

pFANGS: Parallel High Speed Sequence Mapping for Next Generation 454-Roche Sequencing Reads

Sanchit Misra¹, Ramanathan Narayanan², Wei-keng Liao³, Alok Choudhary⁴

Electrical Engineering and Computer Science, Northwestern University, Evanston, IL, USA 60208

Email: ¹smi539@eecs.northwestern.edu, ²ran310@eecs.northwestern.edu,

³wklliao@eecs.northwestern.edu, ⁴choudhar@eecs.northwestern.edu

Simon Lin

Robert H. Lurie Comprehensive Cancer Center, Northwestern University, Chicago, IL, USA 60611

Email: S-Lin2@northwestern.edu

Abstract—Millions of DNA sequences (reads) are generated by Next Generation Sequencing machines everyday. There is a need for high performance algorithms to map these sequences to the reference genome to identify single nucleotide polymorphisms or rare transcripts to fulfill the dream of personalized medicine. In this paper, we present a high-throughput parallel sequence mapping program pFANGS. pFANGS is designed to find all the matches of a query sequence in the reference genome tolerating a large number of mismatches or insertions/deletions. pFANGS partitions the computational workload and data among all the processes and employs load-balancing mechanisms to ensure better process efficiency. Our experiments show that, with 512 processors, we are able to map approximately 31 million 454/Roche queries of length 500 each to a reference human genome per hour allowing 5 mismatches or insertion/deletions at full sensitivity. We also report and compare the performance results of two alternative parallel implementations of pFANGS: a shared memory OpenMP implementation and a MPI-OpenMP hybrid implementation.

Keywords—sequence mapping; next generation sequencers; 454 sequencers; parallel computing

I. INTRODUCTION

DNA sequencing is used in a variety of applications in medicine, such as SNP discovery, comparative genomics, gene expression, genotyping, metagenomics and personal genomics. Recent developments in Next Generation Sequencing (NGS) technology have resulted in affordable desktop-sized sequencers with low running costs and high throughput. These sequencers produce small fragments (reads) of the genome being sequenced as a result of the sequencing process. For example, the Illumina-Solexa system can generate 50 million sequences of length 30-50 nucleotides in just 3 days [7]. The Roche-454 system can generate 400,000 sequences of length 250-500 nucleotides in a 7.5 hour run [14]. The ABI-SOLiD system can also generate data at a similar rate [7]. NGS is a rapidly advancing field with a very high rate of increase in throughput. It is speculated that eventually the running costs of sequencing a genome will be as low as \$1000 [16]. This will trigger the use of such systems in laboratories around the world. The computational demands for processing NGS data are

tremendous and far exceed current capabilities. In fact, without substantial advances in high-performance, scalable algorithms, very little progress would be made to extract knowledge from such a rich set of data. Therefore, there is a need to design powerful algorithms and systems which can efficiently handle the computational challenges posed by NGSs.

An important step in many of the applications mentioned above is mapping a set of read sequences to a canonical genomic database. A typical genomic database, for instance, the human genome, can be 3 billion nucleotides in length. The length of the read sequences depends on the sequencing technology. In this article, we focus on the mapping of the longer reads produced by Roche-454 system. A 454 sequencer was recently used for sequencing the DNA sequence of James D. Watson to 7.4 fold redundancy in just two months [17]. The authors used BLAT [8] to map the 454 reads to a reference genome, which is not at par with the sequencing speed. Moreover, BLAT is designed for local alignment, while sequence mapping requires the entire length of a query to be mapped. There have been considerable efforts to develop faster sequence mapping tools which can match the speed of Next Generation Sequencers, but most of them have been for reads generated by Illumina-Solexa machines. Even though 454 sequencers are widely used by researchers, there has not been sufficient research to develop faster tools for mapping 454 reads. To the best of our knowledge, the only algorithms which are specifically designed for 454 data are BWA [9] and FANGS [12]. BWA is a package based on Burrows-Wheeler Transform (BWT). It supports gapped global alignment with respect to queries and is one of the fastest short read alignment algorithms while also finding suboptimal matches. However, [12] demonstrates that BWA suffers from low sensitivity. FANGS dynamically reduces the search space by using q-gram filtering and the pigeonhole principle, to rapidly map 454 reads onto a reference genome. FANGS allows a large number of mismatches and insertion/deletions. It tries to find all matches of a read in the reference genome and maps

nearly 100% of the reads. FANGS is shown to be upto an order of magnitude faster than the state-of-the-art techniques for 454 reads as long as the number of mismatches and insertion/deletions allowed is small. However, the execution time of FANGS increases dramatically with the increase in number of mismatches and insertion/deletions allowed. Therefore, there is a need to design powerful high throughput parallel programs and systems which can efficiently and accurately map 454 reads.

To the best of our knowledge, very little work has been done to parallelize sequence mapping algorithms. The approaches to parallelize high throughput sequence mapping tools typically include running a separate instance of the tool on each compute node and dividing the queries equally among these nodes. If the genome database occupies only a small amount of memory, this approach can give close to linear speedups. However, for large databases, like the human genome, the amount of memory required may not be available on one node. Moreover, the large memory requirement is also prone to having cache-misses and page-faults. While these sequential and parallel tools demonstrate a significant performance improvement over earlier sequence mapping tools, the throughput requirement of NGSs is also increasing rapidly and developing faster tools is constantly needed.

In this article, we describe our high-throughput parallel sequence mapping program pFANGS, a **parallel Fast Algorithm for Next Generation Sequencers**. pFANGS is a parallel implementation of FANGS. We discuss three parallel implementations of FANGS: (a) a shared memory task-parallel implementation using OpenMP, (b) an MPI-OpenMP task-parallel hybrid implementation, and (c) pFANGS: a fully data- and task-parallel MPI implementation (Section V). The first two implementations are based on query segmentation principle. The third implementation fully distributes the computational workload and data among all the processes and employs load-balancing mechanisms to ensure better process efficiency. We present the performance results in Section VI.

In comparison with existing tools, the most significant features of pFANGS are:

- High flexibility. It allows a large number of mismatches and insertions/deletions in mapping.
- High Sensitivity. It tries to find all the matches for each query and maps nearly 100% of the queries.
- Ability to handle large datasets. Using pFANGS, we have mapped approximately 31 Million queries of length 500 each to a reference human genome per hour allowing 5 mismatch or insertion/deletion at full sensitivity.
- Nearly linear scalability. With 512 processors, pFANGS achieves a speedup of upto 225 over the the time taken with 2 processors.

The remainder of the paper is organized as follows. We

give a formal definition of the problem in Section II followed by a background in Section III. Section IV describes the sequential FANGS algorithm in detail. We describe our parallel implementations in Section V followed by results in Section VI and conclusion in Section VII.

II. PROBLEM DEFINITION

The sequence alignment problem has been studied in great detail in literature. However, it has become even more significant in the wake of the new sequencing technologies in the form of Next Generation Sequencers. Consider, for example, using a 454 sequencer [14] to sequence a human genome. It produces a collection of small DNA fragments called reads. These reads are about 250-500 bases in length. Now, we need to search a read, Q , in the database consisting of a reference human genome, G . The database and the reads are from the genomes of different human beings. Moreover, there can be sequencing errors also. Hence, we may not be able to find an exact match of the read Q in the database. However, since both G and Q are from the genomes of the same species, we should be able to find a near-exact match of Q in G . Hence, while searching for Q in G , we only look for alignments which have less than a certain number of mismatches and insertion/deletions.

Given a string S over a finite alphabet Σ , we use $|S|$ to refer to the length of S , $S[i]$ to denote the i^{th} character of S and $S[i, j]$ to denote the substring of S which starts at position i and ends at position j . A q -gram of S is defined as a substring of S of length $q > 0$. A q -hit between two strings S_1 and S_2 is defined as the tuple (x, y) such that $S_1[x : x + q - 1] = S_2[y : y + q - 1]$. The *unit cost edit distance* between two strings S_1 and S_2 is defined as the minimum number of substitutions, insertions and deletions required to convert S_1 to S_2 [15]. We will use $edist(S_1, S_2)$ to refer to the *unit cost edit distance* between S_1 and S_2 . For a string S , we will refer to the natural decimal representation of S over Σ as $dec(S, \Sigma)$. For example, for $\Sigma = \{A, C, G, T\}$, the nucleotides A, C, G, T are mapped to the numbers 0, 1, 2, 3 respectively. Therefore:

$$f(A) = 0, f(C) = 1, f(G) = 2, f(T) = 3,$$

$$\text{And, } dec(S, \{A, C, G, T\}) = \sum_{i=0}^{|S|-1} 4^i f(S[i])$$

This brings us to the formal definition of the sequence mapping problem. We can represent every genomic sequence as a string over the alphabet $\Sigma = \{A, C, G, T\}$. Given a genomic database G of subject sequences $\{S_1, S_2, \dots, S_l\}$, a query sequence (read) Q of length m and an integer n , we are required to find all substrings from G , such that for each substring B , $edist(B, Q) \leq n$. We will denote the integer n as the *maxEditDist* parameter.

III. BACKGROUND

Since the arrival of BLAST [1] in 1990, many hash-table based sequence alignment methods have been proposed.

These include extremely popular tools like BLAT [8] and SSAHA [13]. BLAST has been the most popular tool for sequence alignment. However, it usually takes several hundreds of days for the data generated by the latest sequencers in just a few hours and hence is not a feasible option.

Recently, the advent of Next Generation Sequencers has inspired the researchers to develop high-speed sequence mapping tools. To the best of our knowledge, the only tools for sequence mapping of 454-Roche sequencing reads are FANGS [12] and BWA [9].

Since BLAST is the most popular sequence alignment tool, several attempts have been made to parallelize it. Early attempts at parallelization have used query segmentation approach [4], [5], where individual compute nodes independently search disjoint sets of queries against the whole database. This technique works well when the database can fit in the memory of a compute node. However, this approach suffers from caching and paging overheads when the database requires large amount of memory as the database is randomly accessed. This led to the development of database segmentation [2], [6], [10], [11], where the genomic database is evenly distributed across compute nodes. This reduces the caching and paging overheads as each compute node uses a small amount of memory for its part of the database. Database segmentation divides the database into mutually exclusive parts and assigns one part to each node. Every node searches for the query in its own part of the database and results from all processes are merged in the end. In particular, mpiBLAST [6] uses a master-worker paradigm in which the master gives each worker a batch of queries to process. Once a worker finishes its batch of queries, it notifies the master. If there are more queries to be processed, the master sends another batch of queries to the worker.

Six parallelization methods for short sequence mapping algorithms are proposed in [3]. The methods are general and should work for most hashing and indexing based algorithms. The first three methods are: (i) Partition Read Only (PRO) partitions the reads into equal parts and sends each part to one processing node. Each node keeps its own copy of the index of the whole genome. This method is useful to match very large number of reads to a relatively short reference genome. If the genome is large, the index may not fit in the memory available on one node. (ii) Partition Genome Only (PGO) partitions the genome equally amongst all processing nodes. Each node creates the index of only the assigned part of the genome and processes all reads against it. PGO performs well when the genome size is large and the number of reads is small but does not scale well if the number of reads is large. (iii) Suffix Based Assignment (SBA) assigns a set of suffixes to each processing node and makes them only responsible for genome and read sequences that end with the corresponding suffixes. The other methods are combinations of the first three methods. The authors compare the scalability of the proposed methods

using theoretical analysis and experimentation using SOLiD System Color Space Mapping Tool.

IV. FANGS

All the parallelization methods mentioned above are coarse-grained in the sense that they treat the sequential algorithm as one application and run separate instances of the sequential application on different parts of the queryset and the database. Many of the hashing and indexing based algorithms including FANGS essentially have four stages: 1) Creating the index of the database, 2) Finding hits in the database based on the index, 3) Reducing the number of hits to be processed using a filtering criteria, 4) Processing the list of hits remaining to get the final mapping. In this paper, we exploit this generic structure to achieve a more fine-grained parallelization. We parallelize each stage separately and perform load balancing between stages to achieve very high throughputs. Although we demonstrate our parallelization techniques using FANGS, they should be applicable to most hashing and indexing based algorithms. In this section, we describe the various stages of FANGS.

Preprocessing step: Creation of the q -gram index - We preprocess the sequences in the database by breaking them into non-overlapping q -grams and store the location of each q -gram in the q -gram index. We will refer to the q -gram index as the *index-table*. We refer to the size of these non-overlapping q -grams, q , as *tileSize*. Each q -gram t can be uniquely mapped to a corresponding integer $dec(t, \Sigma)$ as defined in Section II. For each q -gram t , we calculate two values: (1) $tileHead(t) = dec(t[1 : 12], \Sigma)$ and (2) $tileTail(t) = dec(t[13 : q], \Sigma)$. The *index-table* consists of two arrays. The first array *occurrenceTable* stores (i) the location of $t[1 : 12]$ in the database G and (ii) $tileTail(t)$ for each q -gram. Hence, *occurrenceTable* contains the concatenation of lists $L(t[1 : 12]) = \{i, tileTail(G[i : i + q - 1]) | G[i : i + 11] = t[1 : 12]\}$, where t is a q -gram, that is $t \in \Sigma^q$. For each q -gram $t \in \Sigma^q$, the position $tileHead(t)$ in the second array *lookupTable* contains the pointer $p(t)$, which points to the beginning of the corresponding list $L(t[1 : 12])$ in the *occurrenceTable*; and the count $c(t)$ of the number of occurrences of $t[1 : 12]$ in G . Hence the length of the *lookupTable* is $|\Sigma|^{12}$. In order to find hits for a q -gram t , it first indexes the *lookupTable* with $tileHead(t)$. Let $L(t[1 : 12])$ be the corresponding list. The q -hits can be found by traversing through the list and outputting those locations for which $tileTail(t)$ matches. Note that creation of index need to be done only once for a given value of q . After that, we can process any number of queries.

GetHits - The algorithm takes each overlapping q -gram in the query and finds the locations of all occurrence of the q -gram in the database using the *index-table*. The algorithm creates a q -hit with each location and adds it to the hitList. Each q -hit consists of two values - starting position of the q -gram in the query ($qStart$) and in the database ($dStart$).

FindRegions - FANGS uses the following corollary to filter out non-homologous regions.

Corollary: Given a query $Q[1..m]$ and database $G[1..L]$ ($m < L$). For all substrings α of G such that $edist(Q, \alpha) < n$, $\exists x, y$ such that $Q[x : x+q-1] = \alpha[y : y+q-1]$, $Q[x+q : x+2q-1] = \alpha[y+q : y+2q-1]$, \dots , $Q[x+(T-1)q : x+Tq-1] = \alpha[y+(T-1)q : y+Tq-1]$. In other words Q and α share a common substring of length T q -grams, where T is given by:

$$T = \lfloor \frac{\lfloor \frac{m}{n+1} \rfloor - (q-1)}{q} \rfloor$$

The substring α is called a homologous region of Q in G . Using the q -hits, FANGS finds regions in the database which have a common substring of length T q -grams. This filtering criteria significantly reduces the search space for finding homologous regions.

CheckRegions - Each potential homologous region is further processed by using an adaptation of the Needleman-Wunsch algorithm to check if the region actually has an edit distance $\leq n$.

The above algorithm, though very fast, is still not at par with the current sequencing speeds. Hence there is a need to parallelize the algorithm. For the human genome, the algorithm needs 1GB memory for the *index-table*. Moreover, it also needs to keep the database G in memory as it needs the database to create the *index-table* and also to examine the candidate homologous regions in the end. Hence, the algorithm requires about 4.5GB memory for mapping reads to a reference human genome. This large amount of memory usage can potentially lead to a number of cache-misses and page-faults due to random access and hence slows down the execution. Hence, we need to distribute both the index and the database across processor nodes in order to run it efficiently on a cluster.

V. PARALLEL APPROACHES

In this section, we investigate three parallelization approaches for FANGS: (a) a shared memory task-parallel implementation using OpenMP, (b) an MPI-OpenMP task-parallel hybrid implementation, and (c) a fully data- and task-parallel MPI implementation called pFANGS.

A. Shared Memory Parallel Implementation

Since the human genome database occupies significant amount of memory, a shared memory parallel implementation seems like a natural choice as we can load both index and database in the shared memory. The target platform is the parallel machine equipped with multiple CPU cores sharing a large sized main memory. We adopt the query segmentation strategy in which each thread takes a subset of the queries and processes them independently. In other words, we parallelize the outermost loop of the sequential algorithm. The algorithm divides the queries equally

amongst all threads. All the threads access the same copy of the database and the *index-table* stored in shared memory. Each thread uses FANGS to perform the alignments and stores the results in the *localOutputList* data structure. The *globalOutputList* is shared across all threads. Once a thread finishes processing all its queries it acquires exclusive access to the *globalOutputList* and concatenates its *localOutputList* to it. After all the children threads have merged their results to the global list, the parent thread writes all the outputs to the output file.

Accessing the shared data structures, such as genome database, *index table*, and *globalOutputList*, must be serialized in order to achieve data atomicity and cache coherence. This can become a major performance bottleneck as the number of threads increase.

B. MPI-OpenMP Hybrid Implementation

In order to overcome the drawbacks of the shared memory approach, we have also designed an MPI-OpenMP hybrid approach. This approach targets the parallel computers equipped with multiple SMP compute nodes interconnected with a high speed communication network and the memory in each node is not directly accessible to a remote node. In this hybrid approach, the *index-table* is built independently in each compute node. All processes running on the same node share the *index-table* by accessing the shared memory. The queries are evenly assigned to the MPI processes across all compute nodes. The alignment outputs produced at each node are saved locally, which are later sent to the root process. The root process concatenates all the partial results and writes to the output file.

There is a single MPI process running on each compute node and OpenMP is used to enable thread parallelism using all cores in each node. Even though the memory size per processing core is small, the combined shared memory of all cores on a node is sufficient to hold both the database and the *index-table*. Compared to the shared-memory method, this approach alleviates the congestion problem by reducing the number of processes accessing the shared memory. However, since we are using more than 4GB of the memory on each node, the problem of cache misses is still unsolved.

C. pFANGS: Fully Data and Task Distributed MPI Implementation

The above hybrid implementation may not be very scalable as it requires about 4.5GB memory per node. In this section we describe a completely task and data parallel MPI implementation, named pFANGS.

The idea is to distribute the entire database and the *index-table* equally among all MPI processes. Recall that the genomic database is available as a set of sequences $\{S_1, S_2, \dots, S_i\}$. As a preprocessing step, for each sequence S_i , we remove $|S_i| \bmod q$ nucleotides from the end so that the length of each sequence is a multiple of q . Then we

store all the sequences in a file named *genomeFile* by concatenating the sequences. To keep track of the positions of these sequences, we also maintain a metadata file that stores the name and length of each sequence.

PFANGS starts by having each process read an equal contiguous portion of the *genomeFile*. For example, if the length of the file is L and there are p processes, process 0 will read the first $\frac{L}{p}$ nucleotides, process 1 will read the second $\frac{L}{p}$ nucleotides and so on. Recall that the length of the *lookupTable* is 4^{12} . In order to create the *index-table* in a distributed manner, process 0 is responsible for the first $\frac{4^{12}}{p}$ entries of the *lookupTable* and the corresponding part of the *occurrenceTable*, process 1 is responsible for the second $\frac{4^{12}}{p}$ entries of the *lookupTable* and so on. Each process creates non-overlapping q -grams from its chunk of the database and sends each q -gram to the process responsible for the corresponding part of the *index-table*. After receiving all the q -grams, each process creates its part of the *index-table*. This way we create the *index-table* in a distributed manner. Each process discards its chunk of the database read after the creation of the *index-table*. In the query processing phase, the queries are equally divided to all processes. This phase is divided into five stages.

GetHits - Each MPI process takes its assigned queries and finds all the overlapping q -grams. Each q -gram is represented using three numbers: *tileHead*, *tileTail* and *queryId*; and stored in an array. The array is sorted according to the *tileHead* value. Hence, the q -grams whose corresponding *index-table* entries are on one process are located in contiguous locations in the array. The process then sends each q -gram to the process which has the corresponding part of the *index-table*. It also receives q -grams which correspond to its part of the *index-table*. For all these q -grams, it hashes the corresponding hits using the *index-table*. Each hit consists of three values: the database location (*dStart*) the query location (*qStart*) and the *queryId*.

RedistributeHits - We redistribute the hits across all processes such that after redistribution, all hits corresponding to one query are on one process and hits are approximately evenly divided across all processes.

FindRegions - Every process processes all the hits for each assigned query to obtain the candidate homologous regions. A candidate homologous regions consists of three numbers: (1) *dBegin*, (2) *dEnd* and (3) *queryId*, where *dBegin* and *dEnd* are the start and end positions of the candidate region in the database.

RedistributeRegions - The candidate homologous regions are redistributed across all processes such that the number of regions on each process is approximately equal and regions with close *dBegin* values are on the same process. In order to do this, we perform a global bucket sort on all the regions across all processes based on the *dBegin* value. Then we divide the sorted list of regions into equal parts and assign

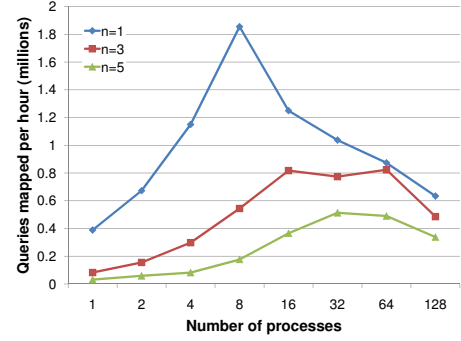


Figure 1. Number of queries mapped per hour to a reference human genome using the shared memory OpenMP implementation for different number of processes and different values of *maxEditDist*, n . Each query is 500 nucleotides long.

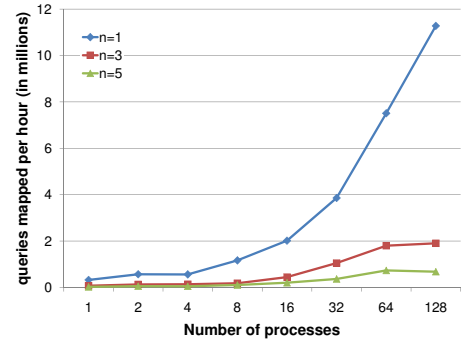


Figure 2. Number of queries mapped per hour to a reference human genome using the MPI-OpenMP hybrid implementation for different number of processes and different values of *maxEditDist*, n . We have used 2 shared memory cores per node. Each query is 500 nucleotides long.

one part to each process.

CheckRegions - Each process takes the list of regions (*regionList*) assigned to it. The regions are already sorted according to *dBegin* values as a result of the global bucket sort. Each process reads the genome database from minimum of *dBegin* values to maximum of *dEnd* values of all the assigned regions. Then it checks each candidate region one by one to see if the edit distance is indeed less than *maxEditDist*. All the homologous regions, which satisfy the criteria, are sent to the root process. The root process concatenates the results from all processes and writes them to the output file. Since the regions are processed in increasing order of *dBegin* values and regions with close *dBegin* values are on the same process, the disk IO cost due to random access is minimized.

VI. EXPERIMENTS AND PERFORMANCE ANALYSIS

In our experiments, the human genome database is used. The queries to be used to search against the database were randomly sampled from the human genome into reads of length 500. The number of queries is set to 10000 per process. In this section, we will use the word "process" to

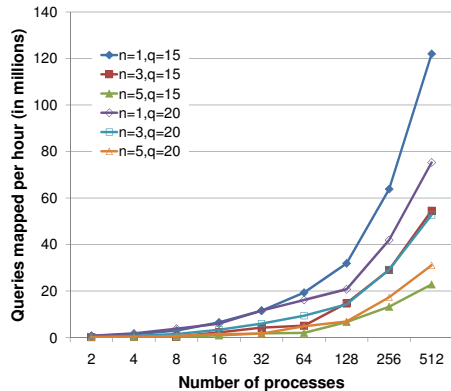


Figure 3. Number of queries mapped per hour to a reference human genome using the data and task parallel distributed memory MPI implementation for different number of processes and different values of $maxEditDist$, n and $tileSize$, q . Each query is 500 nucleotides long.

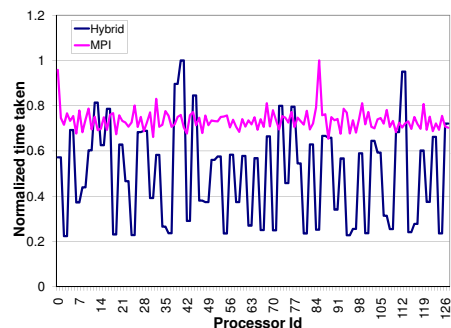


Figure 4. Comparison of load on each process for hybrid and MPI implementations. In order to obtain this plot, we normalized the time taken by each process by the maximum time taken by a process. The variances of load for hybrid and MPI implementations are 0.046 and 0.0021 respectively.

refer to MPI process as well as OpenMP threads, so that the presentation is consistent. The comparison of the speed, sensitivity and accuracy of FANGS with the other existing tools is given in [12]. The paper shows that sequential FANGS is upto an order of magnitude faster than the state-of-the-art techniques for 454-Roche reads of length 500 allowing 5 mismatches or insertion/deletions. To the best of our knowledge, there is no other published parallel implementation for mapping 454-Roche sequencing data. Hence, in this paper, we compare the parallel implementation with the sequential implementation of FANGS. All the parallel implementations give the same output as FANGS and hence retain the sensitivity and accuracy of FANGS (data not shown here).

The experiments of using the shared-memory approach were performed on the NCSA SGI Altix SMP machine (Cobalt). Cobalt has two SMP nodes with 512 1.6 GHz Intel Itanium 2 processors each. The machine has 4GB of memory per processor. The SMP machine is running SGI ProPack 5 and Intel 10.1 C compiler. Figure 1 shows the

number of queries mapped per hour by using the shared memory parallel implementation. Since all the processes are accessing the shared database and *index-table*, memory IO becomes a major bottleneck when the number of processes increases. It is clear that the shared memory approach does not scale. It improves the performance up to a certain number of threads beyond which the performance starts to come down. Such performance saturation and degradation are commonly seen on the SMP parallel machines, due to the contention on the system bus as well as the system overhead of cache coherence control.

The performance results of the MPI-OpenMP hybrid implementation were collected from the NCSA IA-64 Teragrid cluster (Mercury). Mercury consists of 887 IBM cluster nodes, 128 of which have dual 1.3 GHz Intel Itanium 2 processors and with 12GB of memory per node on which we ran our experiments. The cluster runs SuSE Linux and uses Myricom's Myrinet cluster interconnect network. Figure 2 presents the performance results. The hybrid approach scales much better than the pure shared-memory approach, as each node has its own copy of the database and the *index-table*. Only a limited number of threads share each copy of the hash table using OpenMP. However, even though there are a small number of OpenMP threads on each node, they still have to contest for memory access thereby resulting in a sub-optimal speedup. Recall that our algorithm needs about 4.5GB memory to execute. Such large memory requirement with random data accesses can cause significant cache misses.

The distributed memory implementation was also evaluated on Mercury. Figure 3 and Table II show the results. It can be clearly seen that the distributed memory implementation scales very well. As the database and *index-table* (See Table I for performance results of creation of the *index-table*) are distributed across all processes, the memory requirement on one process is smaller, thereby reducing the number of cache-misses and page-faults. Another important thing to note is that, for the shared memory and hybrid implementations, we statically divided the queries equally across all processes assuming that the amount of load is equal for equal number of queries. Since our sequence mapping algorithm is heuristic based, the actual run times and result sizes for queries are highly irregular and difficult to predict. For the 128 process case, figure 4 displays the load on each process for the hybrid and MPI implementations. It is clear that there is significant load-imbalance for the hybrid implementation, while the load for the MPI implementation is much better balanced. For the MPI implementation, two of the processes (0 and 85) always take significantly more time in the *CreateIndex* and *GetHits* stage. Our initial investigations reveal that the imbalance in load is due to the irregular nature of the genomic databases.

Figure 5 shows the breakdown of time spent on each stage of the query processing phase for various values of

| # proc | Time taken (seconds) | | | | | |
|--------|----------------------|-------|-------|----------|-------|-------|
| | $q = 15$ | | | $q = 20$ | | |
| | $n=1$ | $n=3$ | $n=5$ | $n=1$ | $n=3$ | $n=5$ |
| 2 | 300.5 | 299.1 | 299.0 | 225.9 | 231.1 | 225.0 |
| 4 | 159.2 | 160.1 | 161.7 | 117.6 | 117.0 | 123.8 |
| 8 | 87.0 | 87.1 | 80.1 | 58.5 | 63.5 | 64.6 |
| 16 | 46.8 | 45.9 | 46.9 | 39.3 | 34.8 | 34.5 |
| 32 | 31.4 | 32.5 | 31.4 | 27.3 | 21.9 | 25.9 |
| 64 | 27.6 | 28.0 | 27.5 | 15.2 | 15.2 | 17.1 |
| 128 | 21.6 | 21.9 | 21.9 | 12.9 | 14.0 | 14.1 |
| 256 | 20.7 | 20.7 | 21.6 | 11.8 | 12.2 | 14.5 |
| 512 | 21.4 | 21.3 | 21.8 | 13.6 | 12.6 | 14.1 |

Table I

ABSOLUTE TIME TAKEN FOR THE CREATION OF INDEX TABLE FOR DIFFERENT VALUES OF n AND q , FOR DIFFERENT NUMBER OF PROCESSORS.

maxEditDist and number of processes used. As the value of *maxEditDist* increases, more candidate homologous regions are generated by *FindRegions* algorithm since the value of T gets smaller. Moreover, *CheckRegions* stage has higher computational complexity as compared to other stages. Hence for larger values of *maxEditDist*, *CheckRegions* stage consumes more than 85% of the overall execution time. Note that each region can be examined independently of all other regions. We dynamically balance the load across processes by redistributing the candidate homologous regions evenly across processes to achieve better process efficiency. Also note that the percentage of time spent on the *RedistributeHits* stage increases as the communication time increases with the increase in the number of processes. This stage may become a bottleneck and hinder scalability as the number of processes increase. To avoid this, for larger number of processes, we divide them into disjoint subsets of 128 processes each. The queries are equally divided among these subsets. Each of these subsets work independently by creating their own copy of the index. As a result of this, the percentage of time spent on the *RedistributeHits* stage does not increase as the number of processes increase beyond 128. Notice from Figure 3 that we can process up to 31061118 queries per hour for $n = 5$ using 512 processors. Since each query is of length 500, this means we can map 454-Roche reads with a total of $31061118 * 500 = 15.53$ Billion nucleotides per hour against a reference human genome. Hence, with 512 processors, we are able to map 454/Roche reads of 5.17x coverage of a human genome to a reference human genome per hour allowing 5 mismatches or Indels at full sensitivity. In other words, we can map 5.17 human genomes per hour.

VII. CONCLUSION

Advances in sequencing techniques necessitate the development of high performance, scalable algorithms to extract biologically relevant information from these datasets. In this paper, we investigate different parallel implementations of a fast sequence alignment tool FANGS. Firstly we develop

| # proc | Speedup with respect to two processes | | | | | |
|--------|---------------------------------------|-------|-------|----------|-------|-------|
| | $q = 15$ | | | $q = 20$ | | |
| | $n=1$ | $n=3$ | $n=5$ | $n=1$ | $n=3$ | $n=5$ |
| 2 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 |
| 4 | 2.1 | 1.0 | 2.0 | 2.2 | 2.2 | 2.0 |
| 8 | 4.3 | 1.7 | 4.4 | 4.8 | 4.5 | 3.9 |
| 16 | 9.5 | 9.1 | 9.0 | 7.5 | 10.1 | 9.1 |
| 32 | 16.8 | 17.5 | 17.5 | 14.5 | 17.9 | 10.2 |
| 64 | 28.1 | 21.0 | 19.6 | 20.3 | 28.1 | 29.4 |
| 128 | 46.4 | 60.8 | 65.5 | 26.0 | 42.8 | 41.9 |
| 256 | 92.9 | 119.6 | 130.7 | 52.5 | 87.5 | 105.1 |
| 512 | 177.6 | 224.2 | 225.4 | 94.3 | 158.1 | 188.5 |

Table II

SPEEDUP FOR PROCESSING STAGE OF pFANGS WITH RESPECT TO TIME TAKEN BY TWO PROCESSES FOR DIFFERENT VALUES OF n AND q , FOR DIFFERENT NUMBER OF PROCESSES.

query segmentation based OpenMP and MPI-OpenMP hybrid implementations and discuss their limitations. We then develop a highly optimized data- and task-distributed MPI implementation with intelligent load-balancing techniques that avoid problems of memory bandwidth and cache misses. Our experimental evaluation shows that this technique results in excellent load-balance and process efficiency and hence yield close to linear speedups.

With the advent of new technologies, we will need even faster sequence mapping tools to stay at par with the increasing sequencing speed. With the development of better parallel algorithms, we can setup huge processing centers which contain a large number of sequencers producing reads and huge clusters working in tandem to rapidly process them to extract a variety of information. The Next Generation Sequencers along with high-speed sequence processing systems will enable us to realize the dream of personal genomics. This can help us in using a patient's DNA in diagnosing a disease or even knowing in advance whether a person's DNA encodes a risk of a certain disease.

ACKNOWLEDGMENT

This work was supported in part by NSF CCF-0621443, NSF SDCI OCI-0724599, NSF CNS-0551639 and IIS-0536994, DOE SCIDAC-2: Scientific Data Management Center for Enabling Technologies (CET) grant DE-FC02-07ER25808.

REFERENCES

- [1] S. F. Altschul, W. Gish, W. Miller, E. W. Myers, and D. J. Lipman. Basic local alignment search tool. *J Mol Biol*, 215(3):403–410, October 1990.
- [2] R. Bjornson, A. Sherman, S. Weston, N. Willard, and J. Wing. Turboblast(r): A parallel implementation of blast built on the turbohub. *Parallel and Distributed Processing Symposium, International*, 2:0183, 2002.
- [3] D. Bozdag, C. C. Barbacioru, and U. V. Catalyurek. Parallel short sequence mapping for high throughput genome sequencing. *Parallel and Distributed Processing Symposium, International*, 0:1–10, 2009.

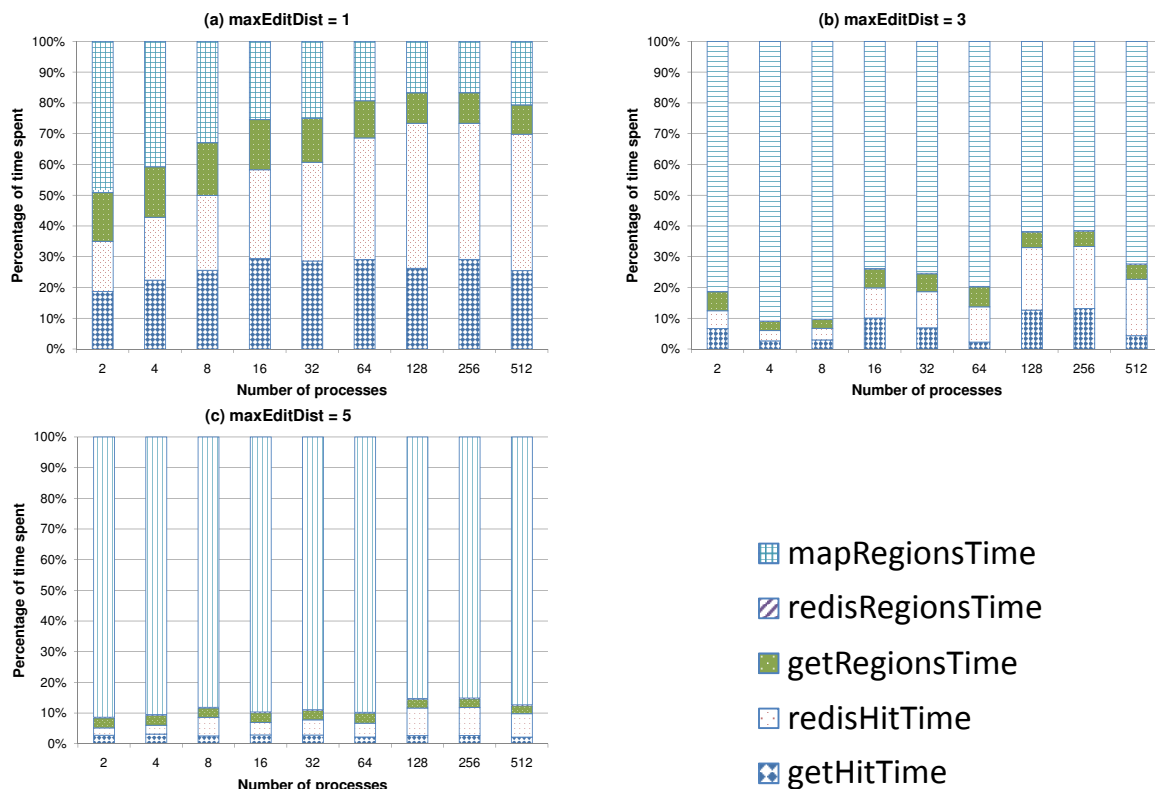


Figure 5. Percentage of time spent each stage of the distributed memory MPI implementation for different number of processes and different values of maxEditDist, n

- [4] R. C. Braun, K. T. Pedretti, T. L. Casavant, T. E. Scheetz, C. L. Birkett, and C. A. Roberts. Parallelization of local blast service on workstation clusters. *Future Gener. Comput. Syst.*, 17(6):745–754, 2001.
- [5] E. Chi, E. Shoop, J. Carlis, E. Retzel, and J. Riedl. Efficiency of shared-memory multiprocessors for a genetic sequence similarity search algorithm. *Technical Report TR97-005*, 1997.
- [6] A. E. Darling, L. Carey, and W. C. Feng. The design, implementation, and evaluation of mpiblast.
- [7] H. Jiang and W. H. Wong. Seqmap : mapping massive amount of oligonucleotides to the genome. *Bioinformatics*, pages btn429+, August 2008.
- [8] W. J. Kent. Blat—the blast-like alignment tool. *Genome Res*, 12(4):656–664, April 2002.
- [9] H. Li and R. Durbin. Fast and accurate long read alignment with burrows-wheeler transform. *Bioinformatics*, January 2010.
- [10] H. Lin, X. Ma, P. Chandramohan, A. Geist, and N. Samatova. Efficient data access for parallel blast. In *Parallel and Distributed Processing Symposium, 2005. Proceedings. 19th IEEE International*, page 72b, 2005.
- [11] D. Mathog. Parallel blast on split databases. *Bioinformatics*, 19:1865–1866, 2003.
- [12] S. Misra, R. Narayanan, S. Lin, and A. Choudhary. Fangs: High speed sequence mapping for next generation sequencers. In *Proceedings of ACM Symposium of Applied Computing (ACM SAC)*, Sierre, Switzerland, 2009.
- [13] Z. Ning, A. J. Cox, and J. C. Mullikin. Ssaha: A fast search method for large dna databases. *Genome Res.*, 11(10):1725–1729, 2001.
- [14] K. L. Patrick. 454 life sciences: Illuminating the future of genome sequencing and personalized medicine. *Yale J Biol Med.*, 80(4):191–4, Dec 2007.
- [15] K. R. Rasmussen, J. Stoye, and E. W. Myers. Efficient q-gram filters for finding all epsilon-matches over a given length. *Journal of Computational Biology*, 13(2):296–308, 2006.
- [16] C. Shaffer. Next-generation sequencing outpaces expectations. *Nature Biotechnology*, 25(2):149, February 2007.
- [17] D. A. Wheeler, M. Srinivasan, M. Egholm, Y. Shen, L. Chen, A. Mcguire, W. He, Y.-J. Chen, V. Makhijani, G. T. Roth, X. Gomes, K. Tartaro, F. Niazi, C. L. Turcotte, G. P. Irzyk, J. R. Lupski, C. Chinault, X.-Z. Song, Y. Liu, Y. Yuan, L. Nazareth, X. Qin, D. M. Muzny, M. Margulies, G. M. Weinstock, R. A. Gibbs, and J. M. Rothberg. The complete genome of an individual by massively parallel dna sequencing. *Nature*, 452(7189):872–876, April 2008.