

more complex genome information but also can represent the distinct gene and assemble the genome of the individual or the population in a highly accurate performance. (?) Specifically, rather than using the hash table for the de Bruijn graph topology, Muggli et al. [6] introduced the efficient VARI representation to the colored de Bruijn graph, which collaborates with BOSS representation by saving the memory space though sacrifice the speed of operations. In this proposal, we want to introduce a creative way of reusing the colored de Bruijn graph to decrease the requirement for genome information storage and improvement the accuracy of the colored de Bruijn graph.

2 Related work

Rainbowfish is a succinct representation of the color information, and uses rank and select operations to lookup the color class corresponding to k-mers in the de Bruijn graph. Ranka(i) returns the number of occurrences of symbol a on the range [0, i], whereas selecta(i) returns the position of the ith occurrence of symbol a. And it is also a representation of colored de Bruijn graph. It is introduced by Iqbal et al. (Nature Genetics, 2012 as a variant of the de Bruijn graph, where each edge is related to some set of colors. And it is aimed at “detecting and genotyping simple and complex genetic variants in an individual or population”. The memory usage of the colored de Bruijn graph representation adopted in Cortex precludes this approach from being adopted when the underlying genomes and color sets become too large. While, Because of the bottleneck of space and time, it is a big problem for maintenance and navigation of the De Bruijn Graph in genome assembly. In recent years, a lot of research have been done to solve this problem [2]. For example, Holley et. al. proposed the Bloom Filter Trie, which also is a succinct data structure for the colored de Bruijn graph. SplitMEM is a related algorithm to create a colored de Bruijn graph from a set of suffix trees representing the other genomes.

3 Approach

In this section, we have a short introduction about de Bruijn graph and BOSS presentation first, next, we describe the colored de Bruijn graph and VARI representation, then introduce a succinct representation of the color information structure called rainbowfish. Finally proposing a possible approach for Reusing Colors in Rainbowfish so that storage of the colors of the graph can be further reduced by reducing the number of colors used.

3.1 BOSS representation

Within the last two decades, assembling a genome from enormous amount of reads from various DNA sequences has been one of the most challenging and important computational problems in molecular biology. Though the problem is proved to be NP-hard [3], many algorithms have been proposed for the problem. Most old-time algorithms (especially for the long Sanger reads) construct “overlap graph” after finding the overlapping pairs of reads. But this strategy is difficult to apply when using a huge amount of data from more recent epoch-making next-generation sequencer(NGSs). Due to the vast amount of genome data sequencing by NGS machine, the comparing step for all the pairs of reads becomes extremely difficult. Also, Most NGSs cannot read long DNA fragments, and their read lengths are not long enough to detect overlaps between reads. That’s why “De Bruijn Graph” being introduced. A de Bruijn graph is a graph where each edge represents a k-mer (a substring of length k) that exists in the reads, and one vertex that connect this edge labelled as prefix (k-1)-mer (a substring of length k-1), the other vertex that connect this edge labelled as suffix (k-1)-mer (a substring of length k-1). Glue all vertices that have the same label. Intuitively, de Bruijn graph can be constructed more efficient that

overlap graph, but the overlap phase still a bottleneck for de Bruijn graph due the large amount requirement of memory usage. Then, the introduce of BOSS representation of de Bruijn graph only require $m(2 + \log \sigma)$ bit to store where m is the number of edges in the de Bruijn graph and σ is the alphabet size (i.e. $\sigma = 4$ in the case of DNA). The size of this presentation is not affected by the value of k and is very small. [4]

3.2 VARI representation

Around 2012, with the introduction of colored de Bruijn graph [5], a variant structure of de Bruijn graph, which is used for detecting and genotyping simple and complex genetic variants in an individual or population. A succinct representation of colored de Bruijn Graph was introduced. [6] And this data structure for colored de Bruijn graph is based on BOSS representation for de Bruijn graph. It reduces a large number of memory usage to store the coloured de Bruijn graph, just like BOSS save a lot of memory usage to store de Bruijn graph.

First, let’s have a quick introduction about what is a colored de Bruijn graph. The basically structure of colored de Bruijn graph is the same as the classical de Bruijn graph, however, each vertex ((k-1) mers) and edge (k-mers) is associated a list of colors corresponding the sample that vertex or label exists. For example, there is a set of n samples, and each sample corresponding to a de Bruijn graph. Given a set C of n colors $c_1, c_2, c_3, \dots, c_n$ where c_i corresponds to sample i where all its relevant k-mers and (k-1) mers are colored with c_i . A colored de Bruijn graph may consist many graphs and its corresponding color, except overlapping nodes.

3.3 Rainbowfish

In another succinct representation of coloured de Bruijn graph, which can be called Rainbowfish, it also adopts the BOSS representation of the de Bruijn graph topology, and compress the stored color information of a graph even further to achieve a more space usage. [7]

We basically use the same structure as Rainbowfish, but we came up with a new strategy to reuse colors in the construction of the graph.

Here, we briefly introduce the definition of rank and select operation that are used to look up the color class corresponding to k-mers in the colored de Bruijn graph.

Rank and select are operations that are commonly used for navigating within succinct data structures. For a bit vector $B[0, 1, \dots, n-1]$, rank(i) returns the number of 1s in the prefix $B[0, 1, \dots, i]$ of B. select(j) returns the position of jth 1, which is, the smallest index i such that rank(i) = j. For example, for the 12-bit vector $B[0, \dots, 11] = 100101001010$, RANK(5) = 3, because there are three bits set to one in the 6-bit prefix $B[0, \dots, 5]$ of B, and SELECT(4) = 8, because B[8] is the fourth 1 in the bit vector.

Now, let’s have short introduction for rainbowfish:

There are two fundamental observations for Rainbowfish’s compact representation of color information:

First, it is often the case that many of the k-mers in a colored de Bruijn graph share the same set of colors. So, let’s define an equivalence relation over the set of k-mers in de Bruijn graph, and denote Col(.) as the function that maps each k-mer to its corresponding set of colors. If and only if $Col(k_1) = Col(k_2)$, we say $k_1 \sim k_2$ (k_1 and k_2 “are color-equivalent”), and we refer the set of colors shared by k_1 and k_2 color class.

Second, the color classes distribution frequency is not uniform. It would be useful if we use a small number of bits represent color classes occurring frequently, and use larger number of bits represent color classes occurring less frequent. (The approach is similar to Huffman code)

In Rainbowfish, the color class representation has three components which are derived from color matrix. Color matrix is used to store <k-mers, color-set> pairs from colored de Bruijn graph. It stores the mappings between labels and color classes in an equivalence class table(ECT). As labels are assigned sequentially, this is simply an array of bit vectors

This is a footnote



Text Text Text Text Text Text Text Text Text Text Text Text Text
Text Text Text Text Text Text Text. Figure 2 shows that the above method
Text Text Text Text

Text Text Text Text Text Text Text. ? might want to know about text
text text text

[1]Idury, R., Waterman, M.: A new algorithm for DNA sequence assembly. *Journal of Computational Biology* 2, 291–306 (1995).

- 138016.