

Concept

MAIN IS NOT TRUE

NGS NGSnext-generation-sequencing

DBG DBGde-Bruijn graph

OLC OLCoverlap-layout-consensus

1.1 Pacbio reads

Daligner finds alignments between long, noisy reads. Pacific Biosciences has commercially launched its first sequencer in 2011. It is able to output reads with an average of 1000 bases, which is significantly longer than NGS (NGS) reads [?]. In 2014, a new polymerase-chemistry combination was released, called P6-C4. This version can output average read lengths of 10000-15000 bases, and its longest reads can exceed 40000 bases [?]. While the drawback is that these reads have an error rate of 12-15%, this can be compensated by the distribution of these errors [?]. First, the set of reads is a nearly Poisson sampling of the sampled genome. This implies that there exists a coverage c for every target coverage k , such that every region of the genome is covered k times [?]. Secondly, the work of Churchill and Waterman [?] implies that the accuracy of the consensus sequence of k sequences is $O(e^{-k})$, which goes to 0 as k increases. This means that if the reads are long enough to handle repetitive regions, in principle a near perfect de novo assembly of the genome is possible, given enough coverage.

Important points for de novo DNA sequencing are: what level of coverage is needed for high quality assembly? And how to build an assembler that is able to deal with high error rates and long reads? Most previous assemblers work with NGS reads, which are much shorter and have much lower error rates. Some algorithms used in these assemblers, such as DBG (DBG) [?] would grow too large for high error rates and long reads. Since Daligner was build, new methods of using DBG with long reads have been developed, but they rely on a short read based DBG to correct errors in long reads [?][?].

1.2 Daligner

The first step in an OLC (OLC) assembler is usually finding overlaps between reads [?]. BLASR [?] was the only long read aligner at the time, and inspired Daligner. It uses the same filtering concept, but with a cache-coherent threaded radix sort to find seeds, instead of a BWT index [?]. The most time-consuming step is extending the seed hit to find an alignment. To do this, Daligner uses a novel method which

adaptively computes furthest reaching waves of the older $O(nd)$ algorithm [?], combined with heuristic trimming and a datastructure that describes a sparse path from the seed hit to the furthest reaching point.

Daligner performs all-to-all comparison on two input databases \mathcal{A} , with M long reads A_1, A_2, \dots, A_M and \mathcal{B} , with N long reads B_1, B_2, \dots, B_N over alphabet $\Sigma = 4$. It reports alignments $P = (a, i, g)x(b, j, h)$ such that $len(P) = ((g - i) + (h - j))/2 \geq \tau$ and the optimal alignment between $A_a[i + 1, g]$ and $B_b[j + 1, h]$ has no more than $2\epsilon \cdot len(P)$ differences, where a difference can be either an insertion, a deletion or a substitution. Both τ and ϵ are user settable parameters, where τ is the minimum alignment length and ϵ the average error rate. The correlation, or percent identity of the alignment is defined as $1 - 2\epsilon$.

An edit graph for read $A = a_1a_2\dots a_m$ and $B = b_1b_2\dots b_n$ is a graph with $(m+1)(n+1)$ vertices $(i, j) \in [0, M] \times [0, N]$. It also has three types of edges:

- deletion edges $(i - 1, j) \rightarrow (i, j)$ with label $\begin{bmatrix} a_i \\ - \end{bmatrix}$ if $i > 0$.
- insertion edges $(i, j - 1) \rightarrow (i, j)$ with label $\begin{bmatrix} - \\ b_j \end{bmatrix}$ if $j > 0$.
- diagonal edges $(i - 1, j - 1) \rightarrow (i, j)$ with label $\begin{bmatrix} a_i \\ b_j \end{bmatrix}$ if $i, j > 0$.

An alignment between $A[i + 1, g]$ and $B[j + 1, h]$ is described as a sequence of labels from vertex (i, j) to (g, h) . A diagonal edge can be either be a match edge, when $a_i = b_j$, or a substitution edge. If a match edge has weight 0, and the other edges have weight 1, the weight of the total path is the number of differences in the alignment it represents. To find suitable alignments, we have to find a read subset pairs P such that $len(P) \geq \tau$ and the weight of the lowest scoring path between (i, j) and (g, h) in the edit graph of A_a and B_b is not more than $2\epsilon \cdot len(P)$.

The $O(ND)$ algorithm tries to find progressive waves of furthest reaching (f.r.) points until the endpoint is reached. The goal is to find longest possible paths starting at a starting point $\rho = (i, j)$ with 0 differences, then 1 difference, then 2 and so on. After d differences, the possible paths can end in diagonals $\kappa \pm d$, where $\kappa = i - j$ is the diagonal of the starting point. The furthest reaching point on diagonal k that can be reached from ρ with d differences is called $F_\rho(d, k)$. A collection of these points for a particular value of d is called the d -wave emanating from ρ , and defined as $W_\rho(d) = \{F_\rho(d, \kappa - d), \dots, F_\rho(d, \kappa + d)\}$. $F_\rho(d, k)$ will be referred to as $F(d, k)$, where ρ is implicitly understood from the context.

In the $O(ND)$ paper it is proven that:

$$F(d, k) = Slide(k, \max\{F(d-1, k-1) + (1, 0), F(d-1, k) + (1, 1), F(d-1, k+1) + (0, 1)\}) \quad (1.1)$$

where $Slide(k, (i, j)) = (i, j) + \max\{\Delta : a_{i+1}a_{i+2}\dots a_{i+\Delta} = b_{j+1}b_{j+2}\dots b_{j+\Delta}\}$.

A slide is a path of sequential match edges. The f.r. d -point on diagonal k is calculated by finding the furthest of

- the f.r. $(d-1)$ -point on $k - 1$ followed by an insertion

- the f.r. $(d-1)$ -point on k followed by a substitution
- the f.r. $(d-1)$ -point on $k + 1$ followed by a deletion

and then continuing as far as possible along the slide. A point (i, j) is furthest when its anti-diagonal $i + j$ is greatest. The best alignment between reads A and B is the smallest d such that $(m, n) \in W_{(0,0)}(d)$, where m and n are the length of reads A and B. The $O(ND)$ algorithm computes d -waves from starting point $(0, 0)$ until the end point (m, n) is reached. The complexity of this algorithm is $O(n + d^2)$ when A and B are non-repetitive sequences [?]. Because seeds are not always at the beginning, so waves are computed in both forward and reverse direction. The latter is easily done by reversing the direction of edges in the edit graph.

1.3 Seeding

To find suitable starting points for the edit graphs, seeding is done. A seed is a section where $A[i, g]$ and $B[j, h]$ have a certain high similarity that indicates that these reads probably originate from the same part of the genome. Finding a seed means finding matching k -mers for every readpair (a, b) with $a \in \mathcal{A}$ and $b \in \mathcal{B}$. Previous methods to match k -mers include Suffix Arrays [?] and BWT indices [?].

Daligner uses a sort-merge procedure:

- Build a list $List_X = \{(kmer(X_x, i), x, i)\}_{x,i}$ of all k -mers for database $X \in \{\mathcal{A}, \mathcal{B}\}$, where $kmer(R, i)$ is the k -mer $R[i - k + 1, i]$.
- Sort both lists in order of their k -mers.
- Merge the two lists and build $List_M = \{(a, b, i, j) : kmer(A_a, i) = kmer(B_b, j)\}$ of read and position pairs that have the same k -mer.
- Sort $List_M$ lexicographically on a, b and i where a is most significant.

All entries for a certain read pair (a, b) are in a continuous segment of the list. This makes it easy to determine if that read pair has enough k -mers and in the right places to constitute a seed hit. Given parameters h and s , each entry (a, b, i, j) for the current read pair is placed in diagonal bands $d = \lfloor (i - j)/2^s \rfloor$ and $d + 1$. Now determine the number of bases in the A-read covered by k -mers in each pair of neighbouring diagonal bands.

For Daligner, the default k is 14, and assumed error rate is 0.85.