mathcal{L}$  is considered a lower bound on the size of the clone.

In the second phase all F/R links between the same pair of contig ends are analyzed, for each pair of contig ends, with the aim of assigning a weight to the arc that will represent this link. This is done by a simple scoring scheme in which a read can contribute certain “points” to the total weight. When assigning points to links we look for events that confirm the links and that are independent as much as possible from each other. We now describe this scoring scheme in detail. Please note that the scheme relies on several parameters, whose values should be given by the user (defaults are provided).<sup>2</sup> All parameters are indicated using `courier` font. The details are as follows:

- First, each read receives a “uniqueness” value. Uniqueness is defined as follows:

$$\text{uniqueness} = 1 - \frac{\text{secondary score}}{\text{primary score}}$$

With this definition, well-anchored reads (secondary score = 0) have uniqueness = 1. Low values for uniqueness indicate that the read probably overlaps a repeat. However, even a read with a low uniqueness may be in the correct place. So rather than simply discard reads with low uniqueness, we postpone this judgement until we have all F/R pairs that make up a given arc. Each arc receives a uniqueness value derived from the uniqueness of its component reads in the following way. The uniqueness of each F/R pair is taken as the minimum value of the uniqueness of each of its reads. The

<sup>2</sup>The values for parameters LMIN, LMAX, and SMAX can also be specified by the user, but they differ from the other parameters in that they should be specified in the “library file”; see Section 6 for details.

uniqueness of an arc will be the mean value of the uniquenesses of its F/R pairs. Arcs for which the uniqueness is below a certain threshold (`minimum_uniqueness`) are discarded. Note also that F/R pairs of a given arc that have been discarded because of duplication or inconsistencies are not used in uniqueness computation.

- Putative identical F/R pairs are identified, and only one copy is retained (with the aim of excluding redundant, duplicate reads from the same clone). An F/R pair  $(r_x^p, r_y^q)$  is considered identical to another F/R pair  $(s_x^p, s_y^q)$  (where  $x$  and  $y$  are the contigs), if the alignment positions of  $r_x^p$  and  $s_x^p$  differ by at most `minimum_alignment_difference` bp, and if the alignment positions of  $r_y^q$  and  $s_y^q$  differ by at most the same value. Note that this check does not take into account the name of the reads; i.e., two reads could have totally different names and still be considered copies of each other.
- The first F/R pair from a given library contributes `weight_first_{small,large}` points to the total weight; each additional F/R pair from the same library contributes `weight_other_{small,large}` points. We regard reads from the same library as being events that are not as independent as reads from different libraries, and therefore `weight_first` should be larger than `weight_other`.
- If all links are just of one type (all `sclones` or all `lclones`), then the weight of the arc is determined by the above rules. If there are `lclones` mixed with `sclones`, then further processing is done as follows.

First a test is done to check whether each `lclone` link yields a value for  $\mathcal{L}$  that is not too small. Note that having links from `sclones` already indicates that the contigs could be quite close to each other in the genome. This means that one or both ends of `lclones` should be quite a distance away from the contig end. To be consistent with this an `lclone` link must have  $\mathcal{L}$  larger than or equal to `LMIN` (user-defined). If all `lclone` links are consistent with small separation between contigs, then both sets of links are considered valid, and additional `hybrid_bonus` points are awarded to the arc weight. These bonus points come from the empirical observation that  $k$  consistent `sclone` and `lclone` links give more credence to the contig link than  $k$  links of just one type.

If some `lclones` are long enough but some are not, then the ones for which  $\mathcal{L}$  is too small are discarded (and their points are not counted), but additional `hybrid_bonus` points are still awarded to the arc weight.

If all `lclones` seem to be too short, then the arc weight is determined by either `sclone` links only or by `lclone` links only. The set that has larger weight determines the choice.

It is also possible to specify clones that are neither `sclones` nor `lclones`. They are called `medium` clones. If present in an arc with other kinds of clones and consistent with them, such clones do not make an arc hybrid.

## 4 Algorithms

### 4.1 Description

This section describes the algorithms used to find scaffolds in the genome: *Greedy Path* (GP) and *Maximum Weight Path* (MWP).

They share a common basic structure. First, they read an input file describing the problem by enumerating contigs and F/R links. Then, a set  $A_0$  of weighted arcs is built according to the procedure described in Section 3. Starting from a graph  $G = (N, \emptyset)$ , both algorithms try to add arcs to  $G$  in a greedy fashion, scanning  $A_0$  in nondecreasing order by weight. To be actually inserted, each arc must satisfy a given set of conditions. Finally, both algorithms find “good” paths in  $G = (N, A)$ , the resulting graph. The essential difference between the algorithms is the set of conditions an arc  $a$  must satisfy in order to be inserted into  $G$ .

In the description below we use  $(c_x^p, c_y^q)$  to represent an arc, where  $c_x^p$  is end  $p$  (either left or right) of contig/node  $x$ , and  $c_y^q$  is end  $q$  (either left or right) of contig/node  $y$ .

#### 4.1.1 Greedy Path Algorithm

Let  $a = (c_x^p, c_y^q)$  (with  $1 \leq x, y \leq |N|$ ) be an arc we want to insert into  $G$ . For  $a$  to be actually inserted, two conditions must hold: (1)  $c_x$  and  $c_y$  must belong to different connected components, and (2) both  $c_x^p$  and  $c_y^q$  must be free, i.e., there must be no arc incident to any of them in  $G$ .

Note that this strategy is very restrictive. Not only does it avoid cycles, but it also forbids parallel paths. By the end of the algorithm, the graph will become a collection of vertex-disjoint paths, all of which are output.

#### 4.1.2 Maximum Weight Path Algorithm

An arc  $a = (c_x^p, c_y^q)$  (with  $1 \leq x, y \leq |N|$ ) will be inserted into  $G$  by MWP if (1) at least one of  $\{c_x^p, c_y^q\}$  is free, (2)  $a$ 's orientation does not conflict with previously inserted arcs, and (3)  $a$  does not create a cycle. Note that conditions (2) and (3) will be relevant only when  $c_x$  and  $c_y$  belong to the same connected component. There is no really good reason for requiring condition (1) except that when the condition is not satisfied arc  $a$  has weight less than or equal to the weight of *two* other competing arcs. When strict inequality holds, this arc is less reliable.

Once the graph is built, the algorithm outputs a number of paths of different classes. For each connected component, the procedure is the following. First, find the heaviest *main path*  $p_1$  and output it. Second, report all *alternate paths* between contigs in  $p_1$ , i.e., paths that start and end in contigs that appear in  $p_1$  but do not use any arc used by  $p_1$ . Then, remove all arcs of  $p_1$  from the graph and repeat the process, finding the second heaviest main path  $p_2$  and the corresponding alternate paths, and so on for other main paths. As an additional constraint, no contig in  $p_i$  may have appeared in  $p_j$ ,  $j < i$ . In other words, all main paths are vertex-disjoint. The connected component will be fully processed only



when every arc is either part of some main path  $p_i$  or has both ends in main paths (in this case, the arc is said to *link* the paths and is listed accordingly in the output file).

## 4.2 Data Structures and Running Times

The first step of both GP and MWP is to build the set  $A_0$  of arcs from the set  $L$  of F/R links read from the input file. The links that constitute the arcs can be determined in  $O(|L| \log |L|)$  total time: sort  $L$  according to link ends and build the arcs from consecutive elements in  $L$ . Since the weight of each arc  $a$  is heuristically determined in linear time w.r.t. the number of links in  $a$ ,  $O(|L| \log |L|)$  is indeed the complexity of building  $A_0$ .<sup>3</sup> Sorting this set according to arc weights can be done in  $O(|A_0| \log |A_0|)$  time. An empty graph  $G$  can be built in  $O(|N|)$  time,  $|N|$  being the number of contigs.

Up to this point, the overall complexity of both algorithms is  $O(|N| + |L| \log |L|)$ , since  $|L| \geq |A_0|$ . The following subsections analyze the remaining operations.

### 4.2.1 Greedy Path Algorithm

For each arc  $a$ , two tests must be made for it to be inserted. First, the ends of  $a$  must be in different connected components. Using a forest-based implementation of a union-find data structure, we can check all arcs in  $O(|A_0| \cdot \alpha(|N| + |A_0|, |N|))$  total time [10]. Second, the ends of  $a$  must be free. Testing this is trivial and can be done in  $O(1)$  time. The paths in the graph  $G$  build by GP can be found in  $O(|N|)$  total time, since  $G$  is actually a collection of vertex-disjoint paths. Therefore, the overall complexity of the algorithm is still  $O(|N| + |L| \log |L|)$ .

### 4.2.2 Maximum Weight Path Algorithm

Although MWP's conditions are not as strict as GP's, testing them is a more complex task. Since contigs in the same connected component may be linked by a new arc, a union-find data structure is not enough. Nevertheless, the algorithm does require a union-find data structure. Since we sometimes need to know the complete list of contigs that are part of a connected component, the list implementation [3, section 22.2] is the better choice in this case.

Let  $a = (c_x^p, c_y^q)$  be the arc we want to insert. The easier condition to test is whether at least one of  $\{c_x^p, c_y^q\}$  is free: this can be done in  $O(1)$  time. If the arc passes this test and if its ends belong to different connected components (which can also be determined in  $O(1)$  time using the list implementation of the union-find data structure), the insertion is performed.

On the other hand, if the contigs are already connected, their relative orientation must be consistent with that suggested by  $a$ . In order to make this test in  $O(1)$  time for each arc, we keep, at all times, the relative orientation of all contigs. We start the algorithm assuming that every contig has the same orientation (whether it is LR or RL is irrelevant).

---

<sup>3</sup>In the worst case weight determination could be quadratic in the number of reads, because of the duplicity check. This is kept linear because we define as a constant the maximum number of reads to be checked for duplicity.

This information is updated as needed, namely when we insert an arc joining different connected components. For simplicity, we will use an example to explain how this is done. Assume that arc  $a = \{c_1^r, c_2^r\}$  is inserted between two different connected components and that both  $c_1$  and  $c_2$  have orientation LR. The evidence of arc  $a$  indicates that  $c_1$  and  $c_2$  cannot have the same orientation in the genome. Therefore, one of the contigs must be changed to RL. Actually, if it is already linked to other contigs, the whole connected component is flipped: RL contigs become LR and vice-versa. Although any of the connected components could be flipped, we achieve better performance by choosing the one with fewer contigs. This ensures that at most  $O(|N| \log |N|)$  contig flips will be necessary during the algorithm. (The entire flipping procedure can be interpreted as a subroutine of the *union* operation of the list-based union-find data structure.)

The third and most costly condition to test is whether the new arc creates a cycle. Since it would link  $c_x^p$  to  $c_y^q$ , we have to check whether there is a path in  $G$  joining the ends of  $c_x$  and  $c_y$  not involved in the link (if  $p \neq q$ , it would be a path starting in  $c_y^p$  and ending in  $c_x^q$ ). This requires  $O(|A|)$  time in the worst case.

Once the graph  $G$  is built, MWP finds and reports the paths. Since the resulting graph is directed and acyclic, maximum weight paths can be found in  $O(|A|)$  time (classic computer science result; see for example [3, section 25.4]). There will be at most  $O(|N|)$  main paths, since they are all vertex-disjoint. Therefore, all  $O(|N|)$  main paths can be found in  $O(|N||A|)$  time. The algorithm finds at most one alternate path starting at each contig, so  $O(|N||A|)$  is enough to find all alternate paths. Links can be easily found in  $O(|A|)$  time.

All steps considered,  $O(|L| \log |L| + |A|^2 + |N||A| + |N| \log |N|)$  is the worst case running time of MWP. An arc can be inserted only if one of its contig ends is free. Since there are  $2|N|$  contig ends, we have  $|A| = O(|N|)$ . Therefore, the complexity of MWP can be rewritten as  $O(|N|^2 + |L| \log |L|)$ .

## 5 Final Remarks

As already remarked, the program described has been successfully used in bacterial genome projects [2]. The values for the library file parameters were as follows:

- LMAX = 55 kb
- LMIN = 30 kb
- SMAX = 5 kb

The values for the other parameters were the defaults indicated in their descriptions in Section 6. The scaffold program takes a few seconds on a Compaq DS20 workstation, on a genome project with hundreds of contigs. This is insignificant compared to the time spent in assembly (measured in hours). The program is available from the authors upon request.

We finish the paper outlining some ways in which the program can be improved. The authors are currently working on these improvements.

- The most problematic part of the program is the requirement that long repeats be manually masked in order for the program to find correct scaffolds. It is possible

to make the program detect these long repeats based on F/R pair information, and then proceed to break misassembled contigs in appropriate positions. The resulting scaffold(s) would then be made not only out of contigs but also out of contig pieces.<sup>4</sup>

- The arc weighting scheme proposed is arbitrary. Although the authors made an effort to parameterize most components, it is desirable to have a weighting scheme that is based on probability models related to the biological phenomena that cause misassemblies. We believe that a much sounder weighting scheme can be devised along these lines.
- Even if the program is able to detect repeats and even if it has a more defensible weighting scheme, some assemblies may result in inconsistent paths (for example, with bifurcations leading to paths of inconsistent length). The program should therefore try to assign some confidence level to each possible path or path section. The scaffolds presented should be described also in terms of the quantitative confidence the program has in their various parts. Such an output is much more helpful than a simple list of scaffolds. The revised arc weighting scheme alluded to above would naturally form the basis for this confidence analysis.
- Taking into account the observations above, the program should be built around just one algorithm instead of offering the user a choice of two algorithms. In programs that attempt to solve problems of a biological nature, it can be the case that some of the choices a program offers to users simply represent a lack of confidence on its possible results. It can thus be said that sometimes fewer choices mean better programs.
- The program presented here, although used successfully in real genome projects, has not been tested against other similar programs. When a new version of this program becomes available, we plan to test it against other programs on real data so as to validate its approach.

## References

- [1] J. K. Bonfield, K. F. Smith, and R. Staden. A new DNA sequence assembly program. *Nucleic Acids Research*, 23:4992–4999, 1995.
- [2] ONSA consortium. *Xanthomonas axonopodis* pv *citri* genome project. [www.lbi.ic.unicamp.br](http://www.lbi.ic.unicamp.br).
- [3] T. Cormen, C. Leiserson, and R. Rivest. *Introduction to Algorithms*. MIT Press/McGraw-Hill, 1990.

---

<sup>4</sup>Since our program relies on results of an assembly program, one might argue that it would be better to write a completely new assembly program, which could handle scaffolds and repeats as we do or propose to do. Given enough time and human resources, this is certainly the best approach. But with both scarce, it is quicker and easier to use the results of a program such as *phrap*.

- [4] R. D. Fleishmann et al. Whole-genome random sequencing and assembly of *Haemophilus influenzae* Rd. *Science*, 269:496–512, 1995.
- [5] P. Green. Phrap documentation. [www.phrap.org](http://www.phrap.org).
- [6] X. Huang and A. Madan. CAP3: A DNA sequence assembly program. *Genome Research*, 9:868–877, 1999.
- [7] J. Kececioglu and J. Yu. Separating repeats in DNA sequence assembly. *Proceedings of RECOMB'01*, 2001.
- [8] E. W. Myers et al. A whole-genome assembly of *Drosophila*. *Science*, 287:2196–2204, 2000.
- [9] J. C. Roach, C. Boysen, K. Wang, and L. Hood. Pairwise end sequencing: a unified approach to genomic mapping and sequencing. *Genomics*, 26:345–353, 1995.
- [10] R. E. Tarjan. Efficiency of a good but not linear set union algorithm. *Journal of the Association for Computing Machinery*, 2:212–225, 1975.

## 6 Program Documentation

### 6.1 Command Line Options

The program is called `genscaff`. Its command line has the following format:

```
genscaff <intermediary file> <library file> [options]
```

The intermediary file, detailed in section 6.3, lists the F/R pairs and the contigs. The library file, described in section 6.4, contains information about the libraries. The optional parameters are:

- a `<mwp|gp>`: algorithm to be executed, MWP or GP (default is GP);
- v: verbose mode (report the names of all F/R pairs that belong to each arc appearing in the scaffolds);
- V: silent mode (report the name of just one F/R pair in each arc appearing in the scaffolds — this is the default);
- c `<configuration file>`: read options from a configuration file (described in section 6.2).

The program produces two output files, both with the same prefix as the intermediary file. The *complete output file* (extension: `allscaff`), described in section 6.5, contains not only the scaffolds, but also some additional information that may lead to a better understanding of the genome assembly. This file has a very rigid format, which makes it easy to be parsed and processed automatically. The *simplified output file* (extension: `scaff`), described in section 6.6, contains just the scaffolds in a format that is easier to read.

## 6.2 Configuration File

Each line in the configuration file contains a tag followed by its value. Valid tags are:

- algorithm:** Either `gp` or `mwp` (`GP` is the default algorithm).
- minimum\_uniqueness:** Minimum uniqueness an arc must have in order to be considered; must be a real value no greater than 1.0 (default is 0.800).
- small\_large\_threshold:** Hybrid threshold (in base pairs). Clones whose maximum lengths are smaller than this value will be considered `sclones`; clones whose minimum lengths are larger than the threshold will be `lclones` (default is 10,000).
- maximum\_alternate\_difference:** Maximum percentual difference in length allowed between a consistent alternate path and a main path (default is 20.0). Section 6.5.10 details how this parameter is used.
- minimum\_alignment\_difference:** Minimum difference in alignment (in base pairs) two parallel clones must have to be considered different (default is 5). It is enough for one member of the F/R pair to have alignment difference above this threshold for the clone to be considered different.
- maximum\_back\_clones:** Maximum number of other F/R pairs checked to determine if a given F/R pair is unique or not<sup>5</sup> (default is 50).
- weight\_first\_small:** Contribution of the first F/R pair of a library of `sclones` to the total weight of an arc (default is 1.00).
- weight\_other\_small:** Contribution of each F/R pair (other than the first) of a library of `sclones` to the total weight of an arc (default is 0.50).
- weight\_first\_medium:** Similar to `weight_first_small`, but for medium clones (clones that are neither `sclones` nor `lclones`) (default is 1.00).
- weight\_other\_medium:** Similar to `weight_other_small`, but for medium clones (default is 1.00).
- weight\_first\_large:** Similar to `weight_first_small`, but for `lclones` (default is 1.00).
- weight\_other\_large:** Similar to `weight_other_small`, but for `lclones` (default is 1.00).
- hybrid\_bonus:** Value added to the weight of hybrid arcs (default is 1.00).
- verbose:** Choose between silent mode (0) or verbose mode (1). Default is 0 (silent mode).

Not all tags need be in the configuration file: `genscaff` uses the default value for the parameters omitted. Note that `algorithm` and `verbose` can also be set directly in the command line. If there is a conflict between the command line and the configuration file, the former prevails.

---

<sup>5</sup>The purpose of this parameter is to avoid that pathological instances slow down the algorithm. In virtually all cases, there is no need to use a value other than the default.

### 6.3 Intermediary File

This file is in extended DIMACS format, and should contain the information derived from some previous assembly of the genome, as discussed in the text. There are three types of lines, identified by their first character. Lines beginning with **c** are reserved for comments and may be ignored. Lines with **v** represent contigs (the vertices of the graph). Each such line contains only two fields:

**v**  $\langle contig\ label \rangle$   $\langle contig\ length \rangle$

Links are represented in **a** lines, since they will constitute the arcs of the graph. Each line is made up by 14 fields,

**a**  $\langle c_1 \rangle$   $\langle c_2 \rangle$   $\langle r_1 \rangle$   $\langle r_2 \rangle$   $\langle R|L \rangle$   $\langle U|C \rangle$   $\langle length \rangle$   $\langle a_1 \rangle$   $\langle a_2 \rangle$   $\langle p_1 \rangle$   $\langle s_1 \rangle$   $\langle p_2 \rangle$   $\langle s_2 \rangle$   $\langle lib \rangle$ ,

meaning:

- $\langle c_1 \rangle$  and  $\langle c_2 \rangle$ : contig labels;
- $\langle r_x \rangle$ : string representing the read aligned with  $c_x$  (where  $x = 1$  or  $2$ );
- $\langle R|L \rangle$  and  $\langle U|C \rangle$ : together, these fields define the relative orientation of  $c_1$  and  $c_2$  induced by the link (see Section 3);
- $\langle length \rangle$ : lower bound on the length (in base pairs) of the clone (see below);
- $\langle a_x \rangle$ : alignment position (in bp) of the first aligned base (at the 5' end) of  $r_x$  with respect to  $c_x$ . Using  $c_x$ 's length and the orientation of  $r_x$  with respect to  $c_x$  it is possible to determine the distance  $d_x$  of the 5' end of  $r_x$  with respect to the far end of  $c_x$  (if  $r_x$  is not complemented with respect to  $c_x$ ,  $d_x = \text{length}(c_x) - a_x$ ; if  $r_x$  is complemented,  $d_x = a_x$ ). This enables us to compute  $\mathcal{L}$  as  $d_1 + d_2$  ( $\mathcal{L}$  is referred to in Section 3). Usually  $\langle length \rangle = \mathcal{L}$ , unless some other source is used for clone lower bound information.  $a_1$  and  $a_2$  are also used to detect duplicate clones.
- $\langle p_x \rangle$  and  $\langle s_x \rangle$ : primary and secondary alignment scores of  $r_x$  with respect to  $c_x$ ;
- $\langle lib \rangle$ : string representing the name of the library to which the clone belongs.

### 6.4 Library File

This file contains information about the libraries mentioned in the intermediary file. There are just two types of lines, each identified by its first character: **c** lines contain comments and may be ignored; **l** lines describe the libraries. Each **l** line contains the name of the library (an arbitrary string) and three integers representing clone sizes (in bp): minimum, mean and maximum. Minimum and maximum sizes are used to determine whether a clone is an **lclone** or an **sclone**. The mean size is used to estimate the length of the paths found by the algorithms, as described in Section 6.7. A typical **l** line looks like this:

l c01 30000 40000 55000

According to this line, library `c01` contains clones with sizes varying from 30 kb to 55 kb; the mean clone size in this library is 40 kb.

## 6.5 Complete Output File (`allscaff`)

Each line in the output file has its structure defined by the first character. Some types of lines appear in several sections:

- **f** (free line): doesn't have a specific format and should be ignored by parsers (free lines with no text are often used to format the output — there are no blank lines);
- **s** (section line): marks the beginning of a section (the file is divided into sections). Everything after **s** represents the name of the section, as in

```
s INPUT DATA
```

- **d** (data line): has a well-defined format, but it depends on the context (section) in which the line appears;
- **t** (tagged data line): has a well-defined format and contains a globally unique tag right after the letter **t** indicating what piece of information is represented. For instance,

```
t contigs 987
```

indicates that there are 987 contigs in the graph. Global uniqueness is useful to build simple parsers.

Other types of lines only occur in section `SCAFFOLDS`, which reports the scaffolds themselves: **c** (connected component), **p** (main path), **q** (alternate path), **a** (arc) and **l** (link). Their formats are described in Section 6.5.10.

In some sections, the output file reports information on individual arcs. Although an arc may contain more than one F/R pair, it is named after a single clone (the one with the highest uniqueness). In the case of hybrid arcs, an asterisk is appended to the name: while `AOQH63102` represents a simple arc (with one or more clones, but all of the same type), `AOQH63102*` represents a hybrid one.

The remainder of this section describes the individual sections of the output file.

### 6.5.1 INPUT DATA

This section summarizes the input information used to build the scaffolds. It begins with the following **t** lines:

- **file**: name of the intermediary file, as it appeared in the command line;
- **clones**: total number of clones;
- **sclones**: number of small clones;

- **mclones**: number of medium clones;
- **lclones**: number of large clones;
- **contigs**: total number of contigs read;

The remainder of the section reports (also in **t** lines) the values of the parameters described in 6.2: **algorithm**, **minimum\_uniqueness**, and so on. All parameters appear in the output file, regardless of whether they appear in the configuration file or not.

### 6.5.2 CLONE SIZES

This section contains the information read from the library file. There is a **d** line for each library, with four fields: a string representing the name of the library followed by the minimum, the mean and the maximum sizes of its clones. For example,

```
d 07 4000 5000 7000
```

indicates that clones in library 07 have sizes ranging from 4 kb to 7 kb and their average length is 5 kb.

### 6.5.3 UNRELIABLE CLONES

This section presents clones that should not be trusted upon, since they have negative uniquenesses (i.e., primary scores smaller than secondary scores). The first line in this section is a **t** line reporting the total number of unreliable clones (**unreliable**). Then, each such clone is reported in a **d** line with four fields: the name of the clone, the label of the contig to which it is unreliably linked, and the clone's primary and secondary scores. A typical section looks like this:

```
s UNRELIABLE CLONES
t unreliable 3
d AOJJ-ODG12-LA00 510 344 347
d AOQR5508D04 494 515 516
d AOJJ-OIA02-LA00 543 158 260
```

According to the first **d** line in this example, the link between clone AOJJ-ODG12-LA00 and contig  $c_{510}$  is unreliable because its primary score (344) is smaller than its secondary score (347).

### 6.5.4 INCOMPATIBLE LINKS

This section presents all pairs of contigs linked to each other in more than one way (i.e., with different pairs of ends involved). The first line in the section is a **t** line reporting the total number of contig pairs with incompatible links (**incompatiblelinks**). A series of **d** lines follows, each containing at least four fields and at most six. The first two are the labels of the contigs; the others represent all possible ways of linking the contigs found. Consider the following example:



```
d 310 396 LL RL
```

This line refers to contigs  $c_{310}$  and  $c_{396}$ . While at least one clone indicates that  $c_{310}^l$  is linked to  $c_{396}^l$ , there is at least another one in the input file stating that the link is between  $c_{310}^r$  and  $c_{396}^l$ . Clearly, they can't both be correct.

### 6.5.5 UNUSED F/R PAIRS

This section lists all F/R pairs discarded by the algorithm. As described in section 3, there are two reasons for an F/R pair to be discarded:

- it is a *putative duplicate*, i.e., it is very similar to other F/R pairs linking the same pair of contigs;
- its expected length is *inconsistent* with the length of other F/R pairs linking the same pair of contigs.

Each unused F/R pair is described in a **d** line with six fields. A typical example is:

```
d 432 L 456 R A0JJ0704A09 duplicate
```

According to this line, clone A0JJ0704A09, linking the left (L) end of  $c_{432}$  to the right (R) end of  $c_{456}$ , was discarded because it was a putative duplicate. If an F/R pair is discarded because of inconsistency with other pairs, the last field is **inconsistent** (instead of **duplicate**).

### 6.5.6 ARC WEIGHTS

This section reports the sizes of the arcs created (regardless of whether they are actually inserted in the graph or not). The first line is a **t** line with three fields (other than the tag itself, **arcs**): the number of arcs created, the weight of the lightest arc and the weight of the heaviest arc. This is followed by a series of **d** lines with two fields: a weight and the number of arcs with such weight. A typical output looks like this:

```
s ARC WEIGHTS
t arcs 749 1.0 4.0
d 1.0 526
d 1.5 37
d 2.0 174
d 2.5 8
d 3.0 3
d 4.0 1
```

In this example, there are 526 arcs with weight 1.0, 37 with weight 1.5, and so on.

### 6.5.7 HYBRID ARCS

This section reports the number of hybrid arcs, arcs that contain both `sclones` and `lclones`. Consistent and inconsistent hybrid arcs also have their number reported separately. (A hybrid arc is said to be *inconsistent* if some of its `lclones` are incompatible with its `sclones`). The information is shown in `t` lines with tags `hybridarcs`, `consistentlybrid` and `inconsistentlybrid`. A typical section is:

```
s HYBRID ARCS
t hybridarcs 18
t consistentlybrid 13
t inconsistentlybrid 5
```

### 6.5.8 INSERTION STATISTICS

This section contains only `t` lines. They account for the total number of arcs inserted (`inserted`) and not inserted (`notinserted`) into the graph. This last piece information is also presented in a more detailed form, with a `t` line for each of the four reasons for discarding an arc:

- **cycle**: the arc would create a cycle if inserted;
- **degree**: at least one end (GP) or both ends (MWP) of the arc were already connected to other arcs;
- **orientation**: the relative orientation of the ends the arc would induce is incompatible with the orientation induced by previously inserted arcs;
- **uniqueness**: the uniqueness of the arc is lower than the minimum threshold.

Note that `degree`, `uniqueness`, and `cycle` can occur in both algorithms; `orientation` may happen only in MWP. However, all four reasons are always listed, regardless of the algorithm. The quantities are expressed both as absolute values (first field after the tag) and as percentages (second field after the tag). A typical output is:

```
s INSERTION STATISTICS
t inserted 218 29.1
t notinserted 531 70.9
t cycle 1 0.1
t degree 399 53.3
t orientation 0 0.0
t uniqueness 131 17.5
```

### 6.5.9 ARCS NOT INSERTED

This section lists all arcs not inserted into the graph alongside with the reason why they were kept out. A typical line looks like this:

```
d AOJJ1388B03* 390 477 degree
```

The meaning is straightforward: the first field contains the name of the arc and the next two the contigs it is connected to. The fourth field is always one of the four reasons listed in the previous section: `cycle`, `degree`, `orientation`, or `uniqueness`.

When the reason is `cycle`, algorithm MWP adds an `f` line right after the `d` line to list the vertices and edges of the cycle. GP doesn't do this, since it can avoid cycles by using a union-find data structure; finding the actual cycles would be too costly. When the reason is `uniqueness` a fifth field gives the score of the rejected arc.

### 6.5.10 SCAFFOLDS

This is the section that actually presents the scaffolds. It begins with a `t` line reporting the number of connected components in the graph (`components`). Then the components themselves are listed in nonincreasing order of size (number of vertices/contig ends). Components with no arcs (i.e., only one contig) are omitted.

The description of a connected component starts with a `c` line with four fields: the label of the connected component (a sequential number starting at 1), the number of contigs, the number of arcs and the number of main paths built. An example:

```
c 2 117 151 31
```

This connected component (component 2) has 117 vertices (contig ends) and 151 arcs; 31 main paths were found by the algorithm (note that for GP the fourth field will always be 1). This line is followed by a description of the *main paths* in the connected component, listed in nondecreasing order by weight. If the algorithm is MWP, each main path is followed by a list of *links* to previous main paths and a list of *alternate paths*.

A typical path looks like this:

```
p 1 5 102.0 634591
a 480 C 58328 AOJJ-1CC03-LA00* 13.5
a 471 C 30387 AORN1574A10* 9.0
a 472 C 50408 AOJJ-1XF03-LA00 13.0
a 511 U 143452 AOJJ0701D04* 4.0
a 409 U 5417 AOJJ-1GE01-LA00 4.0
a 504 C 110999 AOQH1368F04* 28.5
a 494 U 71509 AOJJ-10G05-LA00* 30.0
a 513 U 136421 END 0.0
l 513 5 AOUT6815A09 424 1
l 419 4 AOJJ0713A10 472 5
```

```

q 274734 275713 C
a 511 U 143452 A0JJ-1YF06-LA00 3.0
a 428 U 6012 A0JJ-1ND12-LA00 3.0
a 504 C 110999 END 0.0

```

**Main path** According to the **p** line, this is the fifth heaviest path in its connected component (component 1). Its weight is 102.0 and its estimated length is 634,591 base pairs (for path length computation, see 6.7 below). The first field in each **a** line represents the label of a contig in the path. This particular one is made up by 8 contigs ( $c_{480}, c_{471}, c_{472}, \dots, c_{513}$ , in this order). The second field is either **C** or **U**, depending on whether the contig is complemented or not in the path (in the example,  $c_{480}, c_{471}, c_{472}$ , and  $c_{504}$  are complemented). The length of the contig (in base pairs) is reported in the third field. The name and the weight of the arc linking a pair of contigs are reported in the fourth and fifth fields, respectively. For instance,  $c_{511}$  and  $c_{409}$  are linked by arc **A0JJ0701D04\***, whose weight is 4.0. Note that the last **a** line is special; the word **END** replaces the arc name to mark the last contig of the path (to make parsing easier, a meaningless **0.0** is placed on the fifth field).

This particular example was produced with the **-V** (silent) parameter. The option **-v** (verbose) would print, right after every **a** line (with the exception of the last one), an **f** line with a list of all active clones in the arc. (Active clones are those that are neither duplicates nor inconsistent, and were thus used in arc weight computations.)

**Links** (MWP only) After the path itself is printed, its links to previously reported main paths are listed in **l** lines. In our example, two links are reported. Arc **A0UT6815A09** links  $c_{513}$ , which belongs to the current path (path 5), to  $c_{424}$  (path 1). The order in which the contigs appear is important. Were  $c_{424}$  to be inserted in path 5, it would appear *after*  $c_{513}$ , since  $c_{513}$  appears before  $c_{424}$  in the line that describes the link. The second link reported is between  $c_{419}$  (path 4) and  $c_{472}$  (current path). In this case,  $c_{419}$  would appear *before*  $c_{472}$  in path 5.

**Alternate paths** (MWP only) Finally, alternate paths (if any) are reported. Each one is introduced by a **q** line, which has three fields. The first two are the length of the alternate path ( $\ell_a$ ) and the length of the corresponding portion of the original main path ( $\ell_o$ ). The last one is either **C** (consistent) or **I** (inconsistent), depending on how  $\ell_o$  and  $\ell_a$  relate to each other. If they differ by no more than 20% ( $\ell_a$  being the base value), the path is said to be consistent; otherwise, the path will be inconsistent. (Actually, 20 is just the default percentage; the actual value, `maximum_alternate_difference`, is user-defined.) Each **t** line is followed by **a** lines describing the path (exactly like in the main path).

In our example, only one alternate path was found. Since its length is 274,734 bp and it corresponds to 275,713 bp in the original path, the path is consistent. The path contains three contigs,  $c_{511}$ ,  $c_{428}$ , and  $c_{504}$ . Note that the first and the last contigs of an alternate path must also belong to the original path, while all others must not.

## 6.6 Simplified Output File (scaff)

This file contains a subset of the information provided by the `allscaff` file. It contains just the scaffolds themselves. Although the format is very similar to the one described in section 6.5.10, there are a few differences that make the file a little longer, but easier to read. The path that illustrated section 6.5.10 is reproduced here in the friendlier format:

```
Path 5 (component 1): 634591 bp, weight 102.00
  C 480 (58328 bp) AOJJ-1CC03-LA00* (weight: 13.5)
  C 471 (30387 bp) AORN1574A10* (weight: 9.0)
  C 472 (50408 bp) AOJJ-1XF03-LA00 (weight: 13.0)
    511 (143452 bp) AOJJ0701D04* (weight: 4.0)
    409 (5417 bp) AOJJ-1GE01-LA00 (weight: 4.0)
  C 504 (110999 bp) AOQH1368F04* (weight: 28.5)
    494 (71509 bp) AOJJ-10G05-LA00* (weight: 30.0)
    513 (136421 bp)
```

```
Link: 513 [path 5] AOUT6815A09 424 [path 1]
```

```
Link: 419 [path 4] AOJJ0713A10 472 [path 5]
```

```
Alternate path (274734 bp, original has 275713 bp)
```

```
  511 (143452 bp) AOJJ-1YF06-LA00 (weight: 3.0)
  428 (6012 bp) AOJJ-1ND12-LA00 (weight: 3.0)
  C 504 (110999 bp)
```

## 6.7 Path Length Computation

The length of a path is given by the sum of component contig lengths plus the sum of gap length estimates. Gap length is estimated as follows. In a given arc, there can be one or more F/R pairs. For each F/R pair corresponding to a clone  $p$ , its gap length is estimated as  $\max(\text{mean\_lib\_size}(p) - \text{length}(p), 0)$ , where  $\text{mean\_lib\_size}$  is the mean clone size for the clone library to which  $p$  belongs (information contained in the library file) and  $\text{length}$  is the lower bound on the size of  $p$ , given in the F/R input line for clone  $p$ . The gap length for an arc is taken as the mean of all individual F/R estimates for that arc.