

Read alignment against a graph-based reference genome

Fuzzy searching in large and complex structures

Esten Høyland Leonardsen

Master's Thesis Spring 2016



Read alignment against a graph-based reference genome

Esten Høyland Leonardsen

7th March 2016

Abstract

Contents

I	Introduction	1
1	Algorithm	3
1.0.1	Definitions	3
1.0.2	The graph	4
1.0.3	Overview	5
1.0.4	Precomputation of the graph	5
1.0.5	Generating the new graph G'	7
1.0.6	Searching G' with the modified PO-MSA	9
1.0.7	Handling invalid threshold values	10
1.0.8	Merging	10

List of Figures

1.1	A graph made from the three sequences “ATATA”, “AGAGA” and “ACAA” with 9 valid complete paths	5
1.4	The left suffix tree corresponding to the graph in 1.3	6
1.3	A small reference graph with left contexts (top) and right contexts (bottom) of length 2 shown	7
1.5	The resulting candidate sets for mapping the string "ATA" against the reference genome from 1.3 with varying T values .	8
1.6	The 4 arrays used by the searching algorithm when using the candidate sets from Fig ?? and $T = 1$	10

List of Tables

Preface

Part I

Introduction

Chapter 1

Algorithm

In this section I will present the algorithm “Fuzzy context-based search” as an approach for mapping text strings against graph-based reference genomes. First the elements involved will be described in order to precisely define the problem. Then the syntax of the reference graphs used are explained in more detail. Finally, the algorithm is presented both through a conceptual overview and in more specific detail. The implementation details refers to the tool *Graph Genome Alignment* (See Supplementary XXX).

1.0.1 Definitions

Definition 1 (Graph-based reference genome (graph))

A pair $G = \{V, E\}$ where V is a set of vertices and E is a set of edges. $|G|$ denotes the number of vertices of G .

Definition 2 (Vertice)

A pair $v = \{b, i\}$ where $b \in \{A, C, T, G\}$ and i is a unique index. Every graph G also has two special vertices $s_G = \{s, 0\}$ and $t_G = \{e, -1\}$ which represents unique start and end vertices.

Definition 3 (Edge)

An ordered pair $e = \{i_s, i_e\}$ where both elements are indexes for vertices.

Definition 4 (Complete Path)

An ordered list P of indexes such that for all consecutive ordered pairs $\{i_x, i_{x+1}\} \in P$ there exists an edge $e = \{i_x, i_{x+1}\}$.

Definition 5 (Path)

An ordered list L of indexes such that for all consecutive ordered pairs $\{i_x, i_{x+1}\} \in L$ there exists a complete path P which starts at $\{i_x\}$ and ends at $i_{x+1}\}$.

Definition 6 (Input sequence)

A string s over the alphabet $\{A, C, T, G\}$. The length of the string is given by $|s|$. An individual character on position x is referenced by s_x

Definition 7 (Mapping score)

A score produced by mapping two characters $c_1, c_2 \in \{A, C, G, T\}$ against a scoring matrix

Definition 8 (Path score)

A score produced by traversing a path P through a graph G to create a linear sequence, scoring gaps according to the gap penalties given by a scoring schema.

Definition 9 (Alignment)

Given a sequence s and a graph G , an ordered list A of indexes such that every $a_x \in A$ is either a valid index for a vertex in G or 0. 0 indicates an unmapped element of the input sequence

Definition 10 (Alignment score)

Given a sequence s , a graph G and an alignment A , the score produced by combining mapping scores for the pairs $\{a_x, s_x\}$ for $0 \leq x < |s|$ with the path score for the path(s) provided by A aligned against both G and s .

Definition 11 (The optimal alignment score problem)

*For any pair $\{G, s\}$, where G is a graph and s is a sequence, find the alignment A which produces the highest possible alignment score. **If multiple max scores: Provide all or chose one?***

Definition 12 (The bounded optimal alignment score problem)

Given a triplet $\{G, s, T\}$ where G and s are as before and T is a numeric value, find the alignment A which produces the highest alignment score, iff the score for A is higher than T . If no such alignment exists, s is unmappable.

1.0.2 The graph

As defined the graphs involved will be graphs where one vertex represents a single nucleobase. Every vertex also contains an identifying index, which maps uniquely to that vertex. A graph G is made by starting out with only the start and end vertices, and iteratively merging in new sequences. Every sequence merged into the graph has a corresponding *complete path* starting in s_G and ending in t_G . There is no correspondence the other way, meaning there can be complete paths from s_G to t_G which does not originate from a single sequence (See fig. ??). There exists no information storing the origin of an edge and all paths are thus seen as equally probable when aligning a sequence. How the new sequences are merged is defined entirely through the alignment procedure which relies in part on the scoring threshold λ and the given scoring schema.

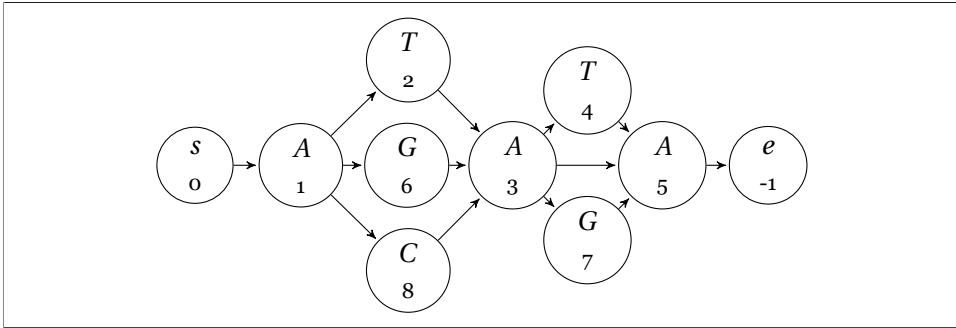


Figure 1.1: A graph made from the three sequences “ATATA”, “AGAGA” and “ACAA” with 9 valid complete paths

1.0.3 Overview

“Fuzzy context-based search” is an algorithm which solves the bounded optimal alignment score problem. The first step of the process is to build a searchable index based on the given graph G . This index is independent from the input sequences to be aligned, and can thus be used for several searches. The alignment itself consists of two steps: Building a new graph G' , from both G and an input string s , and search for an alignment in this newly formed graph. Searching for an alignment means combining nodes, representing bases, into a path which represents a linear sequence. This linear sequence can be aligned against the input sequence with regular string alignment tools and is therefore easily scorable. If the algorithm finds an alignment this is guaranteed to be one of the paths which produces the highest possible alignment score for any path in the graph (See supplementary XXX PROOF). There are some situations where the algorithm results in an empty alignment. These cases will occur when there are no paths in the graph which produces an alignment score higher than the threshold T , and the sequence s is identified as unmappable. When an empty alignment is provided as a basis for merging a new sequence into the graph, this results in a new complete path which is separated from the original vertices (See fig. ??)

1.0.4 Precomputation of the graph

There are two data structures needed for aligning any string against the graph, a suffix tree for left contexts and a suffix tree for right contexts. Before either of the two are built the algorithm needs to decide a length for the contexts. Currently in the tool there are two ways of setting the context length: A user given parameter or an approximation based on the probability of sharing contexts SUPPLEMENTARY?. The length of a context does not impact the quality of the alignments found by the algorithm PROOF IN SUPPLEMENTARY? but will have an impact on the runtime.

When a context length $|c|$ is set, the algorithm can start building the index. Two sets of strings, a left context set and a right context set, is generated for every node in the graph G . The generation of the two sets happen by the same procedure, by swapping around the starting point and the direction of the iteration. When creating left contexts the algorithm starts in the start-node of G and traverses following the direction of the edges, for right contexts the opposite is done. Apart from this the two are equal. To generate the context set $c(n_x)$ for a given node n_x the algorithm looks at every string $c \in c(n_y)$ for every incoming neighbouring node n_y . Every c is modified into a new context string c' by trimming away the last character and prefixing the context with **the character b stored in n_y** . All the generated strings c' is added to $c(n_x)$. As sets per definition does not allow duplicates the impact of a branching occurring in the graph will fade away after exactly $|c|$ steps as the difference is trimmed away (see Fig. 1.3), and thus avoid explosive exponentiality in the context set sizes.

The iteration starts in the node defined as the starting point which has the empty string ϵ as its only context. Whenever a node has finished producing its contexts it enqueues everyone of its outgoing neighbours in a regular FIFO queue. If a node has more actual incoming neighbours than incoming neighbours which are finished generating contexts, the node puts itself back in the queue. The algorithm halts when the queue is empty. Every node has to be visited exactly once to generate its context and as the procedure runs twice to generate both sets the total runtime for the operation is $O(2|G|)$.

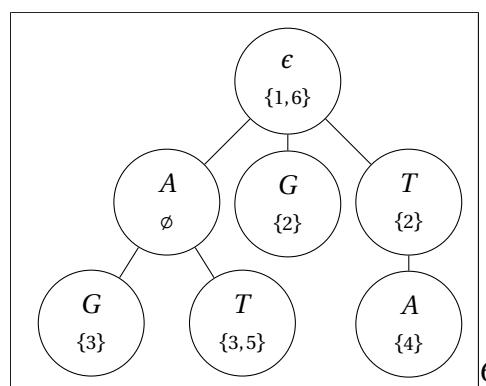
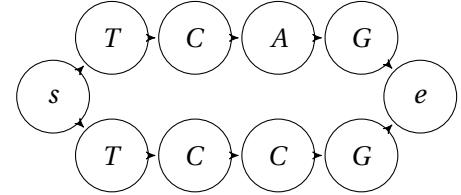


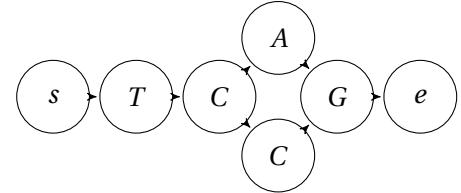
Figure 1.4: The left suffix tree corresponding to the graph in 1.3

	T	C	A	G
T	0	-1	-2	-3
C	-1	0	-1	-2
C	-2	-1	-1	-2
G	-3	-2	-2	-1

(a) The dynamic programming table for aligning the two strings “TCAG” and “TCCG” using modified edit distance as a scoring schema



(b) The result of aligning and merging the sequences with $T = 0$



(c) The result of aligning and merging the sequences with $T = 1$



After generating the two context sets for every node, the elements of each one is inserted into their corresponding suffix tree. In theory every node can have $4^{|c|}$ contexts in each set. When the graph is more or less linear with few branches a more fair approximation is $B * |c|$ where B is the observed branching factor. The current implementation

in the tool uses a naive suffix tree where insertion is $O(|c|)$. This is done for every node in the graph, yielding a total time complexity of $O(|G|B|c|^2)$. A discussion on more efficient suffix structures can be found in **SOMEWHERE IN DISCUSSION**.

Every suffix is stored as a key with the index of its originating node as a value. The total runtime for building a searchable index for a graph is $O(3|G||c|^2B)$

1.0.5 Generating the new graph G'

Creating G' is the process of determining which nodes qualifies as candidate nodes for a given input string s and combining them correctly. In order to determine actual candidates for the given string, the algorithm needs to know how much *fuzziness* to allow. This is a measure which decides how different a read can be from its optimal counterpart in the graph before it is categorized as not mappable. The algorithm takes in a fuzziness parameter λ and decides this by setting a threshold $T = \maxScore(s) - \lambda$. The maximal score is found by mapping the string, be it the entire input string or a context string, against itself with a scoring function provided by the scoring schema. Both λ and T is used throughout the entire process as cutoff variables.

For every character $s_x \in s$ a left-context string and a right-context string is generated by looking at the $|c| + \maxPossibleGapGivenThreshold(T)$ surrounding characters. The two strings are treated as contexts, one left and one right, and used as a basis for a fuzzy search in its corresponding suffix

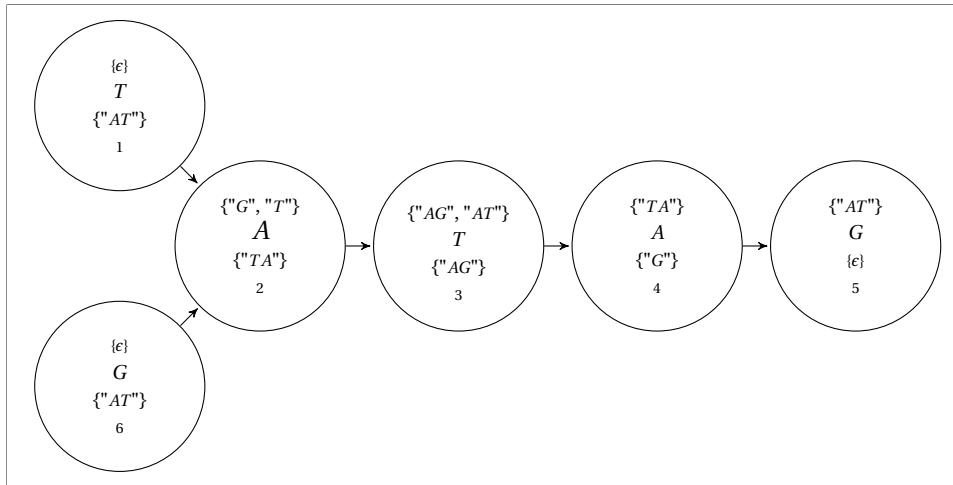
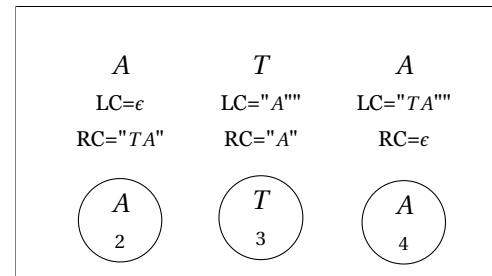


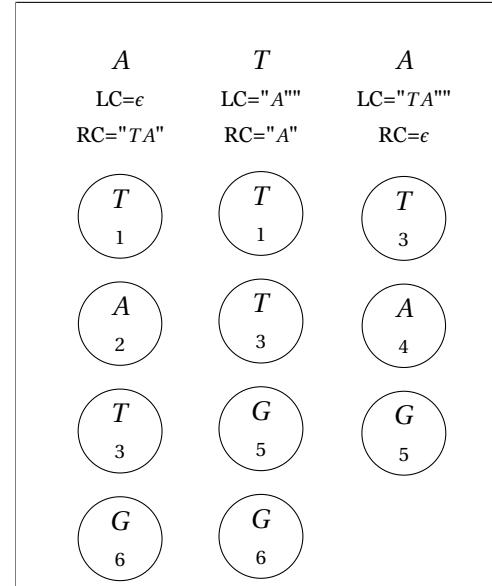
Figure 1.3: A small reference graph with left contexts (top) and right contexts (bottom) of length 2 shown

tree. The search is a recursive function based on PO-MSA. The root node is supplied with a one-dimensional scoring array corresponding to the context string c , which is initialized with all zeroes. Then, for every child, a new scoring array is computed by regular edit distance rules: For each index i take the maximal score for either a gap in the graph, a gap in the string or matching the character c_i with the character contained in the child node (Reference actual code in supplementary?), (more explanation needed?). This newly created array is supplemented to the same recursive function in the child. When a leaf node is reached the last index of the supplied scoring array corresponds to mapping the entire string c against the entire context achieved by concatenating the characters contained in the path traversed by the recursion. If the score is higher than the threshold T every index contained in the node is stored as a pair on the form $\{index, score\}$ to the candidate set. If an index is stored several times, only the pair containing the highest score is saved.

After the fuzzy search is concluded there are two sets of candidates for every index, one containing the nodes matching the left context and an equivalent for nodes matching the right context. These two sets are intersected to produce a final candidate set for the index i , where the score is created by adding together the scores from the two original. When the intersection happens the final set can again be pruned by removing all vertices which has a combined score that is lower than the sum of the maximal score for each context minus λ . When the vertices are found the edges need to be generated in order to finish the graph. Intuitively there should be an edge wherever there is a gap which is traversable without having the gap penalty exceeding λ . In practice this is a step which is done during the next step of the algorithm.



(a) $T = 0$



(b) $T = 1$

Figure 1.5: The resulting candidate sets for mapping the string "ATA" against the reference genome from g.3 with varying T values

The newly formed graph G' can be defined formally:

$G'(G, s, T) = \{V', E'\}$ where V' is an ordered set of sets of length $|s|$ where each set V'_i is a set of nodes such that

$$V'_i = \{v_x \mid v_x \in G \wedge \text{contextScore}(v_x, s_i) \geq T\}$$

and E' is a list of weighted edges such that

$$\begin{aligned} E' = \{e' = \{i_s, i_e, w\} \mid & i_s \in V'_x \wedge i_e \in V'_y \wedge \text{gapPenalty}(y - x) \leq \lambda \\ & \wedge w = \text{distance}(i_s, i_e) \wedge \text{gapPenalty}(w) \leq \lambda\} \end{aligned}$$

where $\text{contextScore}(x, y)$ and $\text{gapPenalty}(x, y)$ are scoring functions provided by the scoring schema and $\text{distance}(x, y)$ is the distance of the shortest path from node x to node y in the graph. (**Mixing up nodes and indexes in the definitions**)

In theory every leaf node has to be visited in order to check the score for every represented context in the tree. In practice the tree can be pruned by cutting off the search whenever the *maximal potential score*(clumsy name) falls below the threshold. The maximal potential score for a node is found by adding together the currently highest score in the scoring array with the maximal matching score for the remainder of the string. This reduces the number of nodes to be searched from $O(4^c)$ to (**something alot smaller. Needs calculations**).

1.0.6 Searching G' with the modified PO-MSA

When the candidate nodes for each position has been chosen there remains one final step: How to combine them in the best possible way to produce an actual path. This corresponds with finding the path through G' which traversal gives an optimal score within the scoring schema. Conceptually this is in many ways similar to a regular PO-MSA search where instead of searching through the reference graph with an input string we are searching through a graph of the string with the graph as our input. Instead of giving every node a score for every index in the string we give every index of the string a score for every candidate node. This presents one problem: We need to decide the gap score for jumping between arbitrary nodes in our graph. In the regular case with strings this is easily done by taking the differences in the indexes of two elements, whereas for graphs we need to find the distance between nodes through searching for shortest paths. We do however know we are not interested in alignments which does not have a score within our given threshold, which means we can bound the search to every node which is reachable without letting the gap penalty fall below λ .

A	T	A
1	1	3
2	2	4
3	3	5
6	4	X
X	5	X
X	6	X

A	T	A
-1	-2	-2
0	-2	0
-1	0	-2
-1	-2	X
X	-3	X
X	-3	X

(a) Indexes		
A	T	A
F	F	F
F	F	F
F	F	T
F	T	X
X	T	X
X	F	X

(b) Scores		
A	T	A
-1:-1	0:1	1:2
-1:-1	0:1	1:2
-1:-1	0:2	1:2
-1:-1	0:2	X
X	0:3	X
X	0:4	X

(c) Gaps		

(d) Backpointers		

Figure 1.6: The 4 arrays used by the searching algorithm when using the candidate sets from Fig ?? and $T = 1$

time. The final score for a node is the maximal attainable score found through searching these pairs. In order to store the scores we need an array of integer arrays the same size as the array of arrays of candidate nodes. To handle non-linear gap penalties we need a boolean array the same size to store gap-information. Lastly, in order to backtrack to find the path yielding the highest score we need a backpointer array consisting of pairs of integers. The resulting space complexity is $O(5\text{SIZEOF CANDIDATESETS})$.

The average case time complexity of the search is $O(\frac{|s|T^4c^2|G|^2}{4^{l_c^2}})$ Reference supplementary.

1.0.7 Handling invalid threshold values

1.0.8 Merging

The implemented search solves the problem through dynamic programming. For every index i look at every node n_x in the candidate set. For each of these nodes make pairs with every node n_y of every candidate set for the preceding indexes j . Once again we know gaps with a penalty larger than λ is not interesting and can bound the distance of the search backwards. For each of these pairs make a score by combining the score in the previous node with the mapping score for the character b stored in n_x mapped against s_i as well as the gap penalty for traversing from index j to i and for traversing the graph from n_x to n_y . For all scoring schemas where the penalty for a gap is larger than or equal to the worst possible mapping score between bases, the two preceding cases will never happen at the same

Bibliography

- [1] Deanna M. Church et al. ‘Extending reference assembly models’. In: (2015).
- [2] Alexander Dilthey et al. ‘Improved genome inference in the MHC using a population reference graph’. In: (2015).
- [3] Zamin Iqbal et al. ‘De novo assembly and genotyping of variants using colored de Bruijn graphs’. In: (2012).
- [4] Schneeberger K. et al. ‘Simultaneous alignment of short reads against multiple genomes’. In: (2009).
- [5] Birte Kehr et al. ‘Genome alignment with graph data structures: a comparison’. In: (2014).
- [6] Christopher Lee, Cathrine Grasso and Mark F. Sharlow. ‘Multiple sequence alignment using partial order graphs’. In: (2001).
- [7] Arthur M. Lesk. *Introduction to Bioinformatics*. Oxford University Press, 2014.
- [8] Shoshana Marcus, Hayan Lee and Michael Schatz. ‘SplitMEM: Graphical pan-genome analysis with suffix skips’. In: (2014).
- [9] Joong Chae Nal et al. ‘Suffix Array of Alignment: A Practical Index for Similar Data’. In: (2013).
- [10] Ngan Nguyen et al. ‘Building a Pan-Genome Reference for a Population’. In: (2015).
- [11] Adam Novak et al. ‘Canonical, Stable, General Mapping using Context Schemes’. In: (2015).
- [12] Benedict Paten, Adam Novak and David Haussler. ‘Mapping to a Reference Genome Structure’. In: (2014).
- [13] PA. Pevzner, H. Tang and MS. Waterman. ‘An eulerian path approach to DNA fragment assembly’. In: (2001).