



HOKKAIDO UNIVERSITY

Title	SLIDESORT : Developing an exact method to find similar pairs with small edit-distance
Author(s)	Shimizu, Kana; Tsuda, Koji
Description	ERATO 湊離散構造処理系プロジェクトシンポジウム (第1回) : 第9回情報科学技術フォーラム(FIT2010) イベント企画セッション. 2010年9月8日 (水). 九州大学伊都キャンパス.
Relation	2010年度科学技術振興機構ERATO湊離散構造処理系プロジェクト講究録. p.376.
Issue Date	2011-06
Doc URL	https://hdl.handle.net/2115/48368
Type	conference presentation
File Information	04.FIT_tsuda.pdf





SLIDESORT : Developing an exact method to find similar pairs with small edit-distance

Kana Shimizu, Koji Tsuda

Computational Biology Research Center, National Institute of Advanced Science and Technology, JAPAN

Contact: shimizu-kana@aist.go.jp

ABSTRACT: In this study, we propose an exact method that finds all similar pairs from a string pool in terms of edit distance without any duplication on small memory in $O(N)$. Using an efficient pattern growth algorithm, our method discovers chains of common k -mers to narrow down the search. Compared to existing methods based on single k -mer, our method is more effective in reducing the number of edit distance calculations. In large short read datasets, our method was 20-3000 times faster than the state-of-the-art method, BWA.

Motivation

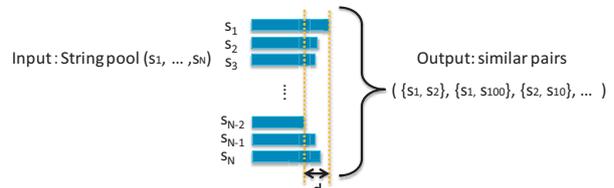
Due to the dramatic improvement of DNA sequencing, it required to evaluate sequence similarities among a huge amount of fragment sequences such as short reads. Finding similar pairs from a string pool is an important first step for many biological sequence analyses. In this work, we would like to address the problem of enumerating all neighbor pair in a large string pool in terms of edit distance. This problem is conventionally called "all pairs similarity search", that will help ...

- ▶ Shot reads analyses
 - preprocess of de novo assembly
 - Clustering reads. (e.g. online tree construction, greedy clustering etc ...)
- ▶ genome analyses
 - Repeat analyses
 - Sequence pattern analyses etc...

Problem setting

Given n strings of similar length (maximal length - minimal length $\leq d$), s_1, \dots, s_N , the task is to find all pairs whose edit distance is at most d ,

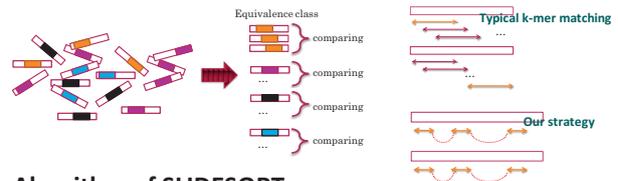
$$E = \{(i, j) \mid \text{EditDis}(s_i, s_j) \leq d, i < j\}$$



Method

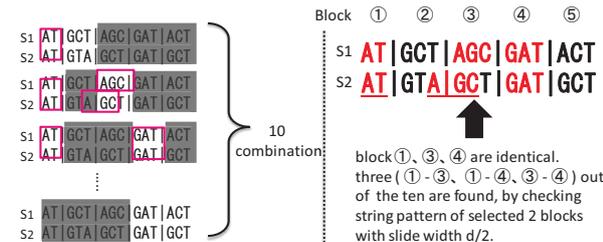
Strategy of SLIDESORT

Like other method based on k -mer matching, SLIDESORT finds common substrings and verify the matches. The key idea of SLIDESORT is to find longer common substrings to narrow down search space. We developed an efficient pattern growth algorithm inspired by multiple sorting algorithm (*Uno 2008 PAKDD*).



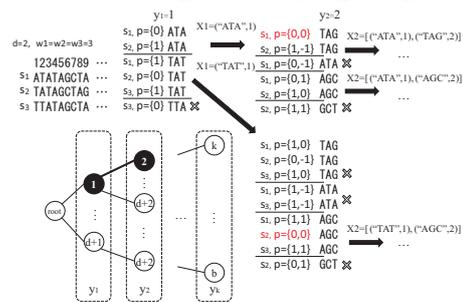
Theory of SLIDESORT

How SLIDESORT can find longer common substrings? Given two sequences that are divided into $k+d$ blocks, if edit distance of the two sequences is at most d , there exist at least k blocks that are the same string patterns between the two within slide width $d/2$. Based on this theory, SLIDESORT regards concatenation of k blocks as a common substring.

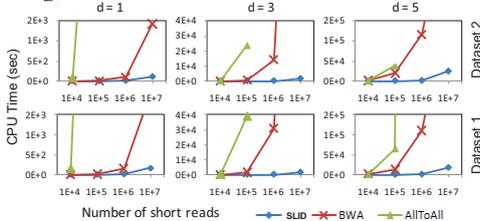


Algorithm of SLIDESORT

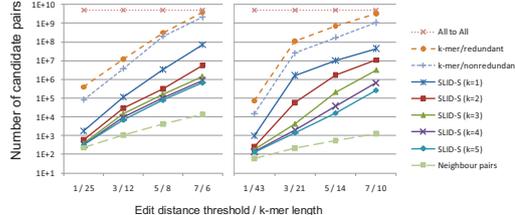
A pattern corresponds to a sequence of substrings. The space of all patterns is organized as a tree and systematically traversed. SLIDESORT uses radix sort to find equivalent strings in pattern growth.



Experiments



Evaluation on two types of short read datasets. The graphs mainly compares performance difference in size of the dataset.



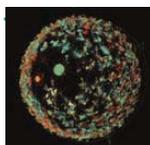
Comparison of number of candidate pairs. The evaluations were done on 100,000 short reads. The proposed method was examined with $k=1, \dots, 5$. 'Neighbour pairs' represent the actual number of neighbour pairs in data. 'k-mer/nonredundant' and 'k-mer/redundant' represent two variants of the single seed method.

Dataset from NCBI SRA
ERR1081 : Read length 51
SRRO20262 : Read length 87
Computation environment
Implemented by C++
Linux PC with Intel Xeon X5570

SLIDESORT is available from www.cbrc.jp/~shimizu/slidesort/

Application

One of an application of SLIDE-SORT is constructing a tree. The right graph is a minimum spanning forest of 112,995 short reads that was generated by processing a pair of sequences (an edge) one-by-one. Most of the case, computation time of each process is smaller than interval time to find next pair.



d	SLIDE-S	BWA	AllToAll
1	0.53	2.34	3.25
3	2.71	2.37	562.63
5	20.64	2.19	19697.67
7	131.88	2.51	10161.15

This work is supported by Grant-in-Aid for Young Scientists (22700319) of Japan Society for the Promotion of Science.