

TITLE1
TITLE2
TITLE3

Tong Dong Qiu

CE-MS-2017

Abstract

DALIGNER or Daligner

TITLE

THESIS

submitted in partial fulfillment of the
requirements for the degree of

MASTER OF SCIENCE

in

COMPUTER ENGINEERING

by

AUTHOR

born in PLACE, COUNTRY

Computer Engineering
Department of Electrical Engineering
Faculty of Electrical Engineering, Mathematics and Computer Science
Delft University of Technology

TITLE

by AUTHOR

Abstract

Laboratory : Computer Engineering
Codenummer : CE-MS-2017-number

Committee Members :

Advisor: , CE, TU Delft

Chairperson: , CE, TU Delft

Member: , CE, TU Delft

Member: , CE, TU Delft

Dedicated to my family and friends

Contents

List of Figures

List of Tables

Acknowledgements

AUTHOR

Delft, The Netherlands

September 19, 2017

MAIN IS TRUE

3.1 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 $O(\epsilon^k)$, 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 [?][?].

3.2 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 $\text{len}(P) = ((g - i) + (h - j))/2 \geq \tau$ and the optimal alignment between $A_a[i + 1, g]$ and $B_b[j + 1, h]$ has no more than $2\epsilon \cdot \text{len}(P)$ 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 $\text{len}(P) \geq \tau$ 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 $2\epsilon \cdot \text{len}(P)$.

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:

$$F(d, k) = \text{Slide}(k, \max\{F(d-1, k-1) + (1, 0), F(d-1, k) + (1, 1), F(d-1, k+1) + (0, 1)\}) \quad (3.1)$$

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

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.

3.3 Seeding: concept

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 includes 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 [?].

Assuming that the k -mer matches are independent, the probability that a k -mer is conserved while sequencing is $\pi = (1 - 2\epsilon)^k$. The number of conserved k -mers in an alignment of τ basepairs is a Bernoulli distribution with rate π , so an average of $\tau \cdot \pi$ k -mers are expected in this alignment. An example: $k = 14$, $\epsilon = 15\%$ and $\tau = 1500$, then $\pi = .7^{14} = 0.0067$ and the average number of conserved k -mers is 10. Only .046% of the expected readpairs have 1 or fewer k -mers, and only 0.26% have 2 or fewer. To filter with 99.74% sensitivity, only readpairs with 3 or more k -mer matches need to be examined.

The specificity of the filter is increased in two ways:

- computing the number of k -mer matches in diagonal bands of width 2^s instead of in the whole reads
- thresholding on the number of bases in k -mer matches, instead of the number of k -mers themselves

The first way decreases the false positive rate because it only allows readpairs that have their k -mer matches relatively close, indicating a smaller region with higher similarity. The second way relies on the fact that 3 overlapping k -mers have a higher probability ($\pi^{1+2/k}$) than 3 disjoint k -mers with $3k$ basepairs (π^3).

The actually find the k -mer matches, 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. Note that only bases in matching k -mers are counted, not the number of k -mers. When there are $k+1$ matching consecutive bases, two k -mers are generated. These are less 'valuable' than two non-overlapping k -mers. If $Count(d) \geq h$ for any diagonal band d , there is a seed hit for each position (i, j) in the band d unless position i was already in the range of a previously calculated local alignment.

The best values for h and s depend on things like ϵ and the read lengths.

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

3.4 Seeding: implementation

Daligner is designed to use multiple threads and use the cache efficiently. Building and merge the lists in steps 1 and 3 is easy, since only one pass is needed for both actions. The elements of the lists are compressed into 64-bit integers. Daligner uses a radix sort [?][?] to sort the lists in steps 2 and 4. Each 64-bit integer is a vector of $P = \lceil hbits/B \rceil$, B -bit digits $(x_P, x_{P-1}, x, \dots, x_1)$ where B is a parameter. The sort needs P passes, where each pass sorts the elements on a B -bit digit x_i . Each pass is done with a bucket sort [?] with 2^B buckets. Instead of a linked list, the integers in the list src are moved in precomputed segments $trg[bucket[b] \dots bucket[b+1] - 1]$ of an array $trg[0 \dots N-1]$ with the same size as src . For the p^{th} pass, $bucket[b] = \{i : src[i]_p < b\}$ for each $b \in [0, 2^B - 1]$.

```
for i = 0 to N-1 do
{
    b = src[i]_p
    trg[bucket[b]] = src[i]
    bucket[b] += 1
}
```

The algorithm takes $O(P(N + 2^B))$ time, but B and P are small fixed numbers so it is effectively $O(N)$. There are a lot of parallel sorting algorithms [?][?], but Daligner uses a new method that needs half the number of passes that traditional methods use. Each thread sorts a contiguous segment of size $part = \lceil N/T \rceil$ of src into trg , where T is the number of threads. This means each thread $t \in [0, T-1]$ has a bucket array $bucket[t]$ where $bucket[t][b] = \{i : src[i] < b \text{ or } src[i] = b \text{ and } i/part < t\}$. To reduce the number of passes, a bucket array for the next pass is filled during the current pass. Each thread counts the number of B -bit digits that will be handled in the next pass by itself and every other thread separately. If the number at index i will be at index j and in bucket b next pass, then the count in the current pass must not be recorded for the thread $i/part$ that currently sorts the number, but for thread $j/part$ that will sort it in the next pass. To do this we need to count the number of these events in $next[j/part][i/part][b]$ where $next$ is a $T \times T \times 2^B$ matrix. If $src[i]$ is about to be moved in the p^{th} pass, then $j = bucket[src[i]_p]$ and $b = src[i]_{p+1}$.

Experiments

5

Bibliography

- [1] K. Davies, “Get smrt: Pacific biosciences unveils software suite with commercial launch,” April 2011. [Online]. Available: <http://www.bio-itworld.com/news/04/29/2011/Pacific-Biosciences-software-commercial-launch.html>

List of definitions

sensitivity True Positive rate, the portion of positives that are correctly identified.

specificity True Negative rate, the portion of negatives that are correctly identified.

A
