# A Domain Decomposition Strategy for Alignment of Multiple Biological Sequences on Multiprocessor Platforms

Fahad Saeed

fsaeed2@uic.edu

Ashfaq Khokhar

ashfaq@uic.edu

*Department of Electrical and Computer Engineering*

*University of Illinois at Chicago*

*Chicago IL, USA*

October 13, 2018

## Abstract

*Multiple Sequences Alignment (MSA) of biological sequences is a fundamental problem in computational biology due to its critical significance in wide ranging applications including haplotype reconstruction, sequence homology, phylogenetic analysis, and prediction of evolutionary origins. The MSA problem is considered NP-hard and known heuristics for the problem do not scale well with increasing number of sequences. On the other hand, with the advent of new breed of fast sequencing techniques it is now possible to generate thousands of sequences very quickly. For rapid sequence analysis, it is therefore desirable to develop fast MSA algorithms that scale well with the increase in the dataset size. In this paper, we present a novel domain decomposition based technique to solve the MSA problem on multiprocessing platforms. The domain decomposition based technique, in addition to yielding better quality, gives enormous advantage in terms of execution time and memory requirements. The proposed strategy allows to decrease the time complexity of any known heuristic of $O(N)^x$ complexity by a factor of $O(1/p)^x$, where $N$ is the number of sequences, $x$ depends on the underlying heuristic approach, and $p$ is the number of processing nodes. In particular, we propose a highly scalable algorithm, Sample-Align-D, for aligning biological sequences using Muscle system as the underlying heuristic. The proposed algorithm has been implemented on a cluster of workstations using MPI library. Experimental results for different problem sizes are analyzed in terms of quality of alignment, execution time and speed-up.*

# 1   Introduction

Multiple Sequences Alignment (MSA) in computational biology provides vital information related to the evolutionary relationships, identifies conserved motifs, and improves secondary and tertiary structure prediction for RNA and proteins. In theory, alignment of multiple sequences can be achieved using pair-wise alignment, each pair getting alignment score and then maximizing the sum of all the pair-wise alignment scores. Optimizing this score, however, is NP-complete [1] and dynamic programming based solutions have complexity of $O(L^N)$, where $N$ is the number of sequences and $L$ is the average length of a sequence. Such accurate optimizations are not practical for even small number of sequences, thus making heuristic algorithms a feasible option. The literature on these heuristics is vast and includes widely used works, including Notredame et al. [2], Edgar [3], Thompson et al. [4], Do et al. [5], Lassmann et al. [6], Sze et al. [7], Schwartz et al. [8] and Morgenstern et al. [9]. These heuristics are complex combination of ad-hoc procedures with some flavor of dynamic programming. Despite the usefulness of these widely used heuristics, they scale very poorly with increasing number of sequences.

The high computational costs and poor scalability of existing MSA algorithms make the design of multiprocessor solutions highly desirable. Also, recent advances in the sequencing techniques such as pyrosequencing [10] are enabling fast generation of large amount of sequence data. For example, at the time of writing this paper, UniProtKB/Swiss-Prot contains 366226 sequence entries, comprising 132054191 amino acids representing 11342 species. Comparing this with less than 50k sequences in 1995, gives a glimpse of exponential growth in biological data. If useful research has to proceed, the performance of multiple alignment systems has to scale up accordingly with the enormous amount of data being generated.

The main goal of the work presented in this paper is to investigate domain decomposition strategies for biological data computations. For the multiple sequence alignment problem discussed in this paper, we use k-mer rank, a metric that depicts similarity of a sequence compared to another sequence, to partition the input data set into load balanced subsets. Subsequently, we show how these decomposed subsets of sequences can be aligned on multiple processors in a distributed fash-

ion, and glued together to get a highly accurate alignment of multiple sequences. Our approach is capable of aligning a large number of sequences (of the order of 20000 sequences [11]), with time complexity scaled down by a factor of $O(1/p)^4$,where $p$ is the number of processors, achieving super-linear speedups on multiprocessors without compromising quality of the alignment.

The rest of the paper is organized as follows. We start with a brief problem statement and background information relevant to our discussions in Section 2. Also, we discuss existing parallel approaches to the MSA problem and identify their limitations. In Section 3, we discuss the proposed domain decomposition based MSA algorithm for aligning protein sequences. This is followed by a rigorous analysis of the computation and communication costs. Section 4 presents the experimental results and analyzes these results in terms of alignment quality, execution time, memory usage, and speedup. Section 5 presents the conclusions and outlines future research.

## 2   Problem Statement and Background Information

We first define the Multiple Sequences Alignment (MSA) problem in simplest form, without indulging with the issues such as scoring functions, which are beyond the scope of this work. Let $N$ sequences be presented as a set $S = \{S_1, S_2, S_3, \cdots, S_N\}$ and let $S' = \{S'_1, S'_2, S'_3, \cdots, S'_N\}$ be the aligned sequence set, such that all the sequences in $S'$ are of equal length, have maximum overlap, and the total alignment score is maximized according to some scoring mechanism suitable for the application.

The method followed in most of the existing multiple alignment systems is that a quick pairwise alignment of the sequences is performed, giving a similarity matrix. This similarity matrix is then used to build a guide tree, which is then used to perform a progressive profile-profile alignment. Note that profile-profile alignments are used to re-align two or more existing alignments. Profile-profile alignment is a useful method as it can be used to gradually add new sequences to already aligned set of sequences, also referred to as progressive alignment. It can also be used to maintain one fixed high quality profile and keep on adding sequences aligned to that fixed pro-

file [4]. These are the basic steps that are followed by almost all distance based multiple sequence alignment methods [12]. Fig. **??** shows a generic MSA scheme, which is also used in Clustalw [4]. The first stage is a pair-wise comparison of the sequences under consideration. The second stage corresponds to the construction of guiding tree, which is later on used in stage three to perform final profile-profile alignments. To improve the alignment score, various iterative methods have been introduced in the later stages, as in the Muscle System [3].

## 2.1   Related Research

There have been numerous attempts to parallelize existing sequential multiple sequences alignment systems. Clustalw [4] is by far the most parallelized multiple sequence alignment system. James et. al. in [13] parallelized Clustalw for PC clusters and distributed shared memory parallel machines. HT Clustal is a parallel solution for heterogeneous multiple sequence alignment and MultiClustal is a parallel version of an optimized Clustalw [14] Zola et al. [15] provided the first parallel implementation of T-Coffee based on MPI. Different modules of the Muscle system have also been parallelized [16]. Other parallelization efforts include parallel multiple sequence alignment with phylogeny search by simulated annealing by Zola et al. [17], Multithreading Clustalw for multiple sequence alignment by Chaichoompu et al. [18] and Schmollinger et al. parallel version of Dialign [19].

Although there seems to be a considerable amount of effort to improve the running times for aligning large number of sequences using parallel computing, it must be noted that all the existing solutions have been aimed at parallelizing different modules of a known sequential system. Therefore the parallelism achieved has been limited to the usage of the function being parallelized. None of the existing parallel alignment approaches has been able to exploit the data parallelism, simply because of the lack of a domain decomposition strategy. A few attempts [20] [21] have also been made to cut each sequence into pieces and compute a piecewise alignment over all the sequences to achieve multiple sequences alignment. In [21], each sequence is 'broken' in half, and halves are assigned to different processors. The Smith-Waterman [20] algorithm is applied to these

divided sequences. The sequences are aligned using dynamic programming technique, and then combined using Combine and Extend techniques [20]. The Combine and Extend methods follow certain models defined to achieve alignment of the combination of sequences. These methods pay little or no attention to the quality of the results obtained. The end results have considerable loss of sensitivity. The constraints in these methods are solely defined by the models used, thus limiting the scope of the methods for wide variety of sequences.

Domain decomposition has been pursued for a large number of application in numerous fields, including elliptic partial differential equations, image processing, graphics simulations, fluid dynamics, astronomical and atmospheric calculations [22]. Most of these applications, take advantage of parallel processing by decomposing the data domains, and using data parallel techniques to achieve high performance.

In this paper, we investigate a data parallel approach to align multiple protein sequences, consequently *decreasing* computational effort in terms of time and memory while *improving* or obtaining quality comparable to other multiple alignment systems.

# 3   Proposed Distributed MSA Algorithm: Sample-Align-D

In this section, we present details of the domain decomposition strategy and the alignment algorithm, referred to as Sample-Align-D. We also analyze the computation and communication complexities of the proposed algorithm.

The proposed domain decomposition strategy draws its motivation from the Sample-Sort approach [23] that has been introduced to sort a very large set of numbers on distributed platforms. The sorting and MSA problems share a common characteristic, i.e., any correct solution requires implicit comparison of each pair of data items. In Sample-Sort, a small sample ($\ll N$) representing the entire data set is chosen over distributed partitions using some sampling technique such as Regular Sampling [24]. Then each item can be independently compared and ranked against this sample. If the sample is a true representative of the underlying data set, it eliminates the need

for explicit comparison of every item with the entire set. This way an $N$ size sorting problem is reduced to solving $p$ independent sorting problems of size $N/p$. We use a similar sampling approach to the domain decomposition of sequences over all the processors. In the case of sorting of integer numbers, numerical values of the numbers are compared to compute the rank of each number. In the case of multiple sequence alignment problem, we need to identify a unique feature of the sequence that could be used to compute the rank of each sequence, in terms of degree of similarity with other sequences in the set. This rank information can then be used to partition the smaller input subsets based on similarity and align smaller subsets of similar sequences independently. We propose to use k-mer distance [25] as a metric to determine the similarity of a sequence with any other sequence. Intuitively, the k-mer distance between any two sequences is based on the relative frequency of repetitive substrings of size $k$ in the sequences. Edgar in [25] showed that k-tuple similarities correlate well with fractional identity, and the small values of $k$ between 4 and 6 work well for biological sequences. For the sake of completeness, in the following, we provide the formal definition of k-mer distance.

**k-mer Rank:** Let's assume that a biological sequence is represented by a string $X$ of $n$ characters taken from an alphabet $\mathring{A}$ that contains $c$ different characters $(a_1, \cdots, a_c)$. For the words of length $k$ (hence named k-mers), there are $\epsilon = c^k$ such different words. We represent the set of k-mers in $X$ by vector $c^X = (c_1^X, \cdots, c_\epsilon^X)$. The distance between string $X$ and any arbitrary string $Y$, of length $m$, is calculated using $c_i^X$ and $c_i^Y$, the count of k-mer occurrences in $X$ and $Y$. Now let $C_i^{XY} = min(c_i^X, c_j^Y)$ denote the common k-mer count.

$$F(X, Y) = \sum_{i=1}^{\epsilon} \frac{C_i^{XY}}{[min(n, m) - k + 1]} \tag{1}$$

$$d^{F(X,Y)} = -log(\Delta + F(X, Y)) \tag{2}$$

where $F$ is the fraction of common k-mers between $X$ and $Y$, and $d^F$ transforms this into a distance. $\Delta$, is a small constant added to prevent logarithm of zero. Based on k-mer distance, we define *k-mer rank* as follows:

$$R_i = \frac{1}{N} \sum_{j=1}^{N} d^{F(i,j)} \tag{3}$$

An intuitive outline of the proposed distributed multiple sequence alignment solution, referred to as Sample-Align-D, is given in Algorithm 1.

---

**Algorithm 1** Sample-Align-D Intuitive Description

---

**Require:** $p$ processor for computation
**Require:** $N$ sequences of amino acids $S_1, S_2, \cdots, S_N$:
**Ensure:** Multiple alignment of $N$ sequences

1. In parallel, calculate the global k-mer rank for each sequence in each processor

2. Redistribute the sequences using the k-mer ranks such that sequences with similar k-mer ranks are accumulated on the same processor.

3. In parallel, align the sequences on each processor using any sequential multiple sequence alignment (MSA) system

4. Calculate the global ancestor using local ancestors produced by the local alignments at each processor in the previous step.

5. In parallel, fine tune the alignment on each processor using the global ancestor

---

## 3.1   k-mer Rank based Decomposing Domain

Our aim is to decompose the data set into subsets such that the sequences within a subset are more similar to each other than the sequences in other subsets. Ideally, this can be accomplished by partitioning the phylogenetic tree in a load balanced fashion such that the partitions also minimize communication across processors. The partitioning of a problem graph can be performed by making a virtual grid over the graph structure [26]. This works well only if the problem graph is uniform. In MSA, the phylogentic trees are rarely uniform. Thus partitioning with naive techniques will lead to non-uniform loads. However, most of the existing partitioning techniques for non-uniform problem graphs [27–31] cannot be used in this case due to progressive alignment dependencies in MSA, as discussed extensively in [15]. In the following, we outline a novel domain decomposition strategy that is based on k-mer rank similarities.

There are three main parameters that may contribute to the computation load in the MSA problem. These include: the number of sequences, the length of the sequences, and the similarity rank (that we call k-mer rank as discussed in the previous sections). However, as our analysis will reveal in the later sections, the lengths of the sequences do not contribute much computationally. Hence we can safely neglect the length of the sequences for load-balanced partitioning, and consider the k-mer rank and the number of sequences for partitioning and mapping. We will use a novel sampling based strategy to compute *global* k-mer ranks.

### 3.1.1  Globalised k-mer Rank

For a highly divergent set of sequences, k-mer rank computed for each sequence locally on each processor using only $N/p$ sequences would be different from the k-mer rank computed using all the $N$ sequences. In order to address this problem, we sample $k$ sequences from each processor such that the k-mer ranks of these $k$ samples represent the ranks of the corresponding set of $N/p$ sequences, yielding a total of $k \times p$ samples. Collectively, it is safe to assume that these $k \times p$ samples represent the entire set of $N$ sequences. The k-mer rank based ordering of these $k \times p$ sequences yields a phylogenetic tree of the samples, which in turn represent all the sequences. Each processor re-computes the k-mer ranks of its sequences using this global sample. Subsequently, redistribution based on this new k-mer rank also ensures that sequences accumulated in each processor are 'similar' to each other.

In Fig. 1, we plot the k-mer ranks computed using samples (referred to as globalized ranks) and using all the sequences collectively (referred to as centralized ranks). As depicted in this figure, the curves have high degree of similarity. The statistics of the two approaches for 500 sequences are presented in Table. 1. As can be seen that the standard deviation for the two sets of ranks is very low ( 0.58). This shows that the k-mer ranks for $N$ sequences computed using a global sample is statistically indistinguishable from the k-mer ranks computed using all the $N$ sequences.
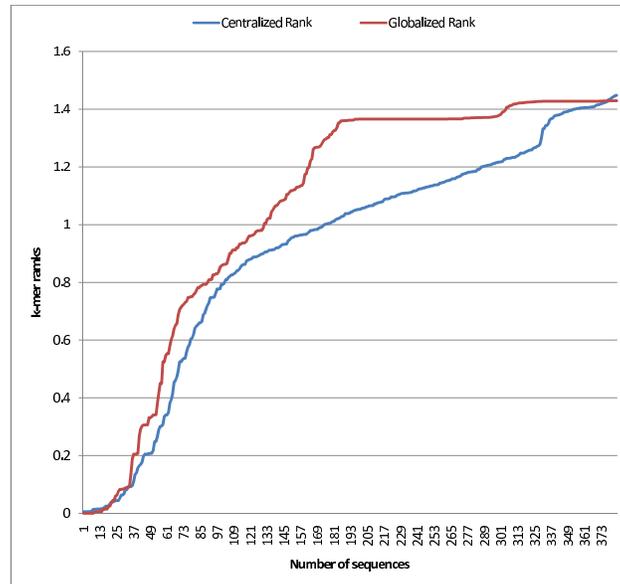
Figure 1: Globalized ranks and Centralised Ranks

Table 1: Comparison of the k-mer ranks computed using global sample (globalized) and using the entire set of sequences (centralized).

| (Maximum, Minimum) Central | (1.44827, 0.0) |
|---|---|
| Average Centralized | 0.722962 |
| (Maximum, Minimum) Globalized | (1.46207,0.0) |
| Average Globalized | 1.11302 |
| Variance w.r.t. Centralized | 0.33190 |
| Standard Dev. w.r.t Centralized | 0.576377 |

### 3.1.2   Redistribution Based on Globalized k-mer Rank

Each processor computes the k-mer ranks of its $w = N/p$ sequences locally using all the $N/p$ sequences, and sorts the sequences based on this local k-mer rank. Here $N$ is the number of sequences in the input and $p$ is the number of processors. From each of the $p$ locally sorted lists, $k = (p-1)$ evenly spaced samples are chosen. The k-mer ranks of these $(p-1)$ samples (pivots) divide the local set into $p$ ordered subsets. The k-mer ranks of these $p-1$ samples from each processor are gathered at the root processor, yielding a set $Y$ of size $p(p-1)$ ranks.

This regular-sampled set $Y$ is sorted to compute the ordered list $Y_1, Y_2, Y_3, \cdots, Y_{p(p-1)}$ determining the range of k-mer ranks over all the processors. Then ranks $Y_{p/2}, Y_{p+p/2}, \cdots, Y_{(p-2)p+p/2}$

are chosen as pivots ($p$ in total) dividing the k-mer rank range into $p$ buckets. These pivots are then broadcast to all the processors. Each processor sends the sequences having k-mer ranks in the range of bucket $i$ to processor $i$. For the bound on the size of the dataset in each processor after redistribution, we refer to the analysis in Section 3.2.

### 3.1.3   The Alignment

Next, a sequential MSA program is executed on each processor. Since our ultimate goal is to have a global alignment of all the $N$ sequences, a procedure has to be devised to concatenate these 'chunks' of *locally* (here *locally* is defined as the chunk of sequences that are aligned on a single processor) aligned sequences so that the *global* alignment of multiple sequences is achieved. Edgar in [32] has observed that multiple sequences alignment for homologous sequences can be obtained by aligning each sequence to the *root* profile. This approach is similar to the one used in the PSI-BLAST, where a *known* profile is used to align any query sequence with the sequences that have generated the profile. This technique may also be categorized as *template* based method, as observed by Notredame in his recent work [12]. We use a similar concept along with domain decomposition of the sequences. We extract the local ancestor from each processor after *locally* aligning each subset in parallel. All of these local ancestors are collected at the root processor and are aligned using a sequential multiple sequence alignment algorithm. The ancestor of all the local ancestors, referred to as the global ancestor, is then broadcast to all the processors. Subsequently, the global ancestor is used to perform a profile-profile alignment. That is, each set of the locally aligned sequences (referred to as profile) in each processor is aligned with the global ancestor profile.

In order to apply pair-wise alignment functions to profiles, a Profile Sum of Pairs (PSP) scoring function must be defined. We use the same PSP score as defined in [33] and [32]:

$$PSP^{xy} = \sum_i \sum_j f_i^x f_j^y \log(p_{ij}/p_i p_j) \qquad (4)$$

Here $x$ and $y$ are the profiles being aligned, $i$ and $j$ are the amino acid in profiles, $p_i$ is the

background probability of $i$, $p_{ij}$ is the joint probability of $i$ and $j$ aligned to each other, $f_i^x$ is the observed frequency of $i$ in column $x$ of the first profile, and $x_G$ is the observed frequency of gaps in that column. The same attributes are assumed for the profile $y$. For our purposes, we will take advantage of PSP functions based on the 200 PAM matrix [34] and the 240 PAM VTML matrix [35]. Some multiple alignment methods implement different scoring functions such as Log expectation (LE) functions, but for our purposes PSP scoring suffices. Of course, future work on decomposition strategies might investigate such functions in this context.

This fine tuning step based on ancestor profile is depicted in Fig. 2. For a highly divergent sequences set, we propose an additional step, in which profiles can be added to the root processor with respect to their similarity rank. This does not change the computation or communication costs, but gives the effect of 'profile-progressive' sort of gluing in the root processor.
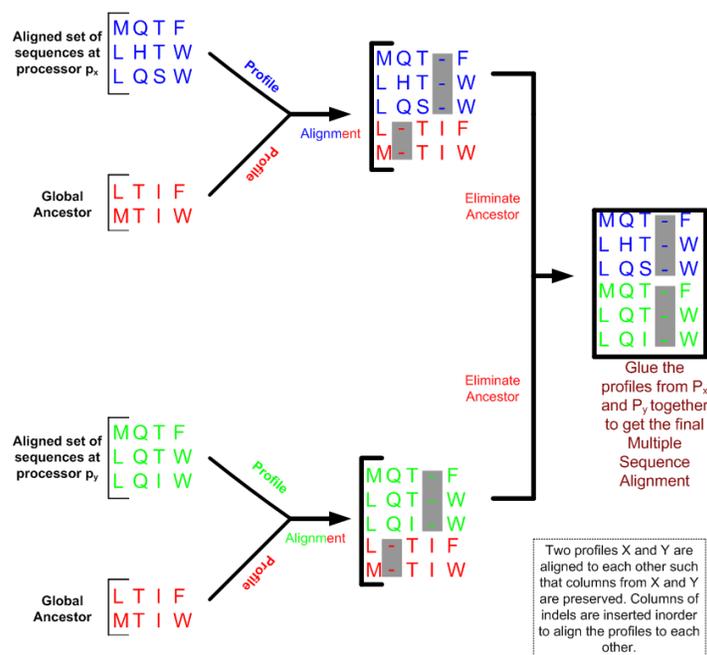


Figure 2: Profile aligning with the ancestor and combining sequence subsets.

The summary of different steps in the Sample-Align-D Algorithm is shown in Fig. 3, and a detailed algorithmic description is given in the Appendix as Algorithm 2.
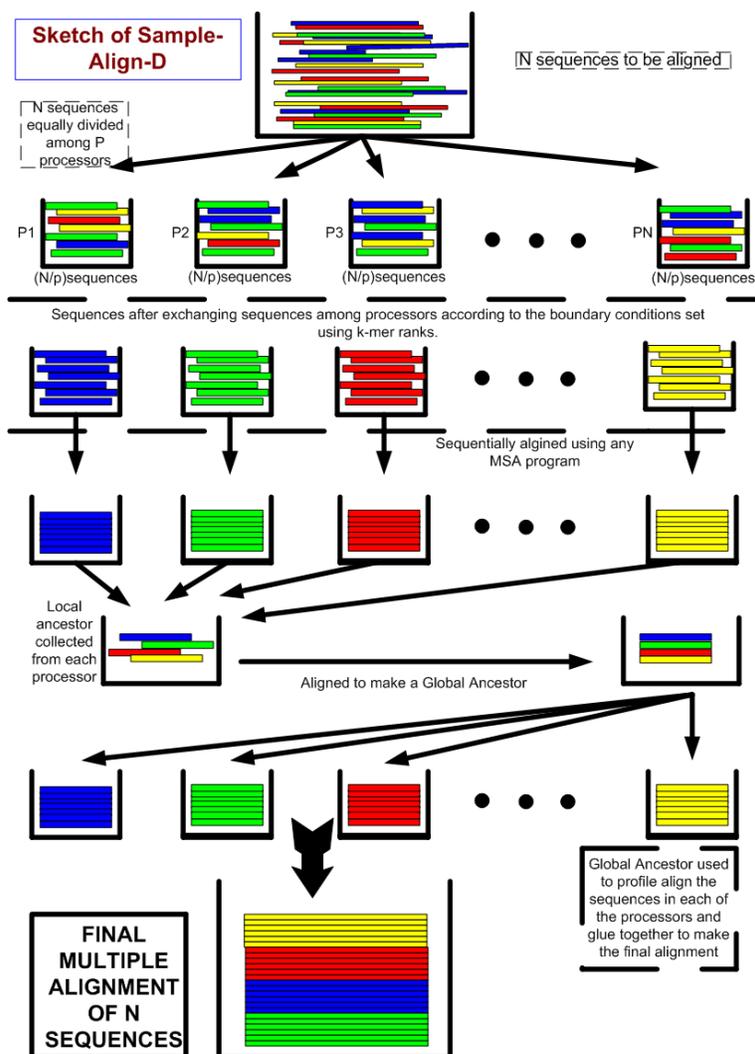
Figure 3: Summary of Sample-Align-D procedure

## 3.2 Analysis of Computation and Communication Costs

For the computation and communication analysis we use a coarse grained computing model such as $C^3$-model [36] and [37]. Also, for analysis purposes, we assume that the Muscle System [3] is being used at each processor as the underlying sequential multiple sequence alignment system. It must be noted that the computation complexity of the alignment step will vary depending on the sequential MSA system used for alignment within each processor.

In the following analysis we assume that each processor has $w = N/p$ sequences, where $N$ is the total number of sequences to be aligned, and $p$ is the number of processors. The average length

Table 2: Computation Costs:

| STEP | O(Time) | O(Space) |
|---|---|---|
| k-mer rank computation on ($w = N/p$) sequences | $w^2 L$ | $w + L$ |
| Sorting of $N/p$ sequences based on k-mer rank | $w \log w$ | $\log w$ |
| Sample $k = p - 1$ sequences | $w$ | $p$ |
| k-mer rank computation of ($k \times p$) sequences in root processor | $p^4 L$ | $p^2 + L$ |
| Sorting of $k \times p$ sample k-mer ranks | $(k \times p) \log(k \times p)$ | $\log(k \times p)$ |
| k-mer rank computation of each of ($w = \frac{N}{p}$) sequences against $k \times p$ samples | $w[(k \times p + 1)^2 L]$ | $w(k \times p + L)$ |
| Muscle executed on ($w = \frac{N}{p}$) sequences in parallel | $w^4 + wL^2$ | $w^2 + L^2$ |
| Ancestor extraction from each of the $p$ processors + export to the root processor | $p^2$ | $p^2$ |
| Muscle executed on local ancestors ($p$ elements) | $(p)^4 + (p)L^2$ | $(p)^2 + L^2$ |
| Profile alignment with all combined aligned sequences on each of the processor | $wL^2$ | $w$ |
| TOTAL Computation Cost (for $w = \frac{N}{p}$) | $O((\frac{N}{p})^4 + (\frac{N}{p})L^2)$ | $O((\frac{N}{p})^2 + L^2)$ |

of a sequence is $L$. In Table 2, we outline the computation cost of each step of the algorithm and its memory requirement.

## 3.3 Communication Cost

The communication overhead is an important factor that dictates the performance of a distributed message passing parallel system. If the communication overhead is much higher than the computation cost, the performance of the system is limited. Fortunately, the communication cost of our system is much less than the cost of the alignment. Essentially, the proposed Sample-Align-D algorithm has two rounds of communication. In the first round, a small set of samples is collected at the root processor and a set of pivots is broadcast from the root processor. In the second round, sequences are redistributed to achieve better alignments and balanced load distribution. For the analysis of the communication costs we have adopted the coarse grained computation model [36] and [22]. However, we ignore the message start up costs and assume unit time to transmit each data byte.

We have assumed the Regular Sampling strategy [24] because of its suitability to our problem domain. Some of the reasons are :

1. The strategy is independent of the distribution of original data, compared to some other

strategies such as Huang and Chow [23].

2. It helps in partitioning of data into ordered subsets of approx. equal size. This presents an efficient strategy for load balancing as unequal number of sequences on different processors would mean unequal computation load, leading to poor performance. In the presence of data skew, regular sampling guarantees that no processor computes more than $(2\frac{N}{p})$ sequences [24].

3. It has been shown in [24] that regular sampling yields optimal partitioning results as long as $N > p^3$, i.e., the number of data items $N$ is much larger than the number of processors $p$, which would be a normal case in the MSA application.

### 3.3.1   First Communication Round

Assuming $k = p - 1$, i.e., each processor chooses $p - 1$ samples, the complexity of the first phase is $O(p^2L)$+ $O(p \log p) + O(k \times p \log p)$, where $O(p^2L)$ is the time to collect $p(p - 1)$ samples of average length $L$ at the root processor, $O(p \log p)$ is the time required to broadcast $p - 1$ pivots to all the processor and $(k \times p \log p)$ is the time required to broadcast $k \times p$ sequences to all the processors.

### 3.3.2   Second Communication Round

In the second round each processor sends the sequences having k-mer rank in the range of bucket $i$ to processor $i$. Each processor partitions its block into $p$ sub-blocks, one for each processor, using pivots as bucket boundaries. Each processor then sends the sub-blocks to the appropriate processor. The sizes of these sub-blocks can vary from 0 to $\frac{N}{p}$ sequences depending on the initial data distribution. Taking the average case where the elements in the processor are distributed uniformly, each sub-block will have $\frac{N}{p^2}$ sequences. Thus this step would require $O(\frac{N}{p})$ time assuming an all-to-all personalized broadcast communication primitive [37]. However, in the following we show that based on regular sampling no processor will receive more than $2\frac{N}{p}$ elements in total in

the worst case. Therefore still the overall communication cost will be $O(\frac{N}{pL})$.

Let's denote the pivots chosen in the first phase by the array: $y_1, y_2, y_3, \cdots, y_{p-1}$. Consider any processors $i$, where $1 < i < p$. All the sequences to be processed by processor $i$ must have k-mer rank $> y_{i-1}$ and $\leq y_i$. There are $(i-2)p + \frac{p}{2}$ sequences of the regular sample which are $\leq y_{i-1}$, implying that there are at least $lb = ((i-2)p + \frac{p}{2})\frac{N}{p^2}$ sequences in the entire data that have k-mer rank $\leq y_{i-1}$. On the other hand, there are $(p-i)p - \frac{p}{2}$ sequences in the regular sample that have k-mer rank $> y_i$. Thus, there are $ub = ((p-i)p - \frac{p}{2})\frac{N}{p^2}$ sequences of $N$ which are $> y_i$. Since the total number of sequences is $N$, at most $N - ub - lb$ sequences will get assigned to processor $i$. It is easy to show that this expression is upper bounded by $2\frac{N}{p}$. The cases for $i = 1$ and $i = p$ are special because the pivot interval for these two processors is $\frac{p}{2}$. The load for these processors will always be less than $2\frac{N}{p}$. Due to page limitations, we refer to [24] for further details of the analysis.

The collection of $p$ local ancestors at the root processor and broadcast of the global ancestor costs $O(L \log p)$ communication overhead each. Therefore the communication cost is: $O(p^2L) + O(p \log p) + O(\frac{N}{pL}) + O(L \log p)$.

The total asymptotic time complexity $T$ of the algorithm would be:

$$CommutationCosts = O(\frac{N}{p})^4 + O(\frac{N}{p})L^2 \tag{5}$$

$$CommunicationCosts = O(p^2L) + O(p \log p) + O(\frac{N}{pL}) + O(L \log p) + O(k \times p \log p) \tag{6}$$

$$T \approx O((\frac{N}{p})^4 + (\frac{N}{p})L^2) + (p^2L) + (\frac{N}{pL}) \tag{7}$$

Next we briefly comment on the scalability of the Sample-Align-D algorithm. We will use the *isoefficiency* metric [22] to show that Sample-Align-D is highly scalable. For this we first define two important terms: problem size defined as the number of basic computation steps to solve a problem on a single processor using the best sequential algorithm; overhead function is defined as the cost, that is not incurred by the fastest known sequential algorithm. We denote problem size with $W$, and overhead function with $T_o(W, p)$.

$$W = \frac{E}{E-1} \times T_o(W, p) \tag{8}$$

where $E$ denotes the efficiency and let $K = \frac{E}{E-1}$. The overhead function of Sample-Align-D:

$$T_o(W, p) \approx (\frac{N}{p}) \times L + p^2 \times L \tag{9}$$

It is easy to show that asymptotically the iso-efficiency of Sample-Align-D is $\Theta(p^2)$, i.e., the number of sequences shall increase by a factor of $p^2$ to maintain the efficiency with increasing number of processors.

## 4   Performance Evaluation

The performance evaluation process has been divided into two parts: the first part deals with the quality assessment, and the second part deals with traditional HPC metrics such as execution time, scalability, memory requirements, etc. The performance evaluation of the Sample-Align-D algorithm is carried out on a Beowulf Cluster consisting of 16 Intel Xeon processors, each running at 2.40GHz, with 512KB cache and 1GB DRAM memory. As for the interconnection network, the system uses Intel Gigabit network interface cards on each cluster node. The operating system on each node is Fedora Core 7(kernel level:2.6.18-1.2798.fc6xen).

### 4.1   Quality Assessment

The quality assessment in our case posed a considerable challenge because most of the existing benchmarks such as BaliBase [38] and Prefab [3] used in the literature are of very small sizes. Therefore they are not effective in evaluating any *sampling* based approach or domain decomposition based distributed approach. Also, other parallel approaches to multiple alignment do not have any decomposition strategy, making the quality of the parallel version similar to the sequential version. Hence, a quality assessment criterion for data parallel multiple alignment methods was

not available. In addition to assessing quality using these traditional benchmarks, we have also formulated a method that can be used to access the quality of the alignment produced by distributed or data parallel MSA approaches.

### 4.1.1   Quality Assessment using Traditional Benchmarks

Traditional benchmarks such as BaliBase and Prefab are quite comprehensive in terms of types of sequences contained in these benchmarks. BaliBase, for example, has five basic categories and covers most of the scenarios when making multiple sequence alignments [2]. For the evaluation of multiple sequence alignment programs, Balibase is divided into 5 hierarchical reference sets:

- Ref1 for equi-distant sequences with various levels of conservation,

- Ref2 for families aligned with a highly divergent "orphan" sequence,

- Ref3 for subgroups with $< 25\%$ residue identity between groups,

- Ref4 for sequences with N/C-terminal extensions, and

- Ref5 for internal insertions.

Tables 3, 4 and 5 compare the quality of Sample-Align-D with different sequential algorithms in terms of quality metrics, Q-Score, and TC-Score, used in Balibase and Prefab benchmarks, respectively. The score for the sequential algorithms have been derived from [3]. The Sample-Align-D was executed on a 4-processor system, therefore corresponds to a 4 factor domain decomposition.

In all the tests for quality assessment using benchmarks, it can be seen that Sample-Align-D preformed very close to the Muscle system. This is because Sample-Align-D was implemented with Muscle System as the underlying sequential MSA algorithm at each processor. Therefore, the quality obtained is limited by the quality of the underlying alignment system.

Table 3: BaliBase Scores

| Method | Q | TC |
|--------|-----|-----|
| Muscle | 0.896 | 0.747 |
| T-Coffee | 0.882 | 0.731 |
| NWNSI(MAFFT) | 0.881 | 0.722 |
| Clustalw | 0.860 | 0.690 |
| **Sample-Align-D** | **0.858** | **0.720** |
| FFTNSI(MAFFT) | 0.844 | 0.646 |

Table 4: Prefab Q-Scores

| Method | Q-Score(All) |
|--------|--------------|
| Muscle | 0.645 |
| **Sample-Align-D** | **0.623** |
| T-Coffee | 0.615 |
| NWNSI(MAFFT) | 0.615 |
| FFTNSI(MAFFT) | 0.591 |
| Clustalw | 0.563 |

Table 5: BaliBase Q-Scores on subsets

| Method | Ref1 | Ref2 | Ref3 | Ref4 | Ref 5 |
|--------|------|------|------|------|-------|
| **Sample-Align-D** | **0.882** | **0.932** | **0.800** | **0.872** | **0.804** |
| Muscle | 0.887 | 0.935 | 0.823 | 0.876 | 0.968 |
| T-Coffee | 0.866 | 0.934 | 0.787 | 0.917 | 0.957 |
| NWNSI | 0.867 | 0.923 | 0.787 | 0.904 | 0.963 |
| Clustalw | 0.861 | 0.932 | 0.751 | 0.823 | 0.859 |
| FFTNSI | 0.838 | 0.908 | 0.708 | 0.793 | 0.947 |

### 4.1.2   Quality Assessment with Increasing Degree of Decomposition

While the quality of alignment produced by the Sample-Align-D algorithm for benchmark data sets is comparable to the other systems, due to the small sizes of these benchmarks, we could not evaluate the quality against several desired parameters that are typical of a distributed environment such as sample size, number of partitions, average length of sequences, etc. In the following, we first outline the assessment procedure that allows us to evaluate the quality while changing different parameters, and then present performance results. In this subsection, we compare only with the Muscle System.



Figure 4: Quality comparison, with varying average length of the sequences.

Using the Rose sequences generator [39] we have generated 23 sets of sequences, with their corresponding 'true' alignment, while changing the following three parameters: the length of the sequences, the number of sequences, and the phylogenetic distance of the sequences. The length of the sequences in our tests varied from 100 to 2000, the number of sequences varied from 100 to 20000 and the average phylogenetic distance varied from 100 to 1000. These values are typical of biological sequences in existing databases. Each of the set was aligned using Sample-Align-D with different number of processors, and 'true' alignment obtained from the Rose system was used as a benchmark. The metrics used for the the assessment are Q-Score [3], TC-Score [38], Modeler score [40], Cline (Shift) score [41], and Sum of Pairs (SP) score. The quality scores obtained by

the Sample-Align-D Algorithm are compared with the Muscle System scores (only Q-scores are reported in this paper).

The experiments were conducted with different data sets while keeping two of the parameters held constants[1]. These parameters include number of sequences, phylogenetic distance, and average sequence length. The experiments were also conducted on different machine sizes to study scalability and the effect of degree of domain decomposition on quality of alignment.

Fig. 4 depicts performance in terms of Q-score with increasing number of processors, while increasing average length of sequences from 200 to 2000. It can be seen that the increase in the number of processors to 16 didn't effect the Q-scores. The Q-scores correlated very well with that of the Muscle System. The scores remained above 0.97 for average length of 2000. There was virtually no difference observed in the SP scores computed for the Muscle system and Sample-Align-D for respective pair of length and number of processors used, as shown in Fig. 5.
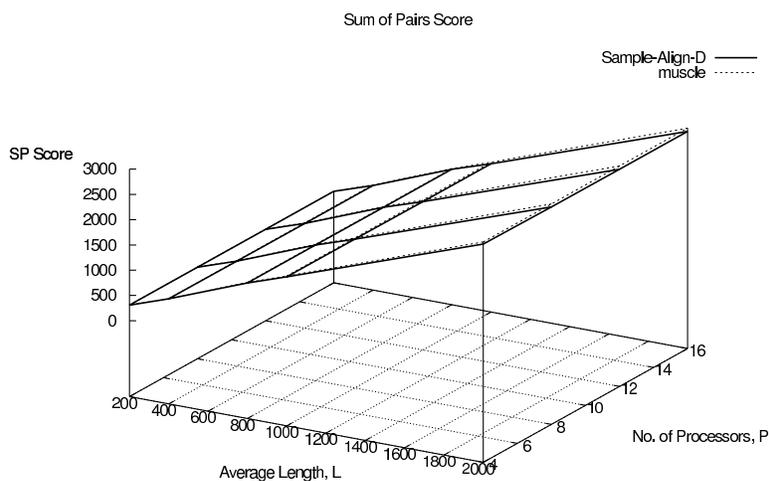


Figure 5: Sum of Pairs (SP) score, with varying average length of the sequences.

The quality with respect to the phylogentic distance is probably the most important criterion. Theoretically, MSA systems should be able to give good multiple alignments for increasing pairwise sequence distance. After all, the relationship between distance species would reveal the phylogenetics of the species. As pointed out in [42], all automatic multiple alignment systems perform

---

[1]Constants held for the experiments are: Length =200, Number of sequence=200 and phylogenetic distance =100

poorly with increasing pair wise distance between the sequences. Different quality metrics calculated while increasing the phylogentic distance from 100 to 1000 are shown in Fig. 6 and Fig. 7 for Q-Scores and TC-Scores, respectively. As can be seen from these figures, the quality in fact decreased with increasing pairwise distance. Our investigation for this criteria, however, was not to rectify the quality issues with increasing pairwise distance, but to see the correlation between the underlying MSA system and the effects of the domain decomposition strategy.



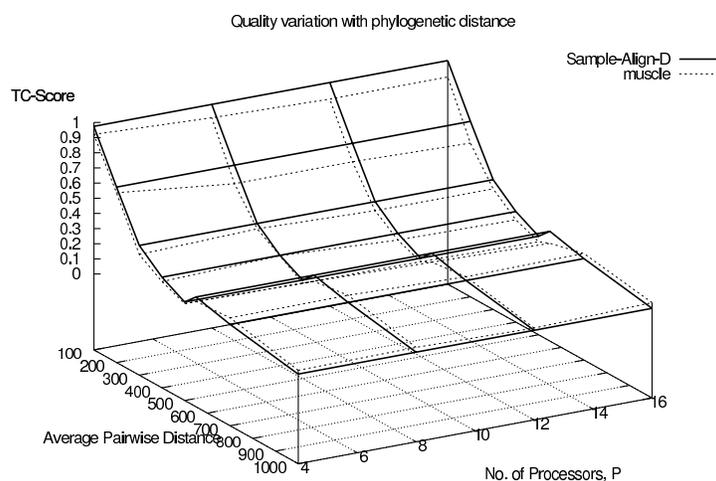Figure 6: Q-score, with varying average pairwise distance of sequences.



Figure 7: TC-score, with varying average pairwise distance of sequences.

As can be seen from Fig. 6, the Q-scores correlated extremely well with the Q-scores of the
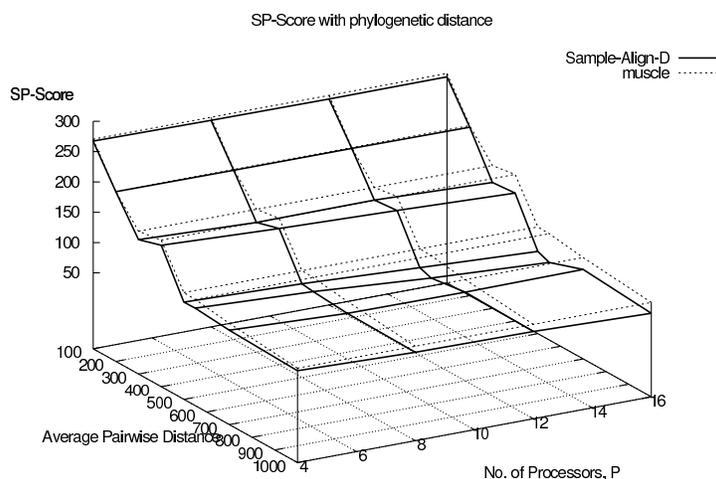
Figure 8: SP Score, with varying average pairwise distance of sequences.

Muscle System. The TC-scores of the Sample-Align-D depicted in Fig. 7 show slight *increase* in the quality when compared with TC scores of the Muscle System. The increase can be attributed to the decomposition strategy which is more inclined towards conserving columns in the multiple alignment, owing to profile alignments. The decrease in the alignment scores in general with increasing phylogenetic distance, is due to the decrease of the alignment quality obtained from the underlying MSA system. Without loss of generality, the quality of a decomposition based MSA can be expected to correlate well with the underlying sequential multiple alignment system that may perform well in terms of quality. Fig. 8 shows the quality performance in terms of SP score.

The assessment of the quality of alignment with increasing number of sequences is also important. As before, we are interested in the relative quality of alignment obtained after decomposition. Fig. 9 shows the quality in terms of Q-Scores for different sizes of the sequences set. The quality of the Sample-Align-D Algorithm based alignments strongly correlates with the quality of the alignments obtained by the Muscle system. It must be noted that we are reporting quality for up to 8000 sequences, because the sequential Muscle System was unable to process larger datasets.
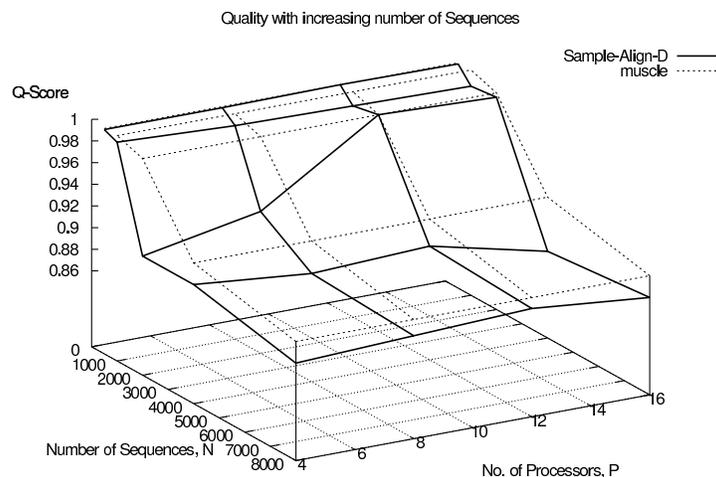
Figure 9: Quality of alignment, with increasing number of sequences.

### 4.1.3 Example of Application to Serine/Threonine Kinases

The purpose of a multiple sequence alignment system is to observe and study the conservation of domain and motifs. To illustrate this, we present here an example Fig. **??** that illustrates the usefulness of our system, in terms of conservation of motifs. We illustrate here the functional features of kinases, also present in the BaliBase as well as used for illustration by Notredame [2]. Each sequence is identified by its SwissProt/UniProt identifier. Some of the identifications have changed, and are illustrated as is viewed in SwissProt at the time of this publication. In the example, there are 3 motifs identified. These motifs are the core blocks identified by BaliBase, and are conserved by Sample-Align-D, marked as red, orange and blue. The motifs are in greater order of difficulty, with red as the least difficult. The blue labeled motif is the most difficult to conserve in the example set, because of the long indel in KIN3-Yeast. As can be seen, Sample-Align-D is able to conserve this most difficult motif, for decomposition factor of 4 as done for benchmarks.

## 4.2 Performance in Terms of HPC Parameters

In this section we analyze performance in terms of execution time, scalability, and memory. The objective of the evaluation is to determine the advantages of the proposed domain decomposition based technique in terms of speedup and reduction in memory requirements. In this section we

also compare Sample-Align-D with the known parallel approaches reported in the literature.

### 4.2.1  Execution Time and Memory

For the sake of coherence in our presentation we generated sequences with same parameters that were used for the quality assessment. We report results for up to 20000 sequences. To the best of authors' knowledge, there are no published reports of aligning this large number of sequences in the literature.
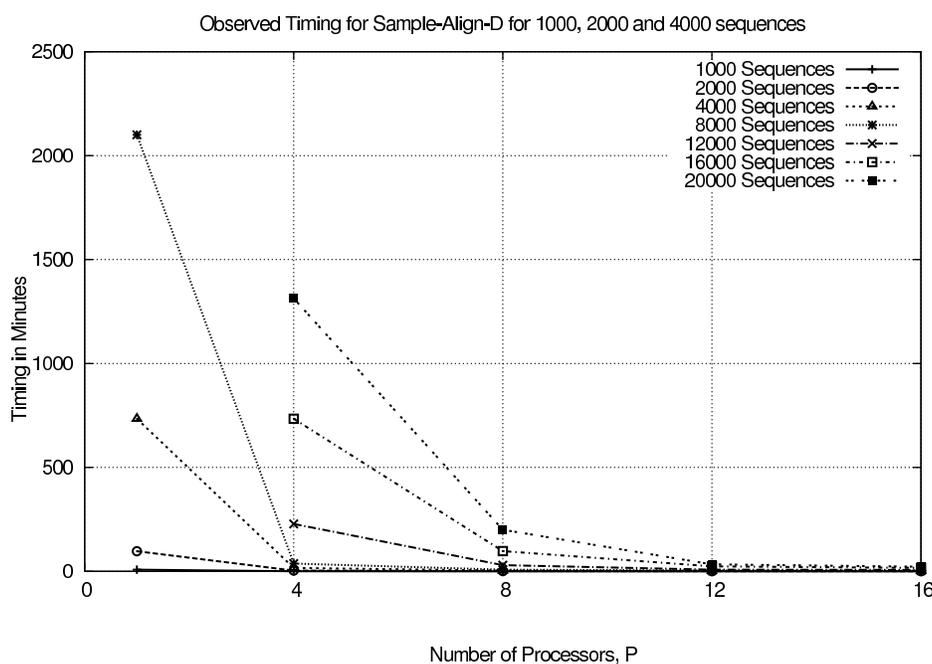


Figure 10: Scalability of the execution time w.r.t. the number of processors.

As shown in Fig. 10, in the case of Sample-Align-D the execution time decreases sharply with the increase in the number of processors. We are able to align 8000 sequences in just 3.9 minutes, compared to 2100 minutes on a sequential Muscle System. The timing for 12000, 16000 and 20000 sequences are also shown for Sample-Align-D. The timing for one node with Muscle is not shown because Muscle Systems was not able to handle this large number of sequences and the resources requirement in terms of memory and time, grew exponentially for these sets of sequences.

As shown in Fig. 11, Sample-Align-D Algorithm exhibits super linear speed-ups (of the order
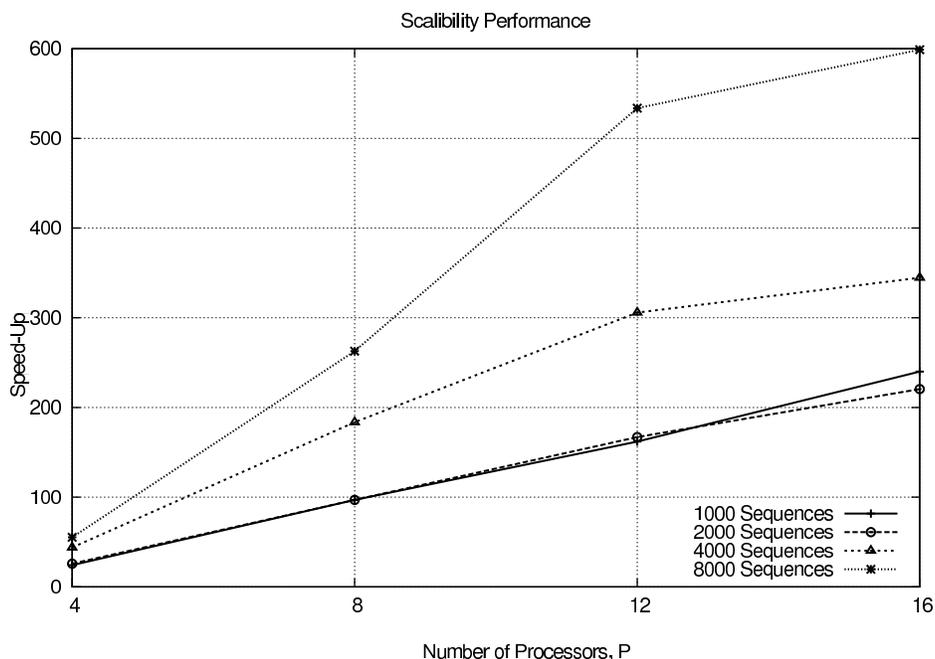
Figure 11: Super-linear speed-ups for Sample-Align-D with increasing number of processors.

of 600) on a 16 processor system. This is primarily because the computation complexity decreases by $O(p^4)$ with the increase in the number of processor, as suggested in our algorithmic analysis section.

The low memory requirements, as predicted by our analysis in Section 3.2, are also evident in our experimental results. The memory requirements while increasing the length of sequences, phylogenetic distance, and number of sequences, are shown in Fig. 12 Fig 13 and Fig. 14 respectively. The most interesting figure is the one that depicts the memory requirements with increasing number of sequences. As can be seen in Fig. 14, with the increase in the number of sequences, the memory requirements for Muscle System is increasing exponentially with 1200 MB required for 8000 sequences. However, the maximum memory required for Sample-Align-D even for 8000 sequences is not greater than 100 MB.
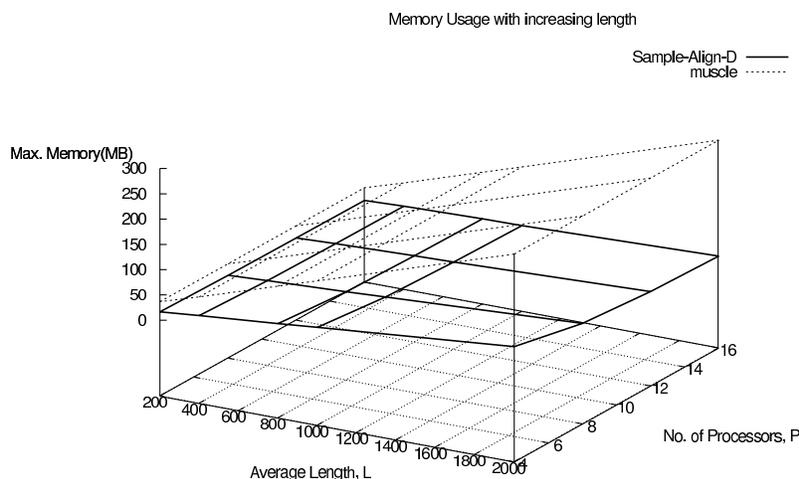
Memory Usage with increasing length

Sample-Align-D ———
muscle ........

Figure 12: Memory usage, with varying average length of the sequences.

Memory Usage with phylogenetic distance
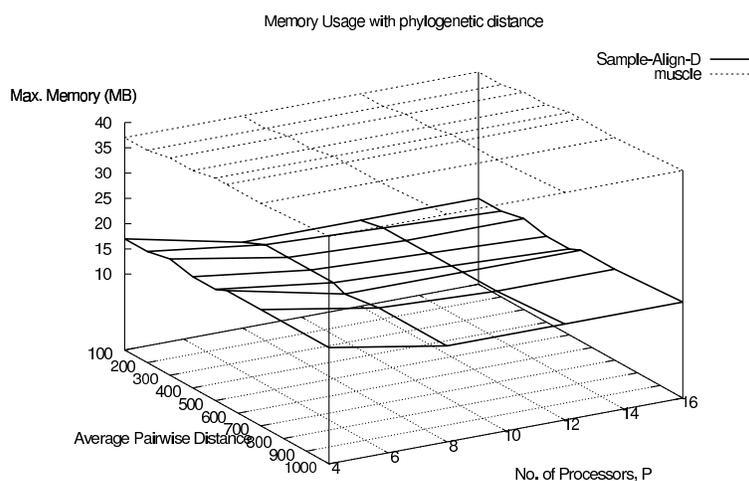
Sample-Align-D ———
muscle ........

Figure 13: Memory usage, with varying average pairwise distance of sequences.

### 4.2.2   Comparison with Existing Parallel MSA Systems

There have been significant efforts towards parallelizing MSA techniques, as discussed in Section 3. We have selected Parallel Clustalw and Parallel T-Coffee, two of the most widely used parallel MSA systems, to compare the performance of the proposed Sample-Align-D Algorithm. The limitation of this aspect of evaluation was also the selection of the common sequence set. This is due to the fact that most of the existing parallel systems are unable to handle large number of sequences because of one or two explicit sequential stages in these solutions.
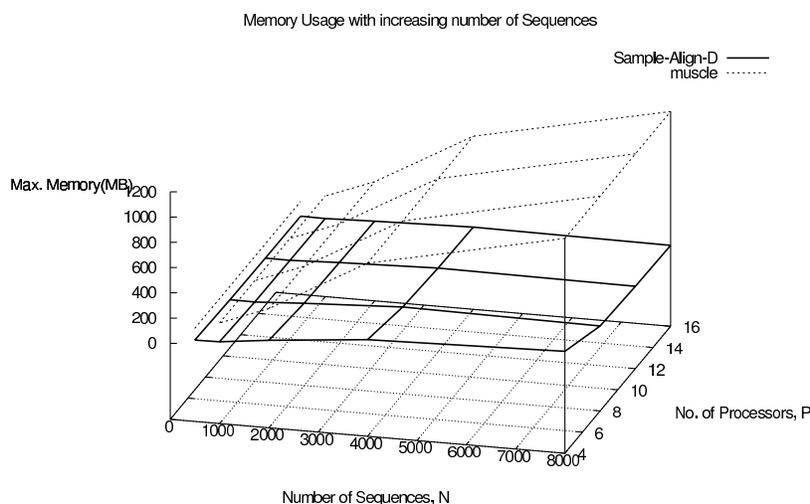
Figure 14: Memory usage, with varying number of sequences.

We chose the data sets that Zola et al. [15] used for the evaluation of Parallel T-Coffee, named PF00500, consisting of 1048 sequences with maximal length of 523 characters. The execution times of different parallel algorithms for the data sets are plotted in Fig. 15. The execution time of Parallel T-Coffee is significantly higher than that of the Sample-Align-D algorithm. For example, on a 16 processor system, it took around 9.1 hours for Parallel T-Coffee, 4.3 minutes for Parallel Clustalw, and only 8.1 seconds for Sample-Align-D. Our experiments show that the performance of Parallel Clustalw degrades significantly compared to Sample-Align-D as the number of sequences in the set increases.

# 5   Conclusion and Discussions

We have described a domain decomposition (data parallel) strategy for multiple sequence alignment of biological sequences. To our knowledge, this is the first attempt to investigate domain decomposition for the multiple sequences alignment problem. This domain decomposition allowed us to devise a highly scalable multiple alignment system. A detailed algorithmic technique based on a novel decomposition strategy was described and rigorous time and space complexity analyses were presented. The proposed strategy decreased the time complexity of any MSA heuristic
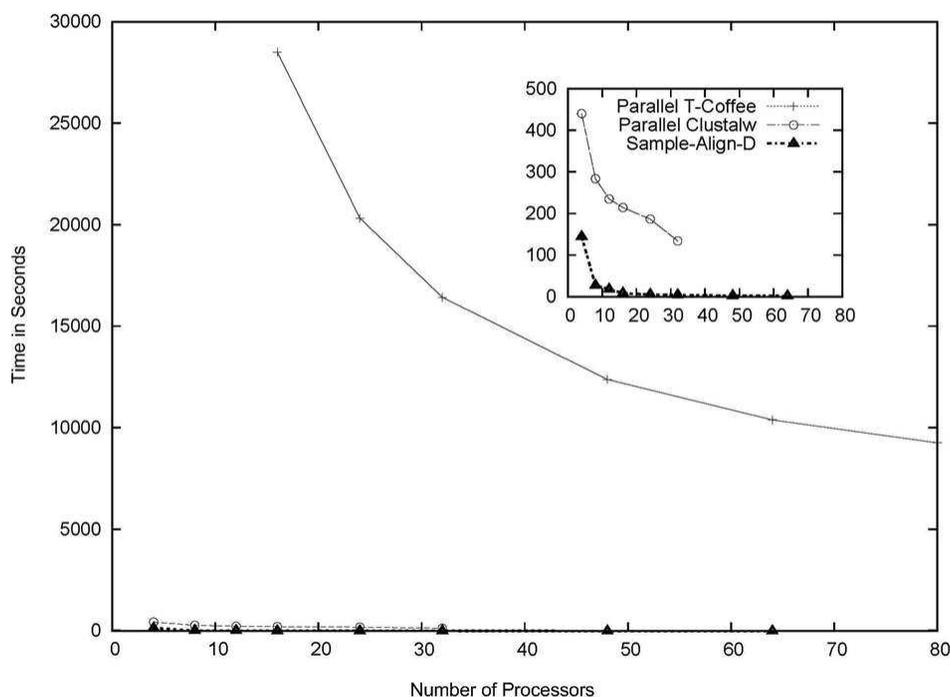
Figure 15: Comparison of execution times of parallel T-Coffee, parallel Clustalw and Sample-Align-D

$O(N)^x$ by a factor of $O(1/p)^x$. Consequently, super-linear speed-ups were achieved and tremendous decrease in memory requirements were observed, as predicted by the complexity analysis. A rigorous quality analysis of the decomposition technique was also introduced, and the effect of decomposition on the quality of the alignment was investigated. The quality analysis allowed us to determine the quality of the alignment relative to that of the underlying sequential MSA system.

A number of research problems remain open and the techniques introduced in this manuscript suggest new directions of research that can be pursued. The open research problems in computational biology and parallel processing that arise from the research presented in this manuscript are as follows:

1. We have presented the decomposition strategy as a parallel computing solution. However, the super-linear speed-ups on multiple processors suggest that the use of the sampling based decomposition strategy on single processor systems would also be able to deliver significant

time and space advantages as evident in [43]. The decomposition strategy in any of the exist-
ing multiple alignment systems with some form of iterative quality 'patch-up' strategy would
be useful. Also, multi-core processors can take enormous advantage of the decomposition
strategy for improved efficiency.

2. In our partitioning strategy, the load balancing scheme is based on the k-mer rank and the
   number of sequences. It would be interesting to develop a more elaborated load balancing
   scheme considering additional factors such as the length of the sequences and the type of
   sequences being considered.

3. The domain decomposition based strategy has been investigated for distance based multiple
   sequences alignment methods such as Muscle and Clustalw. It would be interesting to in-
   vestigate the same or similar partitioning strategy for other type of consistency and profile
   based methods such as T-Coffee, ProbCons, Mafft, DbClustal, MUMMALS etc.

4. We have considered a subset of alignment parameters, for example, PSP scores, 200 PAM
   matrix and the 240 PAM VTML matrix. It would be insightful to consider different mutation
   matrix and other parameters, and investigate the effects of decomposition on quality.

5. Phylogenetic trees are the crux of research on evolutionary biology. However, building phy-
   logenetic trees for large number of species is considerably compute intensive. It would
   be useful to apply decomposition to build distance based phylogenetic trees for multiple
   genomes.

6. As the sizes of the biological sequence archives and structural data are increasing at an expo-
   nential rate, the pattern and motifs searching is getting increasingly more time consuming.
   The application of decomposition strategy could allow to search these keys motifs in such
   databases. The strategy can also find its applications in target and lead identification in drug
   discovery.

7. The strategy can be used to obtain computational analysis without having the need to actually 'import' the sequences locally e.g. multiple sequences alignment can be performed on distant databases, without transferring the entire set of sequences, which some times might be desirable due to proprietary data issues etc.

## APPENDIX

---

**Algorithm 2** Sample-Align-D (sequences $N$)

---

**Require:** $p$ processors for computation and $N$ sequences of amino acids $S_1, S_2, \cdots, S_N$:

**Ensure:** Gaps are inserted in each sequence such that:

- All sequences have the same length and the Score of the global map is maximized according to the chosen scoring function

1. Assume $N/p$ sequences on each of the $p$ processors

2. Locally compute k-mer rank of all the sequences in each processor

3. Sort the sequences locally in each processor based on k-mer rank

4. Choose a sample set of $k$ sequences in each processor, where $k \ll N/p$

5. Send the $k$ samples from each processor to all the processors.

6. Compute the k-mer rank of each sequence against the $k \times p$ samples.

7. Sort the sequences locally in each processor based on the new k-mer rank.

8. Using regular sampling, choose $p - 1$ sequences from each processor and send only their ranks to a root processor.

9. Sort all the $p \times (p - 1)$ ranks at the root processors and divide the range of ranks into $p$ buckets.

10. Send the bucket boundaries to all the processors.

11. Redistributed sequences among processors such that sequences with k-mer rank in the range of bucket $i$ are accumulated at processor $i$, where $0 > i < p + 1$.

12. Align sequences in each processor using any sequential multiple alignment system

13. Broadcast the Local Ancestor to the root processor

14. Determine Global Ancestor GA at the root processor by aligning local ancestors received from all the processors

15. Broadcast GA to all the processors

16. Realign each of the sequences in $p$ processors based on ancestor GA using profile-profile alignment i.e. Each of the profiles of aligned sequences are tweaked using the ancestor profile, with constraints.

17. Glue all the aligned sequences at the root processor.

---

# References

[1] L. Wang and T. Jiang, "On the complexity of multiple sequence alignment," *Journal of Computational Biology*, vol. 1, no. 4, pp. 337–348, 1994.

[2] C. Notredame, D. Higgins, and J. Heringa, "T-coffee: A novel method for multiple sequence alignments," *Journal of Molecular Biology*, vol. 302, pp. 205–217, 2000.

[3] R. C. Edgar, "Muscle: Multiple sequence alignment with high accuracy and high throughput," *Nucleic Acids Research*, vol. Vol. 32, p. No.5, 2004.

[4] J. Thompson, D. Higgins, T. J., and Gibson, "Clustal w: improving the sensitivity of progressive multiple alignment through sequence weighting, position-specific gap penalties and weight matrix choice," *Nucleic Acids Res*, vol. 222, pp. 4673–4690, 1994.

[5] C. Do, M. Mahabhashyam, M. Brudno, and S. Batzoglou, "Probcons: Probabilistic consistency-based multiple sequence alignment," *Genome Research*, vol. 15, pp. 330–340, 2005.

[6] T. Lassmann and E. Sonnhammer, "Kalign - an accurate and fast multiple sequence alignment algorithm," *BMC Bioinformatics*, vol. 6, no. 1, p. 298, 2005.

[7] S.-H. Sze, Y. Lu, and Q. Yang, "A polynomial time solvable formulation of multiple sequence alignment," *Journal of Computational Biology*, vol. 13, no. 2, pp. 309–319, 2006. PMID: 16597242.

[8] A. S. Schwartz and L. Pachter, "Multiple alignment by sequence annealing," *Bioinformatics*, vol. 23, no. 2, pp. e24–29, 2007.

[9] B. Morgenstern, "Dialign: multiple dna and protein sequence alignment at bibiserv," *Nucleic Acids Research*, vol. 32, pp. W33–W36, 2004.

[10] M. Ronaghi, "Pyrosequencing sheds light on dna sequencing," *Genome Res.*, vol. 11, no. 1, pp. 3–11, 2001.

[11] F. Saeed and A. Khokhar, "Sample-Align-D: A high performance multiple sequence alignment system using phylogenetic sampling and domain decomposition," *Proc. 23rd IEEE International Parallel and Distributed Processing Symposium*, April 2007.

[12] C. Notredame, "Recent evolutions of multiple sequence alignment algorithms," *PLoS Computational Biology*, vol. 3, p. e123, Aug 2007.

[13] J. Cheetham, F. Dehne, A. Rau-chaplin, and P. J. Taillon, "Parallel clustal w for pc clusters," in *Proceedings of International Conference on Computational Science and Its Applications (ICCSA 2003)*, pp. 300–309, 2003.

[14] Mikhailov.D., Cofer.H., and Gomperts.R., "Performance optimization of clustalw: Parallel clustalw, ht clustal, and multiclustal," *Silicon Graphics. CA.*, vol. White papers, 2001.

[15] J. Zola, X. Yang, A. Rospondek, and S. Aluru, "Parallel T-Coffee: A parallel multiple sequence aligner," in *ISCA PDCS*, pp. 248–253, 2007.

[16] X. Deng, E. Li, J. Shan, and W. Chen, "Parallel implementation and performance characterization of muscle," *International Parallel and Distributed Processing Symposium (IPDPS)*, 2006.

[17] J. Zola, D. Trystram, A. Tchernykh, and C. Brizuela, "Parallel multiple sequence alignment with local phylogeny search by simulated annealing," *Sixth IEEE International Workshop on High Performance Computational Biology*, March 2007.

[18] K.Chaichoompu, S. Kittitornkun, and S. Tongsima, "Mt-clustal-w: Multithreading multiple sequence alignment," *Sixth IEEE International Workshop on High Performance Computational Biology*, March 2007.

[19] M. Schmollinger, K. Nieselt, M. Kaufmann, and B. Morgenstern, "Dialign p: Fast pair-wise and multiple sequence alignment using parallel processors," *BMC Bioinformatics*, vol. 5, p. 128, September 2004.

[20] T. F. Smith and M. S. Waterman, "Identification of common molecular subsequences," *Journal of Molecular Biology*, vol. 147, pp. 195–197, 1981.

[21] F. Zhang, X. Qiao, and Z. Liu, "A parallel smith-waterman algorithm based on divide and conquer," *ica3pp*, vol. 00, p. 0162, 2002.

[22] V. Kumar, A. Grama, A. Gupta, and G. Karpis, *Introduction to Parallel Computing: Design and Analysis of Parallel Algorithms*. Benjamin-Cummings Pub Co, January 1994.

[23] J. S. Huang and Y. C. Chow, "Parallel sorting and data partitioning by sampling," in *Proceedings of the IEEE Computer Society's Seventh International Computer Software and Applications Conference*, pp. 627–631, 1983.

[24] S. Hanmao and S. Jonathan, "Parallel Sorting by Regular Sampling," *Journal of Parallel and Distributed Computing*, vol. 14, no. 4, pp. 361–372, 1992.

[25] R. C. Edgar, "Local homology recognition and distance measures in linear time using compressed amino acid alphabets," *Nucleic Acids Research*, vol. 32, pp. 380–385, 2004.

[26] M. J. Berger and S. H. Bokhari, "A partitioning strategy for nonuniform problems on multiprocessors," *IEEE Trans. Comput.*, vol. 36, no. 5, pp. 570–580, 1987.

[27] S. H. Bokhari, T. W. Crockett, and D. M. Nicol, "Binary dissection: Variants and applications," Tech. Rep. 97-29, ICASE, 1997.

[28] M. Kaddoura, C.-W. Ou, and S. Ranka, "Partitioning unstructured computational graphs for nonuniform and adaptive environments," *IEEE Parallel and Distributed Technology*, vol. 03, no. 3, pp. 63–69, 1995.

[29] J. R. Pilkington and S. B. Baden, "Dynamic partitioning of non-uniform structured workloads with spacefilling curves:," *IEEE Transactions on Parallel and Distributed Systems*, vol. 7, no. 3, pp. 288–300 (or 288–299??), 1996.

[30] P. E. Crandall and M. J. Quinn, "Non-uniform 2–D grid partitioning for heterogeneous parallel architectures," in *IEEE Symposium on Parallel and Distributed Processing*, (Los Alamitos, CA), pp. 428–435, IEEE, 1995.

[31] M. Willebeek-LeMair and A. Reeves, "Strategies for dynamic load balancing on highly parallel computers," *Parallel and Distributed Systems, IEEE Transactions on*, vol. 4, pp. 979–993, Sep 1993.

[32] R. C. Edgar, "Muscle: A multiple sequence alignment method with reduced time and space complexity," *BMC Bioinformatics*, pp. 1471–2105, 2004.

[33] A. SF, "Amino acid substitution matrices from an information theoretic prospective," *J Mol Biol*, vol. 219(3), pp. 555–565, 1991.

[34] D. T. Jones, W. R. Taylor, and J. M. Thornton, "The rapid generation of mutation data matrices from protein sequences," *BMC Bioinformatics*, vol. 8.No(3), pp. 275–282, 1991.

[35] T. Muller, R. Spang, and M. Vingron, "Estimating amino acid substitution models: A comparison of dayhoff's estimator, the resolvent approach and a maximum likelihood method," *Mol Bio Evol*, vol. 19 No.1, pp. 8–13, 2002.

[36] S. E. Hambrusch and A. Khokhar, "C$^3$: A parallel model for coarse-grained machines," *Journal of Parallel and Distributed Computing*, vol. 32, no. 2, pp. 139–154, 1996.

[37] S. E. Hambrusch, F. Hameed, and A. A. Khokhar, "Communication operations on coarse-grained mesh architectures," *Parallel Computing*, vol. 21, no. 5, pp. 731–751, 1995.

[38] J. D. Thompson, F. Plewniak, and O. Poch, "Balibase: a benchmark alignment database for the evaluation of multiple alignment programs," *Bioinformatics*, vol. 15, no. 1, pp. 87–88, 1999.

[39] J. Stoye, D. Evers, and F. Meyer, "Rose: generating sequence families," *Bioinformatics*, vol. 14, pp. 157–163, March 1998.

[40] J. M. Sauder, J. W. Arthur, and R. L. Dunbrack, "Large-scale comparison of protein sequence alignment algorithms with structure alignments," *Proteins*, no. 40, pp. 6–22, 2000.

[41] M. Cline, R. Hughey, and K. Karplus, "Predicting reliable regions in protein sequence alignments," *Bioinformatics*, vol. 18, no. 2, pp. 306–314, 2002.

[42] T. Lassmann and E. L. Sonnhammer, "Quality assessment of multiple alignment programs.," *FEBS Lett*, vol. 529, pp. 126–130, October 2002.

[43] F. Saeed, A. Khokhar, O. Zagordi, and N. Beerenwinkel, "Multiple sequence alignment system for pyrosequencing reads," *Bioinformatics and Computational Biology (BICoB), LNBI 5462*, pp. 362–375, April 2009.