

SyntenyFinder: A Synteny Blocks Generation and Genome Comparison Tool

Intern: Ilya Minkin

Advisor: Son Pham

We present algorithm for finding synteny blocks in genomes represented as nucleotide sequences. The algorithm is based on colored De Bruijn graphs and graph simplifications. Our method is able to reconstruct synteny blocks previously found by ad hoc methods.

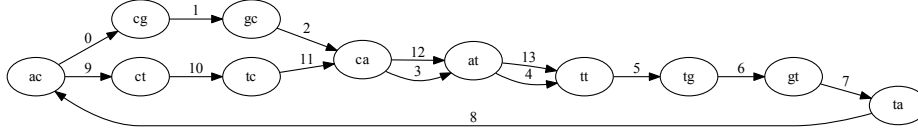
1. Introduction

Recent growth of number of sequenced genomes arises question about their evolution relationship. In order to perform rearrangement analysis, the genomes must be decomposed into conservative segments, called synteny blocks. Currently existing tools for solving this problem, like DRIMM-Synteny [1], require the genomes to be presented as sequences of enumerated local alignments, or *anchors*. Usually, anchors represent homologous genes. At this moment, there are no general purpose tools that can find synteny blocks from the genomes represented as unannotated nucleotide sequences.

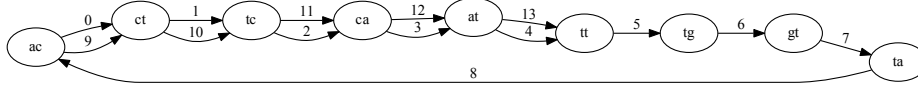
De Bruijn graphs are extensively used in bioinformatics for genome assembly [2, 3]. In this work we address problem of finding synteny blocks from nucleotide sequences. We propose new algorithm for this task based on colored De Bruijn graph.

2. Problem definition

Suppose that we are given a set $S = \{s_1, s_2, \dots, s_n\}$ of chromosomes, and each chromosome is represented as a string over alphabet $\{A, C, G, T\}$. The task of finding synteny blocks is to find a set of so called conserved regions $C = \{C_1, C_2, \dots, C_n\}$, where each conserved region C_i is a set of substrings of chromosomes from S . Such regions are supposed to cover most of the genome for closely related species. All substrings forming a conserved region C_i must be "similar" to each other according to some criterion



(a) De Bruijn graph built from string "acgcattgtactcatt" and $k = 2$. Non-branching paths correspond to multiple copies of the same substrings.



(b) Same De Bruijn graph after simplification. Replacing "acGca" by "acTca" we obtain long non-branching path that corresponds to the syntenic block.

Figure 1: Illustration of De Bruijn graphs and simplification

of similarity. Note that problem of finding syntenic blocks in a set of chromosomes is equivalent to a problem of finding syntenic blocks in one "superchromosome" obtained from concatenating all chromosomes from the set. At this moment there is no generally accepted formal criterion of similarity exist, so the problem of finding syntenic blocks is ill-defined.

In our work we introduce new criterion of similarity based on De Bruijn graphs and graph simplifications.

3. General idea

Informally speaking, k -dimensional De Bruijn graph is a graph with vertices representing all possible strings of length k (called k -mers). In this graph two vertices u and v are connected by a oriented edge from u to v if exists $(k + 1)$ -mer w such that u is the prefix of w and v is the suffix of w .

Given a number k and a string S we can build De Bruijn G graph from the string as follows: for each $(k + 1)$ -mer found in S we add corresponding edge to k -dimensional De Bruijn graph and label it with position of first symbol belonging to the $(k + 1)$ -mer (multiedges with different labels are allowed). In this graph we allow only paths that have consecutive labels on edges, i.e. a sequence of edges $P = (e_1, e_2, \dots, e_n)$ is a path iff the starting vertex of e_{i+1} is the ending vertex of e_i and $P(e_{i+1}) = P(e_i) + 1$. It is easy to see that with such restriction every substring of S corresponds to a valid path in G . Example of such De Bruijn graph built from the string $S = "acgcattgtactcatt"$ and $k = 2$ is depicted on Figure 1a.

Note that two copies of substring "catt" form non-branching path consisting of edges with multiplicity 2 in this graph. Single mismatch in substrings "acGca" and "acTca"

form so-called "bulge", unoriented cycle generated by two valid paths with coinciding ends. If we replace one branch of the bulge by another (say, replace "acGca" by "acTca") then we will obtain long non-branching path (Figure 1b).

This heuristic forms basis of our method – conserved regions in different parts of the genome contain conserved basepairs, but such regions are disrupted by indels and mismatches. These artifacts generate bulges in the graph that spoil non-branching paths. We remove bulges having size less than some predefined constant and thus obtain non-branching paths corresponding to the conserved regions. The process of removing from the graph bulges is called *simplification*.

Conserved regions can be located on opposite strands of DNA. To handle this, we use *colored* De Bruijn graphs [3]. Given string S , for each $(k + 1)$ -mer found in S we add corresponding edge to the graph and color it *blue*, for each $(k + 1)$ -mer found in reverse-complementary counterpart of S we add corresponding edge to the graph and color it *red*. In this graph, coinciding paths with different colors represent synteny blocks located on opposite strands of DNA.

Our algorithm depends on two parameters – k and minimum allowed size of a bulge. It is reasonable to use as high k as possible to keep graph structure simple and avoid connecting regions that are actually not homologous. But even highly conserved regions can lack shared k -mers for $k > 30$. To overcome this, we first find a set of local alignments in the genome (by BLAST), then substitute one aligned subsequence by another for each found significant local alignment. So the pipeline is following:

- 1) Find all alignments in the genome(s) and perform substitution (see above)
- 2) Build De Bruijn graph
- 3) Simplify the graph
- 4) Output synteny blocks as non-branching path in the graph

4. Detailed description

In this section we will present our algorithm in more formal fashion. Suppose that we are given a string $S = (s_0, s_1, \dots, s_{n-1})$ over alphabet $\Sigma = \{A, C, G, T\}$, two numbers k and δ . We denote by $S(i, j)$ the (i, j) substring of S , $S(i, j) = (s_i, s_{i+1}, \dots, s_j)$. Let's denote by \bar{S} string that is the reverse complementary of S .

Colored De Bruijn graph is graph $G_k = (V, E)$ where $V = \Sigma^k$. We define two functions:

- 1) $C : E \rightarrow \{Blue, Red\}$ and call it *color* of an edge
- 2) $P : E \rightarrow \mathbb{N}$ and call it *position* of an edge

For each $i \in \{0, 1, \dots, n - k - 1\}$ we add two oriented edges to the graph:

- 1) $e = (S(i, i + k - 1), S(i + 1, i + k))$, $C(e) = Blue$, $P(e) = i$
- 2) $\bar{e} = (\bar{S}(i, i + k - 1), \bar{S}(i + 1, i + k))$, $C(\bar{e}) = Red$, $P(\bar{e}) = n - i - 1$

A path in G_k is a sequence of edges $P = (e_1, e_2, \dots, e_n)$ iff the starting vertex of e_{i+1} is the ending vertex of e_i , $|P(e_{i+1}) - P(e_i)| = 1$ and $C(e_{i+1}) = C(e_i)$. Let's denote by $Start(P)$ first vertex of the path P and by $End(P)$ the last vertex of P .

A pair of paths $B = \{b_1, b_2\}$ is called a *bulge*, iff following holds:

- 1) $|b_1| < \delta \vee |b_2| < \delta$
- 2) $Start(b_1) = Start(b_2) \vee End(b_1) = End(b_2)$
- 3) b_1 and b_2 have no common vertices except $Start(b_1)$ and $End(b_2)$

5. Experimental results

Results

6. Conclusion

Conclusion

References

- [1] Son K. Pham, Pavel A. Pevzner. DRIMM-Synten: decomposing genomes into evolutionary conserved segments. *Bioinformatics* (2010) 26 (20): 2509-2516.
- [2] Pavel A. Pevzner, Haixu Tang, Michael S. Waterman. An Eulerian path approach to DNA fragment assembly. *Proc. Natl. Acad. Sci. USA*. 2001 Aug 14; 98(17): 9748-53.
- [3] Zamin Iqbal, Mario Caccamo, Isaac Turner, Paul Flicek, McVean. De novo assembly and genotyping of variants using colored de Bruijn graphs. *Nat Genet*. 2012 Jan 8;44(2):226-32. doi: 10.1038/ng.1028.
- [4] Manolis Kellis, Bruce W. Birren, Eric S. Lander. Proof and evolutionary analysis of ancient genome duplication in the yeast *Saccharomyces cerevisiae*. *Nature* 2004 Apr 8;428 (6983): 617-24.