

International Conference on Computational Science, ICCS 2012

Search of maximal unique matches in multicore architectures

Abstract

Maximal Unique Matches are common substrings that match a reference and a query sequence. They are exact, unique and maximal; that is, they cannot be extended in left or right direction without incurring a mismatch. The computation of MUMs in large sequences is a heavy and repetitive task, so there is a fair chance of parallelize and execute this search in multicore architectures. This research resembles an approach to find MUMs in genomic sequences in parallel way. The reference genome is indexed by using a suffix tree in main memory and then parallelized algorithm finds the MUMs of query genome which is readed by several threads. This approach is based on MUMmer, a genome alignment tool, which is able to find Maximal Unique Matches (MUMs).

Keywords:

1. Introduction

The problem of searching maximal unique matching for a minimum length between a reference string and a query string has been identified in several applications, one of them is MUMmer. Altough MUMmer's algorithm can perform searches of maximal unique matches (MUMs) the use of resources:

- High use of main memory to store the reference string.
- A null use of multicore architectures.

If the length of reference and query are very huge, the amount of operations to perform in the search of MUMs increases, see table 1. The use of parallelism could help reduce the execution time for the search of maximal unique

Data structure	L [bp]	Search operations	Search [s]	Memory usage [MB]
Suffix tree	20	9.87×10^{18}	169189,4	48665,12

Table 1: Search of Maximal Unique Matches between a reference sequence (2960,21Mbp) and query sequence (2716,96Mbp)

matches. One approach of parallelism is to take advantage of multicore architectures nowadays. This problem has a time complexity of $O(m + k)$ where m is the length of the query sequence and k is the number of maximal unique matches of some minimum length. This problem is a very high intensive computing task, for every substring in the query sequence the search for a maximum unique match has to be performed.

2. The MUM: an heuristic approach

2.1. Definition of MUM

Although a pair of conserved genes rarely contain the same entire sequence, they share a lot of short common substrings and some of them are indeed unique to this pair of genes. For example the following two sequences, R and Q:

$$\begin{aligned} R &= \underline{\text{ac}} \text{ ga } \underline{\text{ctc}} \text{ a } \underline{\text{gctac}} \text{ t } \underline{\text{ggtcagctatt}} \text{ acttaccgc\$} \\ Q &= \underline{\text{ac}} \text{ tt } \underline{\text{ctc}} \text{ t } \underline{\text{gctac}} \text{ } \underline{\text{ggtcagctatt}} \text{ c } \underline{\text{acttaccgc\$}} \end{aligned}$$

It is clear that sequences R and Q have many common substrings, they are: ac, ctc, gctac, ggtcagctatt, acttaccgc. Among those five common substrings, ac is the only substring that is not unique. It occurs more than once in both sequences. You can also observe that actually a, c, t, and g are common substrings of R and Q. However, they are not maximal, i.e. they are contained in at least one longer common substrings. We are only interested in those that are of maximal length.

Our aim is to search for all these short common substrings. Given genomes R and Q, we need to find all common substrings which are unique and of maximal length. Each of such common substrings is known as Maximum Unique Match (MUM). For almost every conserved gene pairs, there exist at least one MUM which is unique to them. For example, assuming $d = 3$, sequences R and Q in the previous example has four MUMs: ctc, gctac, ggtcagctatt, acttaccgc. Substring ac is not an MUM because its length is smaller than the value of d and it is not unique to both sequences.

$$\begin{aligned} R &= \underline{\text{ac}} \text{ ga } \underline{\text{ete}} \text{ a } \underline{\text{gctac}} \text{ t } \underline{\text{ggtcagctatt}} \text{ acttaccgc\$} \\ Q &= \underline{\text{ac}} \text{ tt } \underline{\text{ete}} \text{ t } \underline{\text{gctac}} \text{ } \underline{\text{ggtcagctatt}} \text{ c } \underline{\text{acttaccgc\$}} \end{aligned}$$

The concept of MUM is important in whole genome alignment because a significantly long MUM is very likely to be part of the global alignment.

2.1.1. Finding MUMs in a suffix tree

The key idea in this method is to build a suffix tree for genome R, a data structure which allows finding, all distinct subsequences in a given sequence.

A suffix tree is a data structure which indexes a whole string in every suffix. It allows finding all matches between a reference sequence (suffix tree) and a query sequence. It can also check for unique matches in the reference very easily just to check if the match ends in a leaf node and by checking the character immediately preceding the start of this match it is possible to determine whether it is a maximal match. To understand the efficient search of matches in a suffix tree, time complexity is $O(m)$ where m is the length of query sequence, it is used a suffix link. A link points from a source node s in suffix tree to another destination node d in suffix tree if the string label from the root node to node d is equal to the label from root node to s with the first character removed. In Figure 1 is shown the suffix tree, green lines are suffix links. For instance, the string label of node i in Figure 1 is tgt and next suffix $i + 1$ is gt. So that is exactly the tree position corresponding to the next position in the query sequence.

In Figure 1 the match gtc which starts at $i + 1$ in the query sequence and node 7 in the suffix tree is not maximal on the left because the preceding character in both strings is a t.

By construction, the location of a match in the suffix tree represents a substring of the subject sequence which maximally matches a prefix of query suffix. Thus it is only necessary to verify that, the substring of the subject sequence is long enough, that it is unique in the subject sequence and that the match is also left maximal. This is done as follows:

1. Does location represent a substring of length at least minimum match length?
2. Does location correspond to a leaf edge? Then then the string represented by the location is unique in the subject sequence.
3. Is the substring left maximal? This is true if one of the following conditions hold:
 - The suffix of the query currently considered is the first suffix, or
 - The string represented by loc is a prefix of the subject string, or

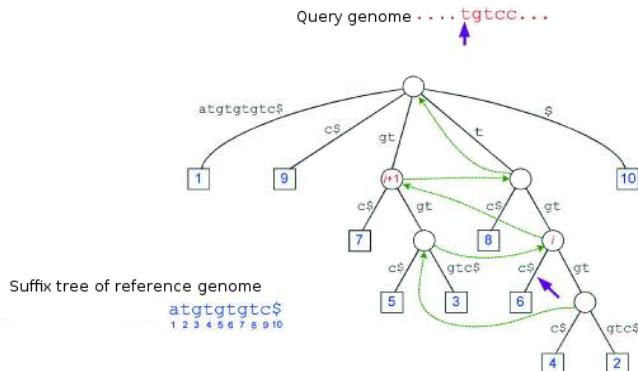


Figure 1: Suffix tree and its use to find matches between a reference and query genome.

- The characters immediately to the left of the matching strings in the subject sequence and the query sequence are different

If all conditions 1-3 are true, then this MUM is stored in a list of MUM-candidates, see Figure 2. It takes the necessary

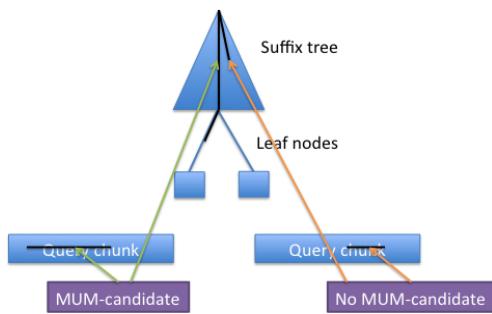


Figure 2: Finding MUMs in a suffix tree.

information about the MUM-candidate:

- Position in reference sequence.
- Position in query genome.
- Length of match.

MUMs can cover 100% of the known conserved gene pairs. Moreover, finding all MUMs can be done in almost linear time.

3. Parallelism technique

General-purpose, commodity CPUs currently have SIMD (Single Instruction Multiple Data) functional units and corresponding SIMD instructions. This kind of CPUs allow to be used in several ways of parallel techniques, such as data-level parallelism.

In addition to using SIMD technology, the availability of multicore architecture makes possible to execute the same task with a different kind of data.

The sequential version of the MUMmer's algorithm trades extra computation for memory and high computation time when executed in a single machine.

Previously the algorithm for sequence alignment was described in detail. Now our own proposal of a parallelization of WGA within multicore architecture is explained. There are two resources to improve in this algorithm:

- Memory usage.
- Running time.

The former was generally improved because it allows being executed in architectures where there is no restriction memory. To improve the performance of the algorithm a data-level parallelism technique is deployed in advance to genome alignment.

Our technique is divided in three phases following:

1. Splitting genome data (chunks) according to the number of available cores using 1 thread per core.
2. Parallel execution of the task of finding MUMs for every chunk where every thread has its own list of MUM-candidate.
3. Get the final list of MUMs from every MUM-candidate list of all threads.

The following Figure 3 shows the process of our data-level parallelism technique. The division of genome data was

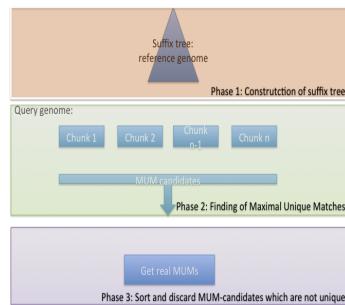


Figure 3: Data-level parallelism in multicore architectures for whole genome alignment.

used using the paradigm of data-level parallelism which consists of a generation of chunks of a query sequence with a fixed size. The main idea behind using a Maximal Unique Match (MUM): it is possible to cover a huge region of a genome when reference and query genomes are very closely related. However to get a MUM, it requires an important feature its uniqueness. Uniqueness can only be found when a whole genome is checked, see Figure 4. If some part of it is only evaluated we could miss the rest of the genome. In other words, after finding MUMs within a chunk it is not possible to determine if the MUM found is or not a "unique" MUM, globally in the query genome, because these MUMs are unique only in the chunk that has been read, the rest of the genome it is not known until all query genome has been read, see Figure 4.

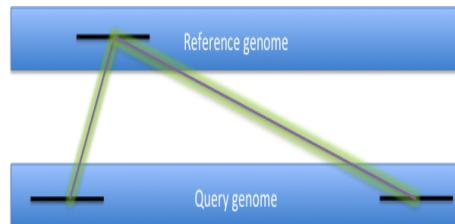


Figure 4: MUM in reference genome but not in query genome.

4. Implementation

Split query genome

As it was previously explained, the approach is to use a fixed size division of query genome in as many chunks as many available cores. To split query genome the algorithm needs to know in advance how many chunks will be used. Then every chunk is computed with two pointers (left and right end) which points out query genome in main memory.

Finding MUMs

The parallelization is carried out with OpenMP. The algorithm to find MUMs is a process which can be executed without any data dependency. However, when we split a query sequence the following problems arise:

- Chunk size is a performance factor.
- Different MUM-candidate in query sequence:
 - Additional MUMs.
 - Lost MUMs.

OpenMP defines the schedule for the loop iterations among the total number of threads. The total number of iterations is the number of chunks created. Every iteration means the whole search of MUMs within a chunk, if the right end of chunk is not the end of the query sequence and there are still nucleotides to match then traversal of suffix tree until it occurs a mismatch or a MUM-candidate is found.

The key factors in this phase are:

- Number of chunks.
- Size of chunks.
- Number of threads.
- OpenMP schedule and its own chunk size.

Get real MUMs

List of MUM-candidates are ordered with quicksort according to position in query sequence. Every thread has found a set of MUM-candidates from previous phase but the all threads don't produce the MUM-candidates in order. That's why a quicksort is required.

After quicksort we need to get the real MUMs. A real MUM is unique in the whole reference and query genome. Those MUM-candidates which are overlapped by bigger MUMs are discarded, see Figure 5.

The final output has the following format:

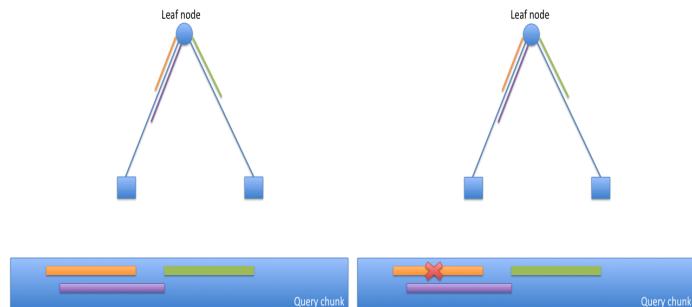


Figure 5: Getting Real MUMs.

```
>Information about the query sequence
Position_in_R Position_in_Q Length_of_MUM
```

5. Experiments and results

To verify that our approach can have a better performance to align a whole genome a set of tests were deployed. These tests were carried out in the following node:

- Hardware:
 - 2 Processor Intel(R) Xeon(R) E5645 @ 2.4GHz of 6 cores each one, 32KB L1 cache, 256KB L2 and 12MB L3 shared cache per socket.
 - RAM: 96 GB
- Software:
 - Linux Kernel 2.6.32-220.el6.x86_64 #1 SMP
 - gcc 4.7.0 with OpenMP support
 - Likwid 2.3.0
- Genomes:
 - Reference: Human chromosome 21 single fasta file
 - Query: Mouse chromosome 16 single fasta file

The main objective of the tests was check the performance of finding MUMs in multicore architectures by using OpenMP (threads). The variables to control were:

1. Number of chunks: fixed for query.
2. Number of threads: from 1 to 12 threads with affinity to 1 thread per core.
3. OpenMP chunk scheduling policy: from 1, 2 and 4 chunks per thread.
4. OpenMP scheduling: static, dynamic and guided.

Moreover, the use of cache was measured with the tool likwid¹, the use of cache included the whole application. The times were collected for the total execution time and the time spent in the parallel region of OpenMP. This parallel region is the algorithm to find MUMs.

One key aspect of the several tests deployed was to assure the execution of every thread in a single core, that is no other thread is competing for the same cache. Affinity is a requirement because we need to ensure the scalability of our proposal.

The construction of the suffix tree for the reference genome of length 35.45[Mbp] takes 20.66[s] to construct and it uses 608[MB] of memory. The minimum length for a MUM was fixed in all tests, 20[bp]. The query sequence of length 95.28[Mbp] was divided in 96 chunks with a fixed size of 992.53[Kbp], the total number of MUMs found were 102844.

Two times were collected during the execution of the experiments: total execution time and execution time to find mums. The experiments were executed 10 times to get accuracy in the measures.

In Figure 6 shows the final time for the application to find MUMs, from this Figure we can conclude that for these sequences the dominant time since 6 threads because the serial part of the application has not been improved; that is phase 1 and phase 3 in Figure 3. However it is shown that phase 2 (finding MUMs) have a near linear speedup because of the forced affinity of one thread per core, in Figure 7 is shown the execution time for this phase with several configurations of OpenMP scheduling policy. The difference between scheduling policy and chunk size is negligible for the sizes of the sequences for performance purposes.

To check the performance of our approach it was measured the use of cache of Level 1 in the multicore architecture. The size of this L1 cache is of 32[KB] for every core. The high cache misses observed in Figure 8 and the replacements in L1 cache in Figure 9 are an obvious conclusion that our application is a memory bound problem. Every cache miss is delay in the execution program, so that a better strategy must be deployed to reduce the execution time, although the theoretical maximum speedup for our application is 4.11 for these test sequences.

In order to reduce the cache misses and indeed improve the performance of our approach the following proposals

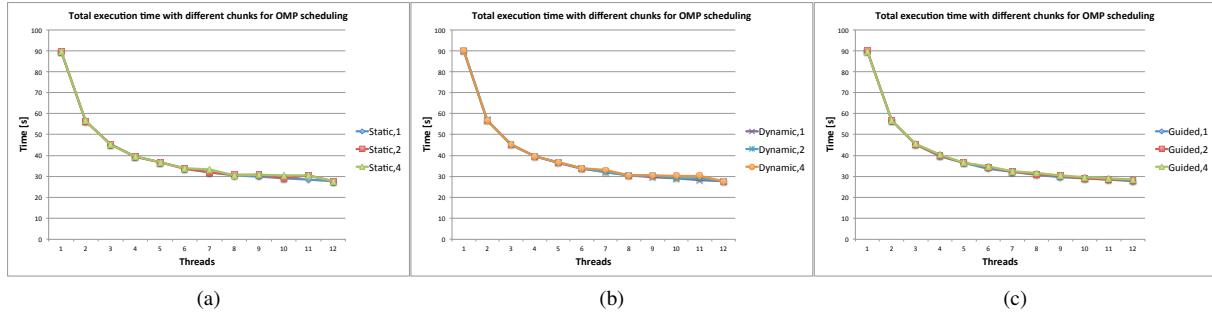


Figure 6: Total execution time for OpenMP scheduling.

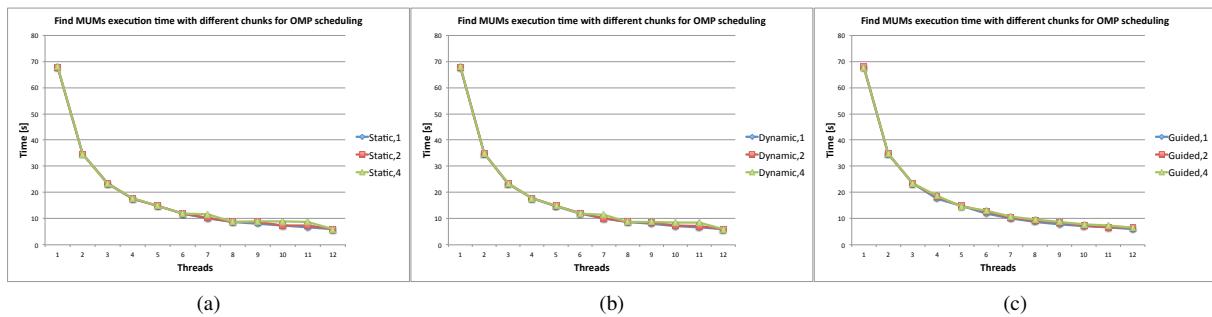


Figure 7: Execution time to find MUMs with OpenMP scheduling.

may be used:

- Prefetching: this may work because we know in advance what information from suffix tree and query sequence we need in every traversal of suffix tree. If we can allow the use of some part of suffix tree be stored and read in cache by several threads which request that subtree.
- Explicit cache management: this is a obvious consequence of previous bullet because we may need to store part of the suffix tree in some cache-level of processor.
- Figure out how much total data will be accessed/required to find next MUM.

References

¹Measure hardware performance counters on Intel and AMD processors.

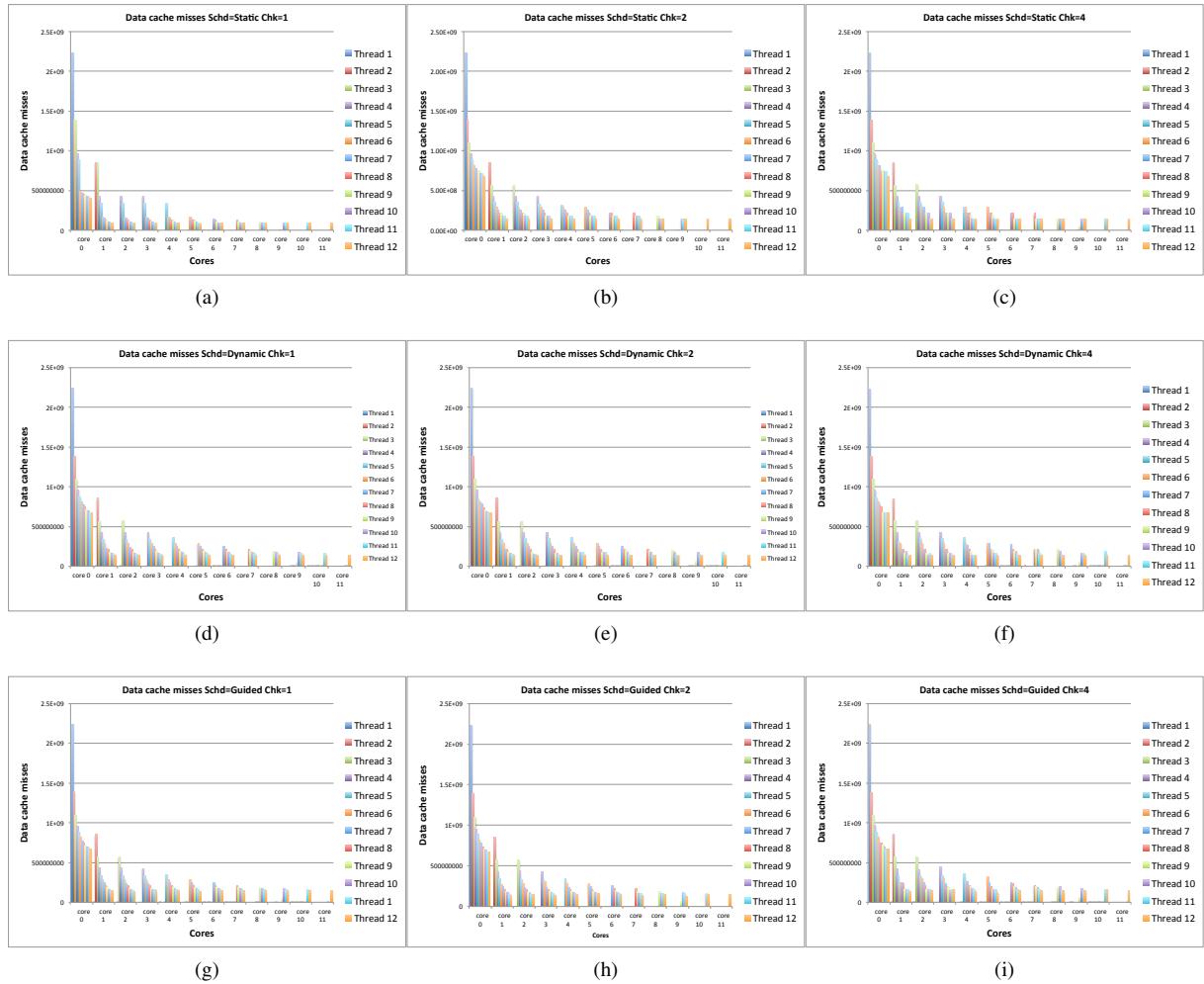


Figure 8: Direct cache misses for OpenMP scheduling and different chunk scheduling.

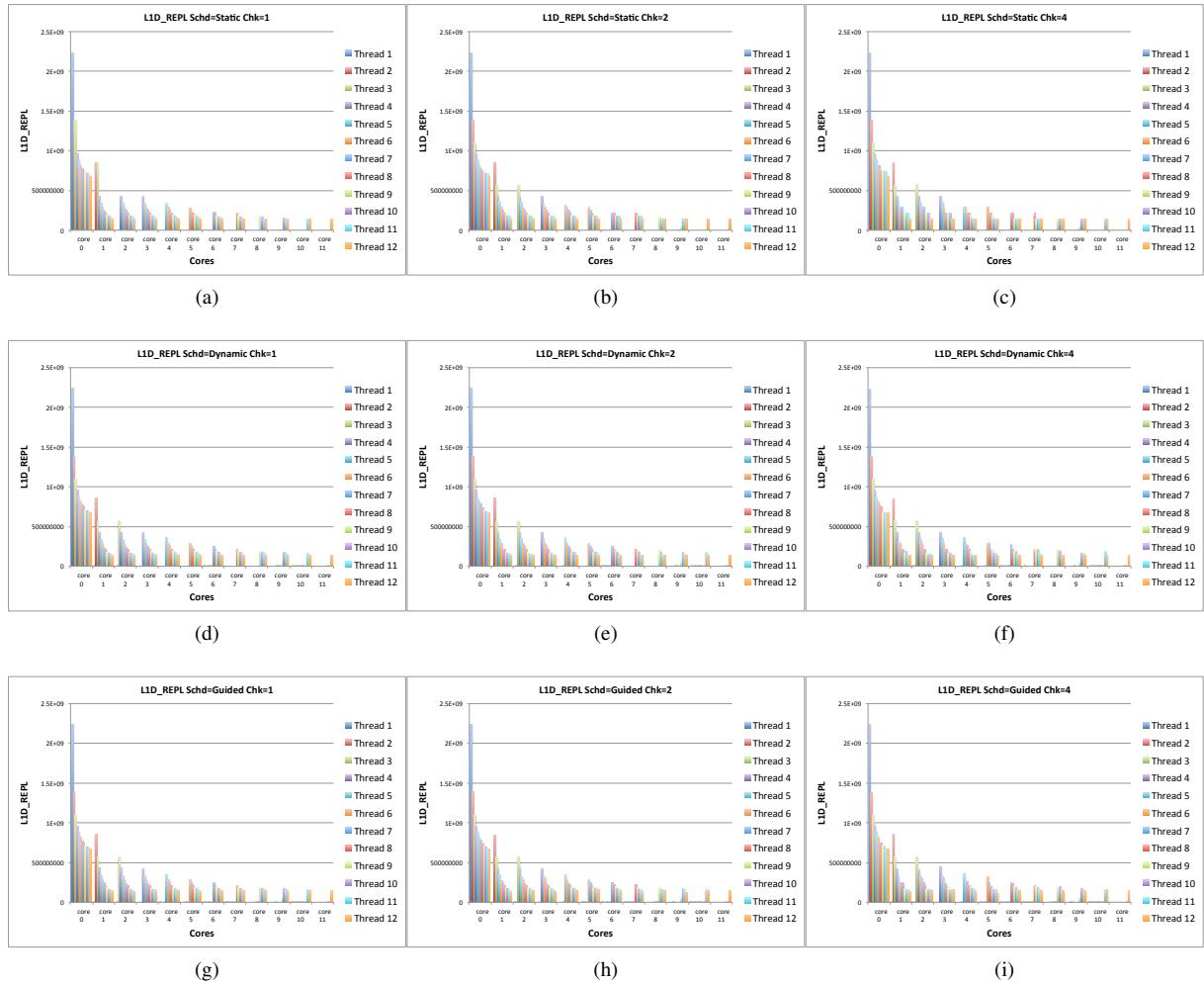


Figure 9: Number of lines that come into the cache for OpenMP scheduling and different chunk scheduling.