

Search for Maximal Unique Matches Multi-core architectures

Julio Cesar Garcia Vizcaino, Antonio Espinosa, Juan Carlos Moure,
Porfidio Hernandez

Universitat Autònoma de Barcelona (UAB), Barcelona, Spain

{jcgarcia,aespinosa}@caos.uab.cat
{juancarlos.moure,porfidio.hernandez}@uab.es

Abstract. Maximal Unique Matches (MUMs) are common substrings that match a Reference and a Query genome. They are exact, unique and maximal. The search for MUMs in large genomes is a heavy and repetitive task, so there is a fair chance of parallelize and execute this search on multi-core architectures. This research evaluates a parallelization of the search for MUMs in genomic sequences within multi-core architectures. The Reference genome is indexed by using a Suffix Tree or Extended Suffix Array in main memory and then a parallelized algorithm is used for MUMs between Reference genome and Query genome. Query genome is read by several threads in chunks of fixed size. This approach is compared on MUMmer, a genome alignment tool, which is able to find Maximal Unique Matches (MUMs). Results show a reduction in search time and less usage of memory.

Keywords: Indexed Search, Bioinformatics, Maximal Unique Matches, Multi-core architectures, Parallelization

1 Introduction

Modern sequencing and computational technologies and advances in bioinformatics has made whole genome sequencing possible. One result of this is the fast alignment of whole genomes. Dynamic programming is used for aligning two large genomes, hundreds of Mbp. One very successful approach for performing a whole genome alignment is based on identifying “maximal unique matches”. This heuristic is based on the assumption that one expects to find many matches occurring in two similar genomes. Maximal unique matches (MUMs) are a natural and surely part of a good alignment of the two sequences and so the search for MUMs is a good starting point for aligning the two sequences.

The problem of searching maximal unique matches for a minimum between a Reference genome and a Query genome has been identified in many applications, one of them is MUMmer [3]. MUMmer’s algorithm searches for Maximal Unique Matches (MUMs), although with a large main memory to store the Reference genome and a null use of multi-processor architectures.

We designed and tested a data-level parallelism to use in multi-processor architectures with two different data structures: Suffix Tree and Enhanced Suffix Array.

The rest of this paper is organized as follows. Section 2 discusses previous work. Section 3 defines some definitions used in this paper. Section 4 describes the implementation of parallelization of search for MUMs. Section 5 shows the results of our parallelization. Section 6 concludes the paper.

2 Related work

Search for Maximal Unique Matches to do Whole Genome Alignment is proposed in [2]. There have been some previous work in the parallelization of matches in genomic data, like [11, 9, 7], however these works search for fixed patterns and read alignment. On the other hand, there have been some experiments in parallelization of Whole Genome Alignment like [8]. The problem of searching for MUMs with a full-text index data structure is also not covered deeply, there is only one approach in [4] but without a reference code to check their implementation and it is more focused with Cache-aware hybrid architectures and Suffix Array. There are other implementations [11, 6, 10] but they search for Maximal Exact Matches with threads instead of Maximal Unique Matches. We propose the use of OpenMP to search for MUMs in multi-core architectures with two full-text index data structures: Suffix Tree and Enhanced Suffix Array.

3 Preliminaries

Let’s assume the Reference genome $R[0, \dots, n - 1]$ of size $|R|$ over an alphabet $\Sigma = \$, A, C, G, T$ which has a sentinel character $\$$ that does not occur nowhere else in the Reference genome and is lexicographically smaller than the characters that occur in the alphabet. The suffixes of the Reference genome are zero indexed by their position in the original Reference by a full-text index data structure like a Suffix Tree or an Enhanced Suffix Array.

Definition 2. A *suffix Tree*, ST , for an n -character string R is a directed graph

Definition 3. *An enhanced suffix array consists of four arrays (common prefix (LCP), child and inverse suffix arrays) that have information saved in a suffix tree. [1].*

We also define a suffix link, which is an important trick to achieve complexity to search for MUMs in both Suffix Tree and Enhanced Suffix

Definition 4. *A suffix link is a pointer from string \overline{aw} to substring*

For reasons of space we recommend the reader to check [1] for full details of Enhanced Suffix Array. A search for a MUM between a Reference genome, $|R| = n$, and a Query genome, $|Q| = m$, can be done in a Suffix Tree in $O(m)$ steps and in an Enhanced Suffix Array in $O(m)$ steps. Once we have the resources used to index a Reference genome, we need to answer the question: Where are the MUMs of R and Q of some minimum length L ? Several algorithms used to answer this question using a Suffix Tree, see [1] and a Enhanced Suffix Array, see Algorithm 2.

Algorithm 1: Search for MUMs in a Suffix Tree.

```

input : R, Q, L
output: List of MUMs of length  $\geq L$ , with start position in R and Q
begin
    ST  $\leftarrow$  buildST(R);
    foreach position  $i \in Q$  do
        length  $\leftarrow$  TraverseSuffixTree(Q[i],ST);
        if isLeafNode(Q[i]) then
            /* Leaf saves the position of a suffix in ST.
            ;
            if R[leaf - 1]  $\neq$  Q[i - 1] and length  $\geq L$  then
                MUMcands  $\leftarrow$  saveMUMcand(R[leaf,i],length);
        /* Get unique MUMs from list of MUM-candidates.
        ;
    MUMs  $\leftarrow$  cleanMUMcand(MUMcands)

```

4 Implementation

The search for MUMs between a Reference and Query genome has

Algorithm 2: Search for MUMs in an Enhanced Suffix Array.

```

input : R, Q, L
output: List of MUMs of length  $\geq L$ , with start position in R and
begin
    ESA  $\leftarrow$  buildESA;
    while Q  $[i] < \text{end}$  do
        traverseESA(Q, Q  $[i]$ , interval, Q  $[i].\text{length}()$ );
        if interval.depth  $\leq 1$  then
            interval.depth = 0; interval.start = 0; interval.end = n-1;
            continue;
        if interval.size() == 1 and interval.depth  $\geq L$  then
            if leftMaximal(Q, Q  $[i]$ , SA  $[\text{interval.start}]$ ) then
                saveMUMcand(SA  $[\text{interval.start}]$ , i, interval.depth);
        repeat
            interval.depth = interval.depth-1;
            interval.start = ISA [SA  $[\text{interval.start}] + 1$ ];
            interval.end = ISA [SA  $[\text{interval.end}] + 1$ ]; i ++; if interval
            0 or suffixlink(inter) == false then
                interval.depth = 0; interval.start = 0; interval.end =
        until interval.depth > 0 and interval.size();
    MUMs  $\leftarrow$  cleanMUMcand(MUMcands)

```

MUM candidates which are not unique in both Reference and Query genome. Our approach, see Figure 1, is divided in three phases: Creation of the Enhanced Suffix Array; Splitting query genome data (chunks) into C chunks, the number of available cores using 1 thread per core and parallelizing the search for MUMs for every chunk, then every thread has to find MUM candidates; and Get MUMs from list of MUM candidates. A data-level parallelism to search for MUMs requires to split the Query genome in chunks and stream each chunk against the full-text index data structure. This search is independent of each other chunk. The final stage to find MUMs must be performed with a full Query genome and with chunks of Reference genome. This is more obvious too. This approach is suitable to be performed in different architectures because chunks can be assigned to each core and search for MUMs in the full-text index data structure which is stored in main memory. A good choice of parallelization involves a data partition technique.

Memory usage is improved, within this approach, with multiple threads against the full-text index data structure, instead of one read at a time with

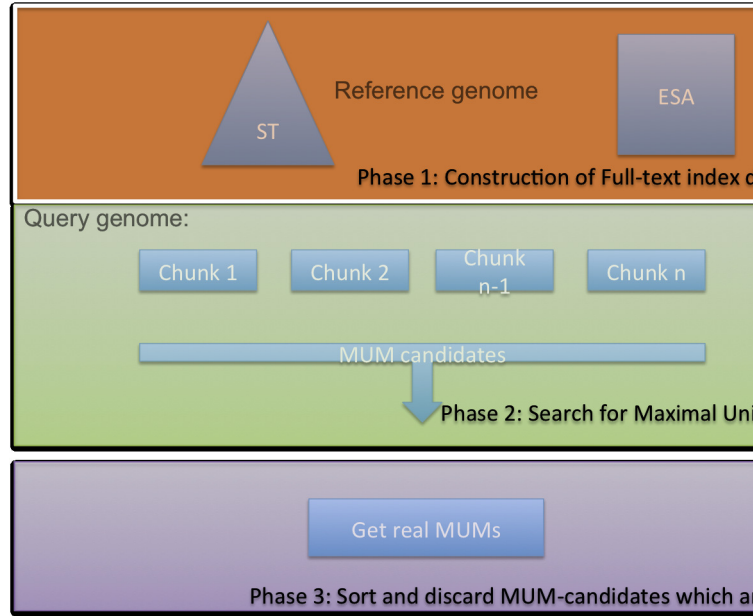


Fig. 1. Approach to search for MUMs in multi-core architectures with Enhanced Suffix Array.

The chunk size was computed with Query genome length divided by the number of available threads. Chunk size is a performance factor because we need a balanced workload among threads.

Phase 1 is the construction of the data structure: Suffix Tree and Suffix Array; this phase is executed in a serial way. Phase 2 requires the parallel execution of the Algorithm (1 or 2) to search for MUMs, the first step of which is the split of the Query genome into as many chunks as available threads. The split defines the start and end position in the Query genome, which is stored in memory and shared by all threads. The chunk size is fixed according to the number of threads. The parallel search for MUMs gets a list of MUM-candidates, which are ordered in each thread by their position in the Query genome. The merge process is performed after all threads have finished their execution and it merges the lists of MUM-candidates and it is ordered by position in the Query genome. This step is needed because the merge process may have unordered lists; the final result is a list of real MUMs.

5 Results

To verify that our approach, can have a better performance to search for MUMs in a Reference and Query genome, in Table 1, we verify that our approach and our approach have the same MUMs. Tests were carried out in a computing node: 2 Processor Intel(R) Xeon(R) E5645 @ 2.4GHz of 6 cores, 32KB L1 cache, 256KB L2 and 12MB L3 shared cache per socket. GCC 4.7.0 with OpenMP support + Linux. The main goal of these experiments

Table 1. List of Genomes used in experiments.

	1	2	3
Reference	4.64 [Mbp]	169 [Mbp]	1031 [Mbp]
Query	5.5 [Mbp]	167 [Mbp]	1357 [Mbp]

check the performance to search for MUMs in multi-core architectures using OpenMP (threads). This performance involves the usage of Memory (execution time to search for MUMs).

Times were collected during the execution of the experiments: construction time of data structure and execution time to search for MUMs with the use of `omp_get_wtime` of OpenMP.

5.1 Construction of Full-text index data structure

This paper compares the performance of two data structures to search for MUMs in multi-core architectures, two features for these structures are execution time and memory footprint. Full-text indexes allow fast access to any length, but they have a great memory and construction cost, which is affected by the type and implementation of the full-text index data structure used. We measure the construction time for the set of genomes in Table 2. Table 2 shows that an ESA has a lower construction time for all the Reference genomes used. Since a full-text index data structure has to be build every time a search for MUMs is performed, a reduction in the construction phase in Full-text index data structure allow us to improve the overall execution time. A future improvement would be to build the full-text index data structure in parallel or load from memory.

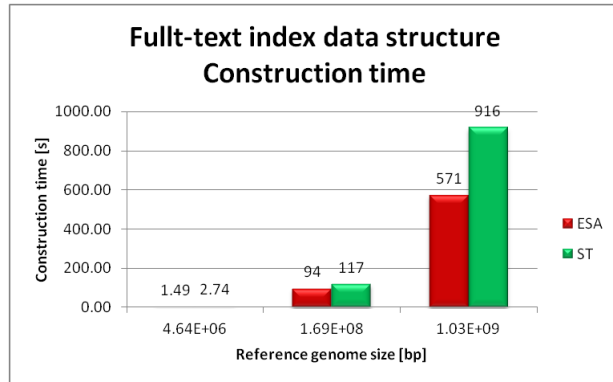


Fig. 2. Construction of Suffix Tree and Enhanced Suffix Array for different reference genome sizes. Our approach stores the whole Reference genome indexed by a Suffix Tree (ST) or an Enhanced Suffix Array (ESA) in main memory.

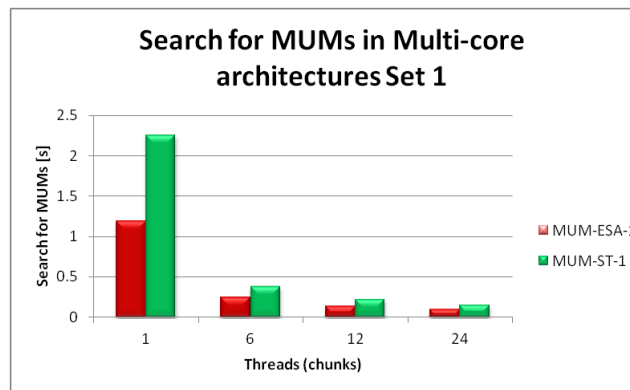


Fig. 3. Search for MUMs time for Ecolik12 and EcoliO157H7 in multi-core architectures with ST and ESA.

EcoliK12 vs. EcoliO157H7 The set 1 in Table 1 was used to search for MUMs in a ST and ESA. In Figure 3, we show that ESA has a better search time than ST. Figure 3. However, the space consumption shows that ESA has a better memory footprint, a reduction of 44% of memory space compare to Suffix Tree (ST) (1.03 GB vs. 1.8 MB). This first experiment shows that it is possible to search for MUMs with two different full-text index data structures with the reduction in space consumption. In the next section, we will compare the search time for MUMs with the space consumption for ESA and ST.

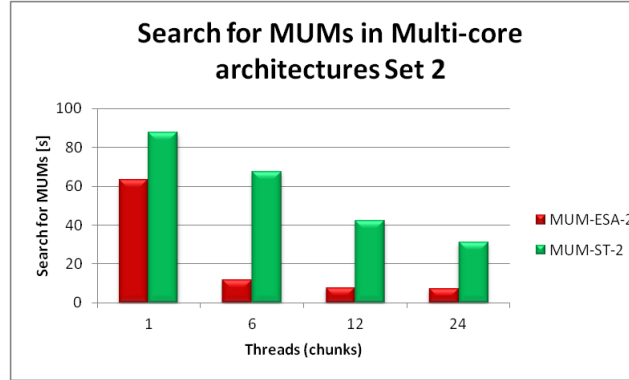


Fig. 4. Search for MUMs time for *D.melanogaster* and *D.pseudoobscura* architectures with ST and ESA.

need to reduce the search time. On the other hand, the space cost is almost constant in both data structures. However, again ESA has a lower cost of memory. The reduction in memory footprint for ESA is near 50% compared to the ST.

Chicken vs. Zebrafish The set 3 in Table 1 have genomes of 1.2 Gbp [Gbp] and we show in Figure 5 that an execution with one thread achieves a better performance in search time in ESA over ST. In order to understand

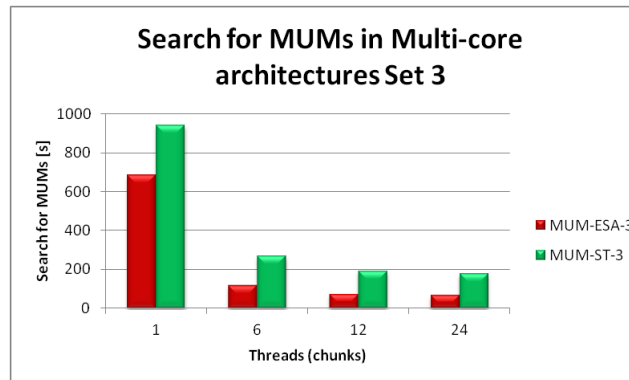


Fig. 5. Search for MUMs time for Chicken and Zebrafish architectures with ST and ESA.

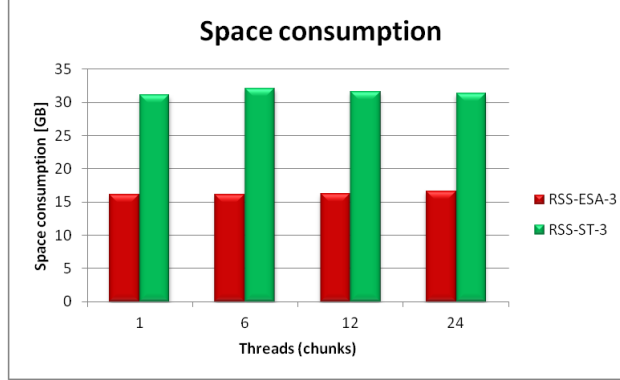


Fig. 6. RAM consumption to search for MUMs for Chicken and Zebrafish architectures with ST and ESA.

6 Conclusions

This paper presents an evaluation of performance to search for MUMs between a Reference genome and Query genome in multi-core architectures using either a Suffix Tree and an Enhanced Suffix Array. The search for MUMs is performed in a Suffix Tree or an Enhanced Suffix Array and the results can help in the process of Whole Genome Alignment. The results show that the MUM-problem can be solved with a ST or ESA and with the use of multi-core architecture we reduce the execution time. There is a reduction in execution time with an Enhanced Suffix Array over a Suffix Tree. We have evaluated a new data structure which provides the same functionality of a Suffix Tree for MUMs, with less use of memory in multi-core architectures.

Improvements involve a better use of CPU when we have more than one thread per core. From the results obtained we conclude that an Enhanced Suffix Array is suitable to solve the MUM-problem. To improve the performance of CPU usage of our approach, the following proposals may be used: Interval management: this may work because we know in advance what intervals to look for in the Reference and Query genome; Data cache management: we may need to store the intervals of ESA in some cache-level of processor. Since we need to get the results of the search by checking the first character of the suffix within the interval, we can use some SSE instructions to check more than one suffix at the same time.

References

1. M. I. Abouelhoda, S. Kurtz, and E. Ohlebusch. Replacing suffix enhanced suffix arrays. *Journal of Discrete Algorithms*, 2(1):53–86, May 2005.
2. A. L. Delcher, S. Kasif, R. D. Fleischmann, J. Peterson, O. White, and S. Salzberg. Alignment of whole genomes MUMMER. *Nucl Acids Res*, 27(12):2322–2326, 1999.
3. A. L. Delcher, S. L. Salzberg, and A. M. Phillippy. Using MUMmer to find similar regions in large sequence sets. *Current protocols in bioinformatics*, edited by *Andreas D Baxevanis et al*, Chapter 10(1934-340X (Electronic version) - Journal Article SB - IM):Unit 10.3, 2003.
4. G. Encarnac and N. Roma. Advantages and GPU Implementation of a Performance Indexed DNA Search based on Suffix Arrays. *Arch Biochem Biophys*, 49:49–55, 2011.
5. D. Gusfield. *Algorithms on strings, trees, and sequences: computational combinatorial biology*. Cambridge University Press, 1997.
6. Z. Khan, J. S. Bloom, L. Kruglyak, and M. Singh. A practical algorithm for finding maximal exact matches in large sequence datasets using sparse suffix arrays. *Bioinformatics*, 25(13):1609–1616, July 2009.
7. C. S. Kouzinopoulos, P. D. Michailidis, and K. G. Margaritis. Parallelization of multiple pattern matching algorithms for biological sequences: performance results. *Systems and Computational Biology - Bioinformatics and Computational Modeling*, 2005.
8. X. Meng and V. Chaudhary. Exploiting multi-level parallelism for high speed using general purpose processors. *October*, 2005.
9. H. Mongelli. Efficient Two-Dimensional Parallel Pattern Matching Algorithms. *IEEE Trans Pattern Anal Mach Intell*, 19(10):1100–1112, 1997.
10. E. Ohlebusch, S. Gog, and A. Kügel. Computing matching statistics for maximal exact matches on compressed full-text indexes. In *SPIRE*, pages 3–14, 2009.
11. M. Oğuzhan Külekci, W.-K. Hon, R. Shah, J. Scott Vitter, and B. W. Berman. A parallel sparse index for genomic read alignment. *BMC genomics*, 12:1, Jan. 2011.
12. M. Vyverman, B. De Baets, V. Fack, and P. Dawyndt. Finding maximal exact matches using enhanced sparse suffix arrays. *Bioinformatics (Oxford, England)*, pages 2–4, Feb. 2013.