

High Performance Parallel Computing with Cloud and Cloud Technologies

Jaliya Ekanayake^{1,2}, Xiaohong Qiu¹, Thilina Gunarathne^{1,2}, Scott Beason¹, Geoffrey Fox^{1,2}
*¹Pervasive Technology Institute, ²School of Informatics and Computing,
Indiana University
{jekanaya, xqiu, tgunarat, smbeason, gcf}@indiana.edu*

Abstract

We present our experiences in applying, developing, and evaluating cloud and cloud technologies. First, we present our experience in applying Hadoop and DryadLINQ to a series of data/compute intensive applications and then compare them with a novel MapReduce runtime developed by us, named CGL-MapReduce, and MPI. Preliminary applications are developed for particle physics, bioinformatics, clustering, and matrix multiplication. We identify the basic execution units of the MapReduce programming model and categorize the runtimes according to their characteristics. MPI versions of the applications are used where the contrast in performance needs to be highlighted. We discuss the application structure and their mapping to parallel architectures of different types, and look at the performance of these applications. Next, we present a performance analysis of MPI parallel applications on virtualized resources.

1. Introduction

Cloud and cloud technologies are two broad categories of technologies related to the general notion of Cloud Computing. By “Cloud,” we refer to a collection of infrastructure services such as Infrastructure-as-a-service (IaaS), Platform-as-a-Service (PaaS), etc., provided by various organizations where virtualization plays a key role. By “Cloud Technologies,” we refer to

various cloud runtimes such as Hadoop(ASF, *core*, 2009), Dryad (Isard, Budiu et al. 2007), and other MapReduce(Dean and Ghemawat 2008) frameworks, and also the storage and communication frameworks such as Hadoop Distributed File System (HDFS), Amazon S3(Amazon 2009), etc.

The introduction of commercial cloud infrastructure services such as Amazon EC2, GoGrid(ServePath 2009), and ElasticHosts(ElasticHosts 2009) has allowed users to provision compute clusters fairly easily and quickly, by paying a monetary value for the duration of their usages of the resources. The provisioning of resources happens in minutes, as opposed to the hours and days required in the case of traditional queue-based job scheduling systems. In addition, the use of such virtualized resources allows the user to completely customize the Virtual Machine (VM) images and use them with root/administrative privileges, another feature that is hard to achieve with traditional infrastructures. The availability of open source cloud infrastructure software such as Nimbus(Keahey, Foster et al. 2005) and Eucalyptus(Nurmi, Wolski et al. 2009), and the open source virtualization software stacks such as Xen Hypervisor(Barham, Dragovic et al. 2003), allows organizations to build private clouds to improve the resource utilization of the available computation facilities. The possibility of dynamically provisioning additional resources by leasing from commercial cloud infrastructures makes the use of private clouds more promising.

Among the many applications which benefit from cloud and cloud technologies, the data/compute intensive applications are the most important. The deluge of data and the highly compute intensive applications found in many domains such as particle physics, biology, chemistry, finance, and information retrieval, mandate the use of large computing infrastructures and parallel processing to achieve considerable performance gains in analyzing data. The

addition of cloud technologies creates new trends in performing parallel computing. An employee in a publishing company who needs to convert a document collection, terabytes in size, to a different format can do so by implementing a MapReduce computation using Hadoop, and running it on leased resources from Amazon EC2 in just few hours. A scientist who needs to process a collection of gene sequences using CAP3(Huang and Madan 1999) software can use virtualized resources leased from the university's private cloud infrastructure and Hadoop. In these use cases, the amount of coding that the publishing agent and the scientist need to perform is minimal (as each user simply needs to implement a *map* function), and the MapReduce infrastructure handles many aspects of the parallelism.

Although the above examples are successful use cases for applying cloud and cloud technologies for parallel applications, through our research, we have found that there are limitations in using current cloud technologies for parallel applications that require complex communication patterns or require faster communication mechanisms. For example, Hadoop and Dryad implementations of Kmeans clustering applications which perform an iteratively refining clustering operation, show higher overheads compared to implementations of MPI or CGL-MapReduce (Ekanayake, Pallickara et al. 2008) – a streaming-based MapReduce runtime developed by us. These observations raise questions: What applications are best handled by cloud technologies? What overheads do they introduce? Are there any alternative approaches? Can we use traditional parallel runtimes such as MPI in cloud? If so, what overheads does it have? These are some of the questions we try to answer through our research.

In section 1, we give a brief introduction of the cloud technologies, and in section 2, we discuss with examples the basic functionality supported by these cloud runtimes. Section 3 discusses how these technologies map into programming models. We describe the applications

used to evaluate and test technologies in section 4. The performance results are in section 5. In section 6, we present details of an analysis we have performed to understand the performance implications of virtualized resources for parallel MPI applications. Note that we use MPI running on non virtual machines in section 5 for comparison with cloud technologies. We present our conclusions in section 7.

2. Cloud Technologies

The cloud technologies such as MapReduce and Dryad have created new trends in parallel programming. The support for handling large data sets, the concept of moving computation to data, and the better quality of services provided by the cloud technologies make them favorable choice of technologies to solve large scale data/compute intensive problems.

The granularity of the parallel tasks in these programming models lies in between the fine grained parallel tasks that are used in message passing infrastructures such as PVM(Dongarra, Geist et al. 1993) and MPI(Forum n.d.) ; and the coarse grained jobs in workflow frameworks such as Kepler(Ludscher, Altintas et al. 2006) and Taverna(Hull, Wolstencroft et al. 2006), in which the individual tasks could themselves be parallel applications written in MPI. Unlike the various communication constructs available in MPI which can be used to create a wide variety of communication topologies for parallel programs, in MapReduce, the “map->reduce” is the only communication construct available. However, our experience shows that most *composable* applications can easily be implemented using the MapReduce programming model. Dryad supports parallel applications that resemble Directed Acyclic Graphs (DAGs) in which the vertices represent computation units, and the edges represent communication channels between different computation units.

In traditional approaches, once parallel applications are developed, they are executed on

compute clusters, super computers, or Grid infrastructures (Foster 2001), where the focus on allocating resources is heavily biased by the availability of computational power. The application and the data both need to be moved to the available computational resource in order for them to be executed. These infrastructures are highly efficient in performing compute intensive parallel applications. However, when the volumes of data accessed by an application increases, the overall efficiency decreases due to the inevitable data movement. Cloud technologies such as Google MapReduce, Google File System (GFS) (Ghemawat, Gobioff et al. 2003), Hadoop and Hadoop Distributed File System (HDFS), Microsoft Dryad, and CGL-MapReduce adopt a more data-centered approach to parallel runtimes. In these frameworks, the data is staged in data/compute nodes of clusters or large-scale data centers, such as in the case of Google. The computations move to the data in order to perform the data processing. Distributed file systems such as GFS and HDFS allow Google MapReduce and Hadoop to access data via distributed storage systems built on heterogeneous compute nodes, while Dryad and CGL-MapReduce support reading data from local disks. The simplicity in the programming model enables better support for quality of services such as fault tolerance and monitoring.

2.1 Hadoop

Apache Hadoop has a similar architecture to Google's MapReduce runtime, where it accesses data via HDFS, which maps all the local disks of the compute nodes to a single file system hierarchy, allowing the data to be dispersed across all the data/computing nodes. HDFS also replicates the data on multiple nodes so that failures of any nodes containing a portion of the data will not affect the computations which use that data. Hadoop schedules the MapReduce computation tasks depending on the data locality, improving the overall I/O bandwidth. The outputs of the *map* tasks are first stored in local disks until later, when the *reduce* tasks access

them (pull) via HTTP connections. Although this approach simplifies the fault handling mechanism in Hadoop, it adds a significant communication overhead to the intermediate data transfers, especially for applications that produce small intermediate results frequently.

2.2 Dryad and DryadLINQ

Dryad is a distributed execution engine for coarse grain data parallel applications. It combines the MapReduce programming style with dataflow graphs to solve the computation tasks. Dryad considers computation tasks as directed acyclic graph (DAG)s where the vertices represent computation tasks and with the edges acting as communication channels over which the data flow from one vertex to another. The data is stored in (or partitioned to) local disks via the Windows shared directories and meta-data files and Dryad schedules the execution of vertices depending on the data locality. (Note: The academic release of Dryad only exposes the DryadLINQ (Y.Yu, Isard et al. 2008) API for programmers. Therefore, all our implementations are written using DryadLINQ although it uses Dryad as the underlying runtime). Dryad also stores the output of vertices in local disks, and the other vertices which depend on these results, access them via the shared directories. This enables Dryad to re-execute failed vertices, a step which improves the fault tolerance in the programming model.

2.3 CGL-MapReduce

CGL-MapReduce is a light-weight MapReduce runtime that incorporates several improvements to the MapReduce programming model such as (i) faster intermediate data transfer via a pub/sub broker network; (ii) support for long running *map/reduce* tasks; and (iii) efficient support for iterative MapReduce computations. The architecture of CGL-MapReduce is shown in figure 1 (Left).

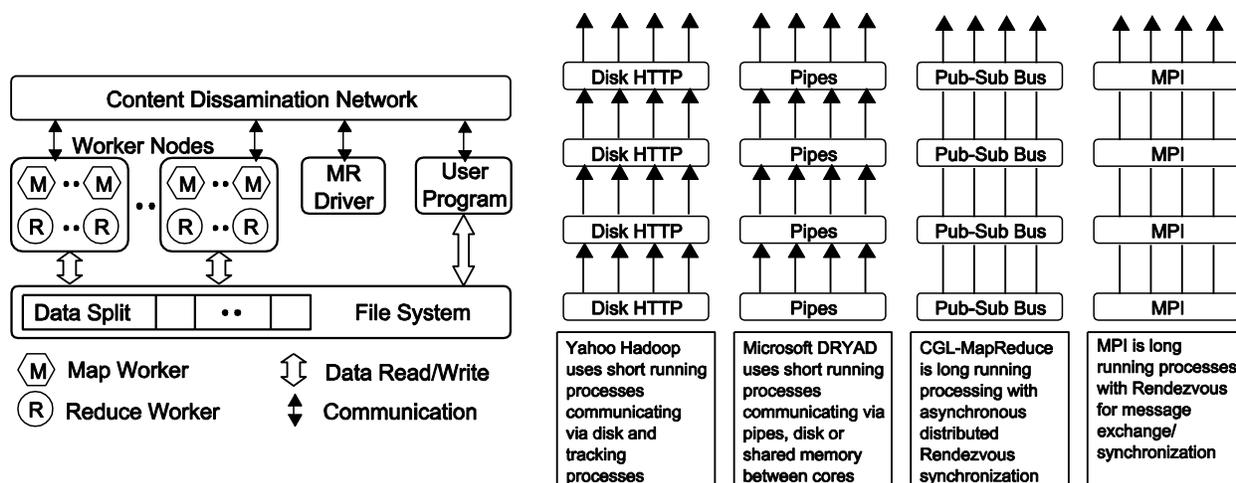


Figure 1. (Left) Components of the CGL-MapReduce. (Right) Different synchronization and intercommunication mechanisms used by the parallel runtimes.

The use of streaming enables CGL-MapReduce to send the intermediate results directly from its producers to its consumers, and eliminates the overhead of the file based communication mechanisms adopted by both Hadoop and Dryad. The support for long running *map/reduce* tasks enables configuring and re-using of *map/reduce* tasks in the case of iterative MapReduce computations, and eliminates the need for the re-configuring or the re-loading of static data in each iteration. This feature comes with the distinction of “static data” and the “dynamic data” that we support in CGL-MapReduce. We refer to any data set which is static throughout the computation as “static data,” and the data that is changing over the computation as “dynamic data.” Although this distinction is irrelevant to the MapReduce computations that have only one *map* phase followed by a *reduce* phase, it is extremely important for iterative MapReduce computations, in which the map tasks need to access a static (fixed) data again and again. Figure 1 (Right) highlights the synchronization and communication characteristics of Hadoop, Dryad, CGL-MapReduce, and MPI.

Additionally, CGL-MapReduce supports the distribution of smaller variable data sets to all the *map* tasks directly, a functionality similar to *MPI_Bcast()* that is often found to be useful in

many data analysis applications. Hadoop provides a similar feature via its distributed cache, in which a file or data is copied to all the compute nodes. Dryad provides a similar feature by allowing applications to add resources (files) which will be accessible to all the vertices. With the above features in place, CGL-MapReduce can be used to implement iterative MapReduce computations efficiently. In CGL-MapReduce, the data partitioning and distribution is left to the users to handle, and it reads data from shared file systems or local disks. Although the use of streaming makes CGL-MapReduce highly efficient, implementing fault tolerance with this approach is not as straightforward as it is in Hadoop or Dryad. We plan to implement fault tolerance in CGL-MapReduce by re-executing failed *map* tasks and redundant execution of *reduce* tasks.

2.4 MPI

MPI, the de-facto standard for parallel programming, is a language-independent communications protocol that uses a message-passing paradigm to share the data and state among a set of cooperative processes running on a distributed memory system. MPI specification (Forum, MPI) defines a set of routines to support various parallel programming models such as point-to-point communication, collective communication, derived data types, and parallel I/O operations.

Most MPI runtimes are deployed in computation clusters where a set of compute nodes are connected via a high-speed network connection yielding very low communication latencies (typically in microseconds). MPI processes typically have a direct mapping to the available processors in a compute cluster or to the processor cores in the case of multi-core systems. We use MPI as the baseline performance measure for the various algorithms that are used to evaluate the different parallel programming runtimes. Table 1 summarizes the different characteristics of

Hadoop, Dryad, CGL-MapReduce, and MPI.

Table 1. Comparison of features supported by different parallel programming runtimes.

Feature	Hadoop	Dryad	CGL-MapReduce	MPI
Programming Model	MapReduce	DAG based execution flows	MapReduce with a <i>Combine</i> phase	Variety of topologies constructed using the rich set of parallel constructs
Data Handling	HDFS	Shared directories/ Local disks	Shared file system / Local disks	Shared file systems
Intermediate Data Communication	HDFS/ Point-to-point via HTTP	Files/TCP pipes/ Shared memory FIFO	Content Distribution Network (NaradaBrokering (Pallickara and Fox 2003))	Low latency communication channels
Scheduling	Data locality/ Rack aware	Data locality/ Network topology based run time graph optimizations	Data locality	Available processing capabilities
Failure Handling	Persistence via HDFS Re-execution of map and reduce tasks	Re-execution of vertices	Currently not implemented (Re-executing map tasks, redundant reduce tasks)	Program level Check pointing OpenMPI (Gabriel, E., G.E. Fagg, et al. 2004), FT MPI
Monitoring	Monitoring support of HDFS, Monitoring MapReduce computations	Monitoring support for execution graphs	Programming interface to monitor the progress of jobs	Minimal support for task level monitoring
Language Support	Implemented using Java. Other languages are supported via Hadoop Streaming	Programmable via C# DryadLINQ provides LINQ programming API for Dryad	Implemented using Java Other languages are supported via Java wrappers	C, C++, Fortran, Java, C#

3. Programming models

When analyzing applications written in the MapReduce programming model, we can identify three basic execution units wiz: (i) map-only; (ii) map-reduce; and (iii) iterative-map-reduce. Complex applications can be built by combining these three basic execution units under the MapReduce programming model. Table 2 shows the data/computation flow of these three basic execution units, along with examples.

Table 2. Three basic execution units under the MapReduce programming model.

Map-only	Map-reduce	Iterative map-reduce
Cap3 Analysis (we will discuss more about this later) Converting a collection of documents to different formats, processing a collection of medical images, and brute force searches in cryptography Parametric sweeps	HEP data analysis (we will discuss more about this later) <i>Histogramming</i> operations, distributed search, and distributed sorting Information retrieval	Expectation maximization algorithms Kmeans clustering Matrix multiplication

In the MapReduce programming model, the tasks that are being executed at a given phase have similar executables and similar input and output operations. With zero *reduce* tasks, the MapReduce model reduces to a map-only model which can be applied to many “embarrassingly parallel” applications. Software systems such as batch queues, Condor(Condor 2009), Falcon (Raicu, Zhao et al. 2007) and SWARM (Pallickara and Pierce 2008) all provide similar functionality by scheduling large numbers of individual maps/jobs. Applications which can utilize a “reduction” or an “aggregation” operation can use both phases of the MapReduce model and, depending on the “associativity” and “transitivity” nature of the reduction operation, multiple reduction phases can be applied to enhance the parallelism. For example, in a histogramming operation, the partial histograms can be combined in any order and in any number of steps to produce a final histogram.

The “side effect free”-nature of the MapReduce programming model does not promote iterative MapReduce computations. Each *map* and *reduce* tasks are considered as atomic execution units with no state shared in between executions. In parallel runtimes such as those of the MPI, the parallel execution units live throughout the entire life of the program; hence, the

state of a parallel execution unit can be shared across invocations. We propose an intermediate approach to develop MapReduce computations. In our approach, the *map/reduce* tasks are still considered side effect-free, but the runtime allows configuring and re-usage of the *map/reduce* tasks. Once configured, the runtime caches the *map/reduce* tasks. This way, both *map* and *reduce* tasks can keep the static data in memory, and can be called iteratively without loading the static data repeatedly.

Hadoop supports configuring the number of *reduce* tasks, which enables the user to create “map-only” applications by using zero *reduce* tasks. Hadoop can be used to implement iterative MapReduce computations, but the framework does not provide additional support to implement them efficiently. The CGL-MapReduce supports all the above three execution units, and the user can develop applications with multiple stages of MapReduce by combining them in any order. Dryad execution graphs resembling the above three basic units can be generated using DryadLINQ operations. DryadLINQ adds the LINQ programming features to Dryad where the user can implement various data analysis applications using LINQ queries, which will be translated to Dryad execution graphs by the compiler. However, unlike in the MapReduce model, Dryad allows the concurrent vertices to have different behaviors and different input/output characteristics, thus enabling a more workflow style programming model. Dryad also allows multiple communication channels in between different vertices of the data flow graph. Programming languages such as Swazall (Pike, Dorward et al. 2005), introduced by Google for its MapReduce runtime enables high level language support for expressing MapReduce computations, and the Pig (ASF, *pig*, 2009) available as a sub project of Hadoop, allows query operations on large data sets.

Apart from these programming models, there are other software frameworks that one can use

to perform data/compute intensive analyses. Disco (Nokia 2009) is an open source MapReduce runtime developed using a functional programming language named Erlang (Ericsson 2009). Disco architecture shares clear similarities with both the Google and the Hadoop MapReduce architectures. Sphere (Gu and Grossman 2009) is a framework which can be used to execute user-defined functions in parallel on data stored in a storage framework named Sector. Sphere can also perform MapReduce style programs, and the authors compare the performance with Hadoop for tera-sort application. All-Pairs (Moretti, Bui et al. 2009) is an abstraction that can be used to solve the common problem of comparing all the elements in a data set with all the elements in another data set by applying a given function. This problem can be implemented using Hadoop and Dryad as well, and we discuss a similar problem in section 4.4. We can also develop an efficient iterative MapReduce implementation using CGL-MapReduce to solve this problem. The algorithm is similar to the matrix multiplication algorithm that we will explain in section 4.3.

MPI and threads are two other programming models that can be used to implement parallel applications. MPI can be used to develop parallel applications in distributed memory architectures whereas threads can be used in shared memory architectures, especially in multi-core nodes. The low level communication constructs available in MPI allows users to develop parallel applications with various communication topologies involving fine grained parallel tasks. The use of low latency network connections between nodes enables applications to perform large number of inter-task communications. In contrast, the next generation parallel runtimes such as MapReduce and Dryad provide small number of parallel constructs such as “map-only”, “map-reduce”, “Select”, “Apply”, “Join” etc., and does not require high speed communication channels. These constraints require adopting parallel algorithms which perform

coarse grained parallel tasks and less communication. The use of threads is a natural approach in shared memory architectures, where communication between parallel tasks reduces to the simple sharing of pointers via the shared memory. However, the operating system's support for user level threads plays a major role in achieving better performances using multi-threaded applications. We will discuss the issues in using threads and MPI in more detail in section 5.4.2.

4. Data Analyses Applications

4.1 CAP3 – Sequence Assembly Program

CAP3 is a DNA sequence assembly program developed by Huang and Madan (1999) that performs several major assembly steps: these steps include computation of overlaps, construction of contigs, construction of multiple sequence alignments, and generation of consensus sequences to a given set of gene sequences. The program reads a collection of gene sequences from an input file (FASTA file format) and writes its output to several output files, as well as the standard output (as shown below).

Input.fsa -> CAP3 -> Stdout + Other output files

The program structure of this application fits directly with the “Map-only” basic execution unit, as shown in Table 2. We implemented a parallel version of CAP3 using Hadoop, CGL-MapReduce, and DryadLINQ. Each *map* task in Hadoop and in CGL-MapReduce calls the CAP3 executable as a separate process for a given input data file (the input “Value” for the *map* task), whereas in DryadLINQ, a “homomorphic Apply” operation calls the CAP3 executable on each data file in its data partition as a separate process. All the implementations move the output files to a predefined shared directory. This application resembles a common parallelization requirement, where an executable script, or a function in a special framework such as Matlab or

R, needs to be executed on each input data item. The above approach can be used to implement all these types of applications using any of the above three runtimes.

4.2 High Energy Physics

Next, we applied the MapReduce technique to parallelize a High Energy Physics (HEP) data analysis application, and implemented it using Hadoop, CGL-MapReduce, and Dryad. The HEP data analysis application processes large volumes of data, and performs a histogramming operation on a collection of event files produced by HEP experiments. The details regarding the two MapReduce implementations and the challenges we faced in implementing them can be found in (Ekanayake, Pallickara et al. 2008). In the DryadLINQ implementation, the input data files are first distributed among the nodes of the cluster manually. We developed a tool to perform the manual partitioning and distribution of the data. The names of the data files available in a given node were used as the data to the DryadLINQ program. Using a homomorphic “Apply” operation, we executed a ROOT interpreted script on groups of input files in all the nodes. The output histograms of this operation were written to a pre-defined shared directory. Next, we used another “Apply” phase to combine these partial histograms into a single histogram using DryadLINQ.

4.3 Iterative MapReduce – Kmeans Clustering and Matrix Multiplication

Parallel applications that are implemented using message passing runtimes can utilize various communication constructs to build diverse communication topologies. For example, a matrix multiplication application that implements Fox's Algorithm (Fox and Hey, 1987) and Cannon's Algorithm (Johnsson, Harris et al. 1989) assumes parallel processes to be in a rectangular grid. Each parallel process in the grid communicates with its left and top neighbors as shown in figure 2 (left). The current cloud runtimes, which are based on data flow models such as MapReduce

and Dryad, do not support this behavior, in which the peer nodes communicate with each other. Therefore, implementing the above type of parallel applications using MapReduce or DryadLINQ requires adopting different algorithms.

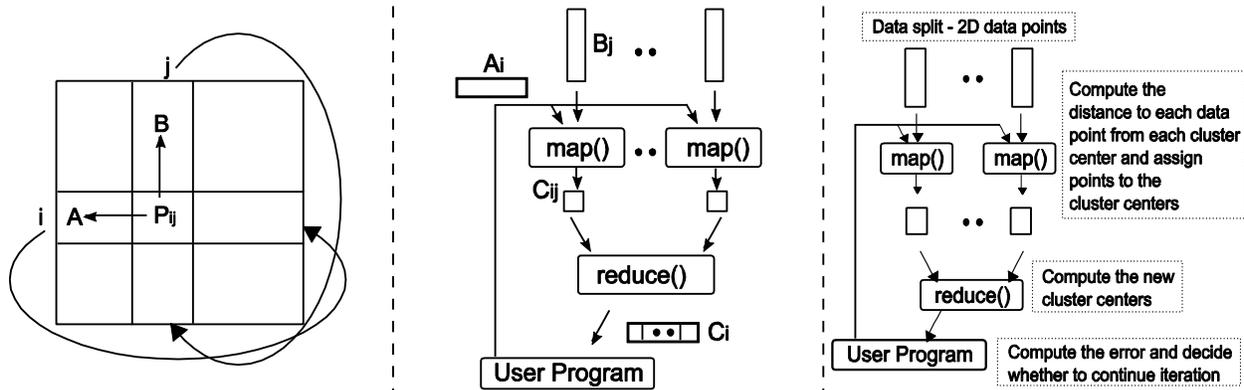


Figure 2. (Left) The communication topology of Cannon's Algorithm implemented using MPI, (middle) Communication topology of matrix multiplication application based on MapReduce, and (right) Communication topology of Kmeans Clustering implemented as a MapReduce application.

We have implemented matrix multiplication applications using Hadoop and CGL-MapReduce by adopting a row/column decomposition approach to split the matrices. To clarify our algorithm, let's consider an example where two input matrices, A and B , produce matrix C , as the result of the multiplication process. We split the matrix B into a set of column blocks and the matrix A into a set of row blocks. In each iteration, all the *map* tasks process two inputs: (i) a column block of matrix B , and (ii) a row block of matrix A ; collectively, they produce a row block of the resultant matrix C . The column block associated with a particular *map* task is fixed throughout the computation, while the row blocks are changed in each iteration. However, in Hadoop's programming model (a typical MapReduce model), there is no way to specify this behavior. Hence, it loads both the column block and the row block in each iteration of the computation. CGL-MapReduce supports the notion of long running *map/reduce* tasks where these tasks are allowed to retain static data in the memory across invocations, yielding better performance for "Iterative MapReduce" computations. The communication pattern of this

application is shown in figure 2 (middle).

Kmeans clustering (Macqueen 1967) is another application that performs iteratively refining computation. We also implemented Kmeans clustering applications using Hadoop, CGL-MapReduce, and DryadLINQ. In the two MapReduce implementations, each *map* task calculates the distances between all the data elements in its data partition, to all the cluster centers produced during the previous run. It then assigns data points to these cluster centers, based on their Euclidian distances. The communication topology of this algorithm is shown in figure 2 (right). Each *map* task produces partial cluster centers as the output; these are then combined at a *reduce* task to produce the current cluster centers. These current cluster centers are used in the next iteration, to find the next set of cluster centers. This process continues until the overall distance between the current cluster centers and the previous cluster centers, reduces below a predefined threshold. The Hadoop implementation uses a new MapReduce computation for each iteration of the program, while CGL-MapReduce's long running *map/reduce* tasks allows it to re-use *map/reduce* tasks. The DryadLINQ implementation uses various DryadLINQ operations such as "Apply", "GroupBy", "Sum", "Max", and "Join" to perform the computation, and also, it utilizes DryadLINQ's "loop unrolling" support to perform multiple iterations as a single large query.

4.4 ALU Sequencing Studies

4.4.1 ALU Clustering

The ALU clustering problem (Batzer and Deininger 2002) is one of the most challenging problems for sequence clustering because ALUs represent the largest repeat families in human genome. There are about 1 million copies of ALU sequences in human genome, in which most insertions can be found in other primates and only a small fraction (~ 7000) are human-specific. This indicates that the classification of ALU repeats can be deduced solely from the 1 million

human ALU elements. Notable, ALU clustering can be viewed as a classical case study for the capacity of computational infrastructures because it is not only of great intrinsic biological interests, but also a problem of a scale that will remain as the upper limit of many other clustering problem in bioinformatics for the next few years, e.g. the automated protein family classification for a few millions of proteins predicted from large metagenomics projects.

4.4.2 Smith Waterman Dissimilarities

We identified samples of the human and Chimpanzee ALU gene sequences using Repeatmasker (Smith and Hubley 2004) with Repbase Update (Jurka 2000). We have been gradually increasing the size of our projects with the current largest samples having 35339 and 50000 sequences and these require a modest cluster such as Tempest (768 cores) for processing in a reasonable time (a few hours as shown in section 5). Note from the discussion in section 4.4.1, we are aiming at supporting problems with a million sequences -- quite practical today on TeraGrid and equivalent facilities given basic analysis steps scale like $O(N^2)$.

We used open source version NAligner (Smith Waterman Software) of the Smith Waterman – Gotoh algorithm SW-G (Smith, Waterman et al 1981; Gotoh, 1982) modified to ensure low start up effects by each thread/processing large numbers (above a few hundred) at a time. Memory bandwidth needed was reduced by storing data items in as few bytes as possible.

4.4.3 The $O(N^2)$ Factor of 2 and structure of processing algorithm

The ALU sequencing problem shows a well known factor of 2 issue present in many $O(N^2)$ parallel algorithms such as those in direct simulations of astrophysical stems. We initially calculate in parallel the Distance $D(i,j)$ between points (sequences) i and j . This is done in parallel over all processor nodes selecting criteria $i < j$ (or $j > i$ for upper triangular case) to avoid calculating both $D(i,j)$ and the identical $D(j,i)$. This can require substantial file transfer as it is

unlikely that nodes requiring $D(i,j)$ in a later step will find that it was calculated on nodes where it is needed.

For example the MDS and PW (PairWise) Clustering algorithms described in Fox, Bae et al. (2008), require a parallel decomposition where each of N processes (MPI processes, threads) has $1/N$ of sequences and for this subset $\{i\}$ of sequences stores in memory $D(\{i\},j)$ for all sequences j and the subset $\{i\}$ of sequences for which this node is responsible. This implies that we need $D(i,j)$ and $D(j,i)$ (which are equal) stored in different processors/disks. This is a well known collective operation in MPI called either gather or scatter.

4.4.4 *Dryad Implementation*

We developed a DryadLINQ application to perform the calculation of pairwise SW-G distances for a given set of genes by adopting a coarse grain task decomposition approach which requires minimum inter-process communication to ameliorate the higher communication and synchronization costs of the parallel runtime. To clarify our algorithm, let's consider an example where N gene sequences produces a pairwise distance matrix of size $N \times N$. We decompose the computation task by considering the resultant matrix and groups the overall computation into a block matrix of size $D \times D$ where D is a multiple (>2) of the available computation nodes. Due to the symmetry of the distances $D(i,j)$ and $D(j,i)$ we only calculate the distances in the blocks of the upper triangle of the block matrix as shown in figure 3 (left). The blocks in the upper triangle are partitioned (assigned) to the available compute nodes and an "Apply" operation is used to execute a function to calculate $(N/D) \times (N/D)$ distances in each block. After computing the distances in each block, the function calculates the transpose matrix of the result matrix which corresponds to a block in the lower triangle, and writes both these matrices into two output files in the local file system. The names of these files and their block numbers are communicated back

to the main program. The main program sort the files based on their block numbers and performs another “Apply” operation to combine the files corresponding to a row of blocks in a single large row block as shown in the figure 3 (right).

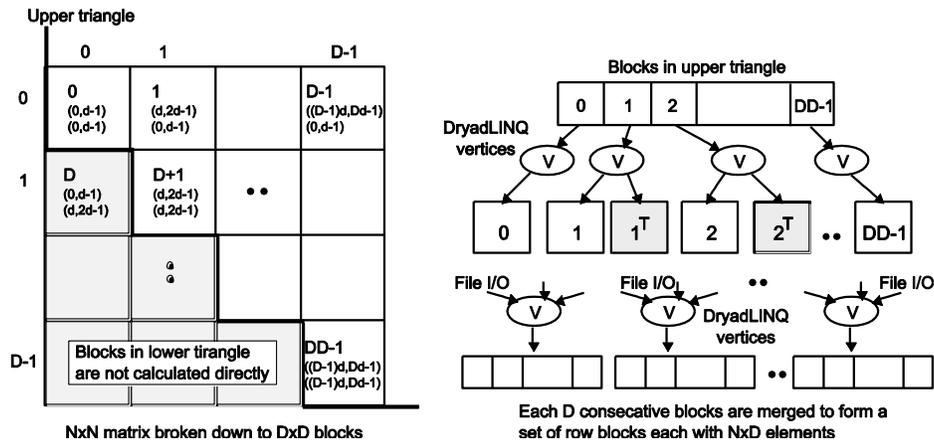


Figure 3. Task decomposition (left) and the DryadLINQ vertex hierarchy (right) of the DryadLINQ implementation of SW-G pairwise distance calculation application.

4.4.5 MPI Implementation

The MPI version of SW-G calculates pairwise distances using a set of either single or multi-threaded processes. For N gene sequences, we need to compute half of the values (in the lower triangular matrix), which is a total of $M = N \times (N-1) / 2$ distances. At a high level, computation tasks are evenly divided among P processes and execute in parallel. Namely, computation workload per process is M/P . At a low level, each computation task can be further divided into subgroups and run in T concurrent threads. Our implementation is designed for flexible use of shared memory multicore system and distributed memory clusters (tight to medium tight coupled communication technologies such threading and MPI). We provide options for any combinations of thread vs. process vs. node as shown in figure 4. The real computation workload per parallel unit is decided by $M / (T \times P \times \# \text{ nodes})$.

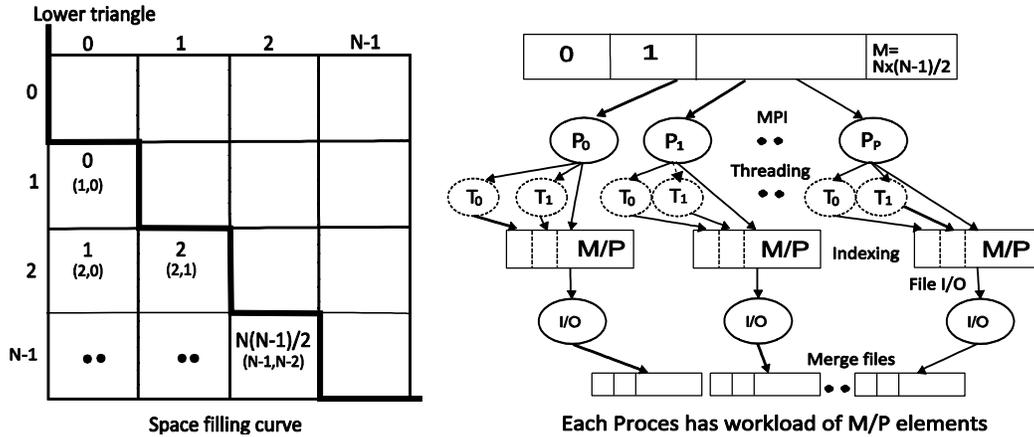


Figure 4. Task decomposition (left) and the MPI (right) implementation of SW-G pairwise distance calculation application.

As illustrated in figure 4, the data decomposition strategy runs a “space filling curve through lower triangular matrix” to produce equal numbers of pairs for each parallel unit such as process or thread. It is necessary to map indexes in each pairs group back to corresponding matrix coordinates (i, j) for constructing full matrix later on. We implemented a special function "PairEnumertator" as the convertor. We tried to limit runtime memory usage for performance optimization. This is done by writing a triple of i, j , and distance value of pairwise alignment to a stream writer and the system flushes accumulated results to a local file periodically. As the final stage, individual files are merged to form a full distance matrix.

5. Evaluations

5.1 Introduction

For our evaluations, we used three compute clusters (details are shown in Table 3) with two 32-node clusters have almost identical hardware configurations and one latest 32-node cluster of 24 core machines with Infiniband connections. DryadLINQ and the MPI application that performs SW-G computation were run on the Windows cluster (Ref B, Ref C), while the Hadoop, CGL-MapReduce, and other MPI applications were run on the Linux cluster (Ref A).

We measured the performance of these applications, and present the results in terms of parallel overhead defined for Parallelism P by

$$f(P) = [P * T(P) - T(1)] / T(1) \quad (1)$$

where P denotes parallelism (e.g. processes, threads, *map* tasks) used, and T denotes time as a function of the number of parallel processes used. T(1) is replaced in practice by T(S) where S is the smallest number of processes that can run the job. We used Hadoop release 0.20, the academic release of DryadLINQ, Microsoft MPI, and OpenMPI version 1.3.2 for our evaluations.

Table 3. Different computation clusters used for the analyses

Feature	Linux Cluster (Ref A)	Windows Cluster (Ref B)	Windows Cluster (Ref C)
# Node	32	32	32
CPU	Intel(R) Xeon(R) CPU L5420 2.50GHz	Intel(R) Xeon(R) CPU L5420 2.50GHz	Intel(R) Xeon(R) CPU E7450 2.40GHz
# CPU / # Cores	2/8	2/8	4/24
Total Cores	256	256	768
Memory	32GB	16 GB	48 GB
Disk	1 Disk of Western Digital Caviar RE 160 GB SATA 7200	2 Disks of 1000GB (1TB) Ultrastar A7K1000 7200	2 HP 146GB 10K 2.5 SAS HP SP HDD
Network	Giga bit Ethernet	Giga bit Ethernet	20 Gbps Infiniband
Operating System	Red Hat Enterprise Linux Server release 5.3 - 64 bit	Windows Server Enterprise - 64 bit	Windows Server 2008 HPC Edition (Service Pack 1)

5.2 The CAP3 and Particle Physics Case Studies

The results of our performance measurements for CAP3 and Particle Physics are shown in figures 5-8.

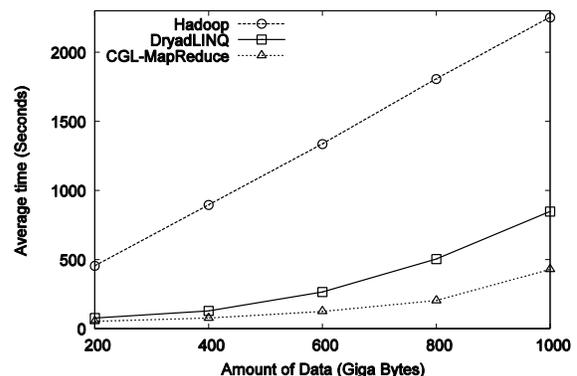
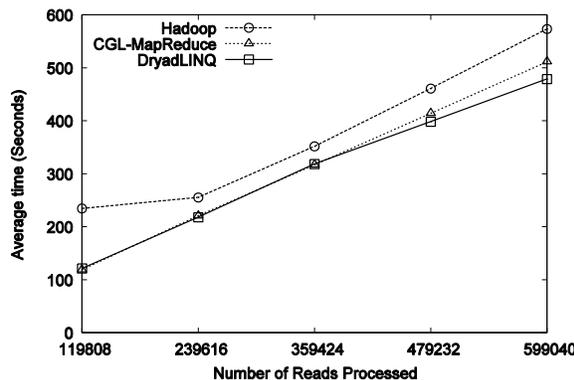


Figure 5. Performance of the CAP3 application – average time (in seconds) against the number of gene reads processed.

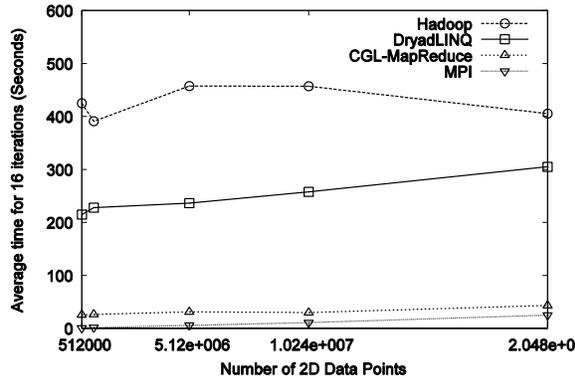


Figure 7. Overhead induced by different parallel programming runtimes for the Kmeans Clustering application – overhead against the number of 2D data points clustered (Note: Both axes are in log scale)

Figure 6. Performance of the HEP data analysis application – average time (in seconds) against the amount of input data processed (in Giga bytes).

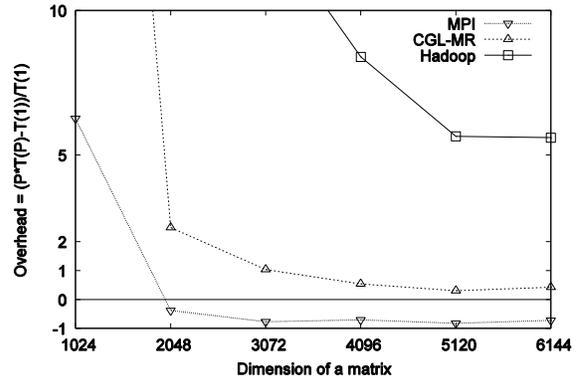


Figure 8. Overhead induced by different parallel programming runtimes for the matrix multiplication application – overhead against the dimension of an input matrix.

From these results, it is clearly evident that the cloud runtimes perform competitively well for both the “Map-only” and the “Map-reduce” style applications. In the HEP data analysis, both the CGL-MapReduce and DryadLINQ access input data from local disks, where the data is partitioned and distributed beforehand. Currently, HDFS can be accessed using Java or C++ clients only, and the ROOT – data analysis framework developed at CERN- interpretable scripts are not capable of accessing data from HDFS. Therefore, we placed the input data in the IU Data Capacitor – a high performance parallel file system based on the Lustre file system, which allows each *map* task in Hadoop to directly access data from this file system. The performance results shows that this dynamic data movement in the Hadoop implementation incurred considerable overhead to the computation, while the ability of reading input data from local disks gave significant performance improvement to both DryadLINQ and CGL-MapReduce, as compared to the Hadoop implementation.

5.3 Kmeans and Matrix Multiplication Case Studies

For iterative class of applications, cloud runtimes show considerably high overheads,

compared to the MPI and CGL-MapReduce versions of the same applications; the results shown in figures 7 and 8 imply that, for these types of applications, we still need to use high performance parallel runtimes or use alternative approaches. (Note: The negative overheads observed in the matrix multiplication application are due to the better utilization of a cache by the parallel application than the single process version). CGL-MapReduce shows a close performance closer to the MPI for large data sets in the case of Kmeans clustering and matrix multiplication applications, an outcome which highlights the benefits of supporting iterative computations and the faster data communication mechanism present in the CGL-MapReduce.

5.4 ALU Sequence Analysis Case Study

5.4.1 Performance of Smith Waterman Gotoh SW-G Algorithm

We performed the Dryad and MPI implementations of ALU SW-G distance calculations on two large data sets and obtained the following results.

Table 4 Comparison of DryadLINQ and MPI technologies on ALU sequencing application with SW-G algorithm

Technology		Total Time (seconds)	Time per Pair (milliseconds)	Partition Data (seconds)	Calculate and Output Distance(seconds)	Merge files (seconds)
Dryad	50,000 sequences	17200.413	0.0069	2.118	17104.979	93.316
	35,339 sequences	8510.475	0.0068	2.716	8429.429	78.33
MPI	50,000 sequences	16588.741	0.0066	N/A	13997.681	2591.06
	35,339 sequences	8138.314	0.0065	N/A	6909.214	1229.10

There is a short partitioning phase for DryadLINQ and then both approaches calculate the distances and write out these to intermediate files as discussed in section 4. We note that merge time is currently much longer for MPI than DryadLINQ while the initial steps are significantly faster for MPI. However the total times in table 4 indicates that both MPI and DryadLINQ implementations perform well for this application with MPI a few percent faster with current implementations. As expected, the times scale proportionally to the square of the number of

distances. On 744 cores the average time of 0.0067 milliseconds per pair that corresponds to roughly 5 milliseconds per pair calculated per core used. The coarse grained Dryad application performs competitively with the tightly synchronized MPI application. It proves once more the applicability of the Cloud technologies for the composable applications.

5.4.2 Threaded Implementation

Above we looked at using MPI with one process per core and we compared this with a threaded implementation with each process having several threads. Labeling configuration as $t \times m \times n$ for t threads per process, m MPI processes per node and n nodes, we compare choices of t , m and n in figure 9.

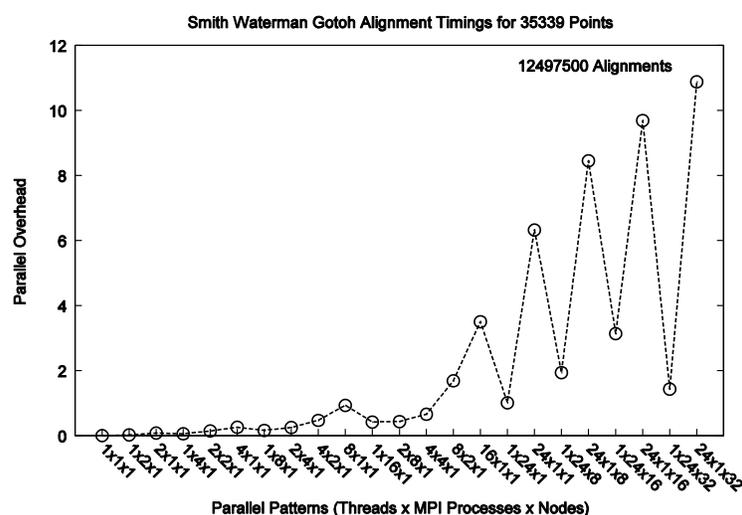


Figure 9. Performance of ALU Gene Alignments for different parallel patterns

The striking result for this step is that MPI easily outperforms the equivalent threaded version of this embarrassingly parallel step. In figure 9, all the peaks in the overhead correspond to patterns with large values of t . Note that the MPI intra-node $1 \times 24 \times 32$ pattern completes the full 624 billion alignments in 2.33 hours – 4.9 times faster than threaded implementation $24 \times 1 \times 32$. This 768 core MPI run has a parallel overhead of 1.43 corresponding to a speed up of 316.

The SW-G alignment performance is probably dominated by memory bandwidth issues, and

we are pursuing several points that could affect this though it is not at our highest priority as SW-G is not the dominant step. We have tried to identify the reason for the comparative slowness of threading. Using Windows monitoring tools, we found that the threaded version has about a factor of 100 more context switches than in the one thread per process MPI version. This could lead to a slow down of the threaded approach and correspond to Windows handing of paging of threads with large memory footprints.

6. The Performance of MPI on Clouds

After the previous observations, we analyzed the performance implications of cloud for parallel applications implemented using MPI. Specifically, we were trying to find the overhead of virtualized resources, and understand how applications with different communication-to-computation (C/C) ratios perform on cloud resources. We also evaluated different CPU core assignment strategies for VMs, in order to understand the performance of VMs on multi-core nodes.

Table 5. Computation and the communication complexities of the different MPI applications used.

Application	Matrix multiplication	Kmeans Clustering	Concurrent Wave Equation
Description	Implements Cannon's Algorithm	Implements Kmeans Clustering Algorithm	A vibrating string is decomposed (split) into points, and each MPI process is responsible for updating the amplitude of a number of points over time.
	Assume a rectangular process grid (figure 1-left)	A fixed number of iterations are performed in each test	
Grain size (n)	The number of points in a matrix block handled by each MPI process	The number of data points handled by a single MPI process	Number of points handled by each MPI process
Communication Pattern	Each MPI process communicates with its neighbors both row wise and column wise	All MPI processes send partial clusters to one MPI process (rank 0). Rank 0 distribute the new cluster centers to all the nodes	In each iteration, each MPI process exchanges boundary points with its nearest neighbors
Computation per MPI process	$O((\sqrt{n})^3)$	$O(n)$	$O(n)$
Communication per MPI process	$O((\sqrt{n})^2) = O(n)$	$O(1)$	$O(1)$

C/C	$O\left(\frac{1}{\sqrt{n}}\right)$	$O\left(\frac{1}{n}\right)$	$O\left(\frac{1}{n}\right)$
Message Size	$(\sqrt{n})^2 = n$	D – Where D is the number of cluster centers. $D \ll n$	Each message contains a double value
Communication routines used	<code>MPI_Sendrecv_replace()</code>	<code>MPI_Reduce()</code> <code>MPI_Bcast()</code>	<code>MPI_Sendrecv()</code>

Commercial cloud infrastructures do not allow users to access the bare hardware nodes, in which the VMs are deployed, a must-have requirement for our analysis. Therefore, we used a Eucalyptus-based cloud infrastructure deployed at our university for this analysis. With this cloud infrastructure, we have complete access to both virtual machine instances and to the underlying bare-metal nodes, as well as the help of the administrators; as a result, we could deploy different VM configurations, allocating different CPU cores to each VM. Therefore, we selected the above cloud infrastructure as our main test bed.

For our evaluations, we selected three MPI applications with different communication and computation requirements, namely, (i) the Matrix multiplication, (ii) Kmeans clustering, and (iii) the Concurrent Wave Equation solver. Table 5 highlights the key characteristics of the programs that we used as benchmarks.

6.1 Benchmarks and Results

The Eucalyptus (version 1.4) infrastructure we used is deployed on 16 nodes of an iDataplex cluster, each of which has 2 Quad Core Intel Xeon processors (for a total of 8 CPU cores) and 32 GB of memory. In the bare-metal version, each node runs a Red Hat Enterprise Linux Server release 5.2 (Tikanga) operating system. We used OpenMPI version 1.3.2 with gcc version 4.1.2. We then created a VM image from this hardware configuration, so that we would have a similar software environment on the VMs once they were deployed. The virtualization is based on the Xen hypervisor (version 3.0.3). Both bare-metal and virtualized resources utilized giga-bit

Ethernet connections.

When VMs are deployed using Eucalyptus, it allows us to configure the number of CPU cores assigned to each VM image. For example, with 8 core systems, the CPU core allocation per VM can range from 8 cores to 1 core per VM, resulting in several different CPU core assignment strategies. In an Amazon EC2 infrastructure, the standard instance type has $\frac{1}{2}$ a CPU per VM instance (Evangelinos and Hill 2008). In the current version of Eucalyptus, the minimum number of cores that we can assign for a particular VM instance is 1; hence, we selected five CPU core assignment strategies (including the bare-metal test) listed in Table 6.

Table 6. Different hardware/virtual machine configurations used for our performance evaluations.

Ref	Description	Number of CPU cores accessible to the virtual or bare-metal node	Amount of memory (GB) accessible to the virtual or bare-metal node	Number of virtual or bare-metal nodes deployed
BM	Bare-metal node	8	32	16
1-VM-8-core	1 VM instance per bare-metal node	8	30 (2GB is reserved for Dom0)	16
2-VM-4-core	2 VM instances per bare-metal node	4	15	32
4-VM-2-core	4 VM instances per bare-metal node	2	7.5	64
8-VM-1-core	8 VM instances per bare-metal node	1	3.75	128

We ran all the MPI tests, on all 5 hardware/VM configurations, and measured the performance and calculated speed-ups and overheads. We calculated two types of overheads for each application using formula (1). The total overhead induced by the virtualization and the parallel processing is calculated using the bare-metal single process time as $T(1)$ in the formula (1). The parallel overhead is calculated using the single process time from a corresponding VM as $T(1)$ in formula (1).

In all the MPI tests we performed, we used the following invariant to select the number of parallel processes (MPI processes) for a given application.

$$\text{Number of MPI processes} = \text{Number of CPU cores used}$$

For example, for the matrix multiplication application, we used only half the number of nodes (bare-metal or VMs) available to us, so that we had 64 MPI processes = 64 CPU cores. (This is mainly because the matrix multiplication application expects the MPI processes to be in a square grid, in contrast to a rectangular grid). For Kmeans clustering, we used all the nodes, resulting in a total of 128 MPI processes utilizing all 128 CPU cores. Some of the results of our analysis highlighting the different characteristics we observed are shown in figures 10 - 17.

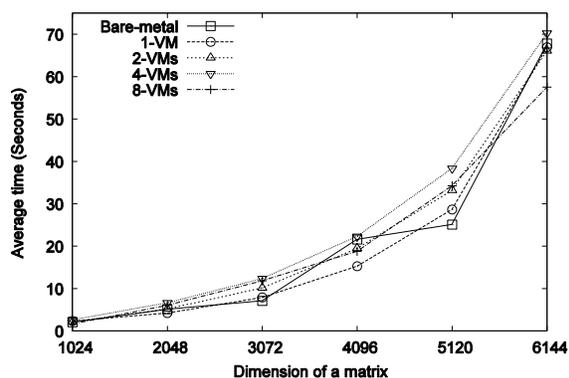


Figure 10. Performance of the matrix multiplication application – average time (in seconds) against the size of a matrix (Number of MPI processes = 64).

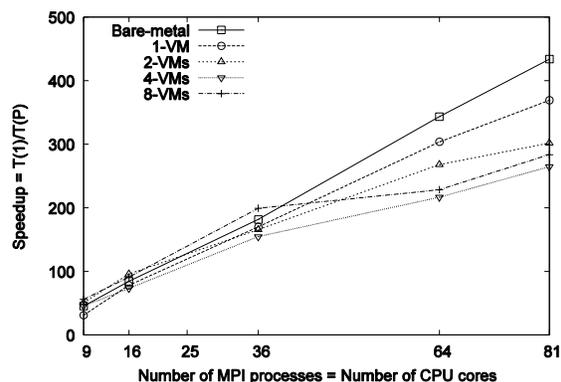


Figure 11. Speed-up of the matrix multiplication application – Speedup against the number of MPI processes = number of CPU cores used (Fixed matrix size = 5184x5184).

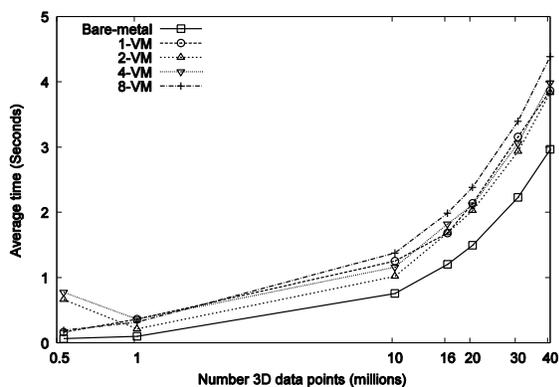


Figure 12. Performance of Kmeans clustering – average time (in seconds) against the number of 2D data points clustered (Number of MPI Processes = 128)

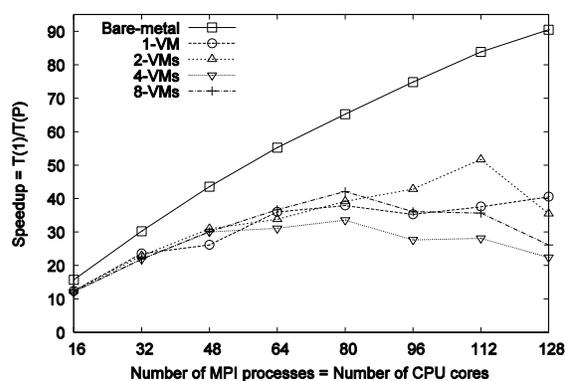


Figure 13. Speed-up of the Kmeans clustering – speedup against the number of MPI processes = number of CPU cores used (Number of data points = 860160)

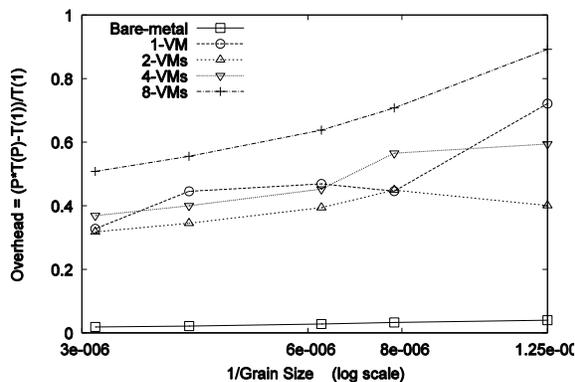


Figure 14. Total overhead of the Kmeans clustering – overhead against the 1/ grain size, grain size = number of 2D data points per parallel task (Number of MPI Processes = 128)

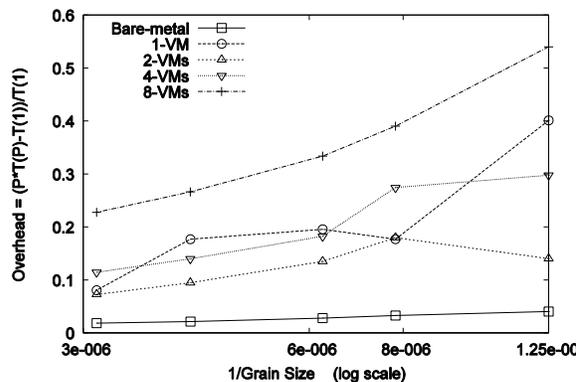


Figure 15. Parallel overhead of the Kmeans clustering – parallel overhead against 1/grain size (Number of MPI Processes = 128)

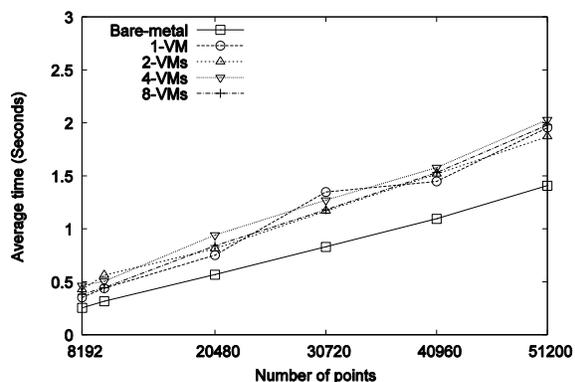


Figure 16. Performance of the Concurrent Wave Solver – average time (in seconds) against the number of points computed (Number of MPI Processes = 128)

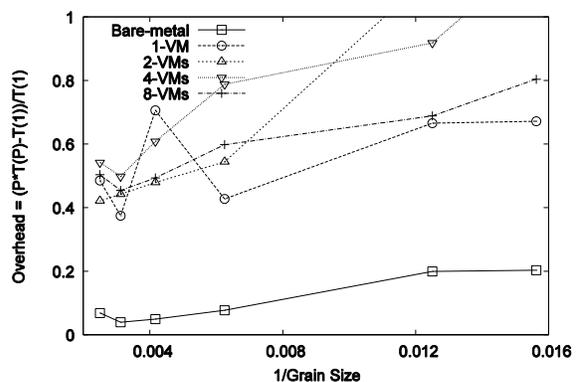


Figure 17. Total overhead of the Concurrent Wave Solver – overhead against 1/grain size, grain size = number of points assigned per parallel task (Number of MPI Processes = 128)

For the matrix multiplication, the graphs show very close performance characteristics in all the different hardware/VM configurations. As we expected, the bare-metal has the best performance and speedup values, compared to the VM configurations (apart from the region close to the matrix size of 4096x4096, where the VM performed better than the bare-metal. We have performed multiple tests at this point, and found that it is due to the cache performances of the bare-metal node). After the bare-metal, the next best performance and speed-ups were recorded in the case of 1-VM per bare-metal node configuration, in which the performance difference was mainly due to the overhead induced by the virtualization. However, as we

increased the number of VMs per bare-metal node, the overhead increased as well. At the 81 processes, the 8-VMs per node configuration shows about a 34% decrease in speed-up compared to the bare-metal results.

In Kmeans clustering, the effect of virtualized resources is much clearer than in the case of the matrix multiplication. All VM configurations show a lower performance compared to the bare-metal configuration. In this application, the amount of data transferred between MPI processes is extremely low compared to the amount of data processed by each MPI process, and also, in relation to the amount of computations performed. Figure 14 and figure 15 show the total overhead and the parallel overhead for Kmeans clustering under different VM configurations. From these two calculations, we found that, for VM configurations, the overheads are extremely large for data set sizes of less than 10 million points, for which the bare-metal overhead remains less than 1 (<1 for all the cases). For larger data sets such as those of 40 million points, all overheads reached less than 0.5. The slower speed-up of the VM configurations (shown in figure 13) is due to the use of a smaller data set (~800K points) to calculate the speed-ups. The overheads are extremely large for this region of the data sizes, and hence, it resulted in lower speed-ups for the VMs.

Concurrent wave equation splits a number of points into a set of parallel processes, and each parallel process updates its portion of the points in some number of steps. An increase in the number of points increases the amount of computations performed. Since we fixed the number of steps in which the points were updated, we obtained a constant amount of communication in all the test cases, resulting in a C/C ratio of $O(1/n)$. In this application also, the difference in performance between the VMs and the bare-metal version was clearer, and at the highest grain size, the total overhead of 8-VMs per node is about 7 times higher than the overhead of the bare-

metal configuration. The performance differences between the different VM configurations became smaller with the increase in grain size.

From the above experimental results, we can see that the applications with lower C/C ratios experienced a slower performance in virtualized resources. When the amount of data transferred between MPI processes is large, as in the case of the matrix multiplication, the application is more susceptible to the bandwidth than the latency. From the performance results of the matrix multiplication, we can see that the virtualization has not affected the bandwidth considerably. However, all the other results show that the virtualization has caused considerable latencies for parallel applications, especially with smaller data transfer requirements. The effect on latency increases as we use more VMs in a bare-metal node.

According to the Xen para-virtualization architecture (Barham, Dragovic, et al. 2003), domUs (VMs that run on top of a Xen para-virtualization) are not capable of performing I/O operations by themselves. Instead, they communicate with dom0 (privileged OS) via an event channel (interrupts) and the shared memory, and then the dom0 performs the I/O operations on behalf of the domUs. Although the data is not copied between domUs and dom0, the dom0 needs to schedule the I/O operations on behalf of the domUs. Figure 18(left) and figure 18 (right) shows this behavior in the 1-VM per node and 8-VMs per node configurations that we used.

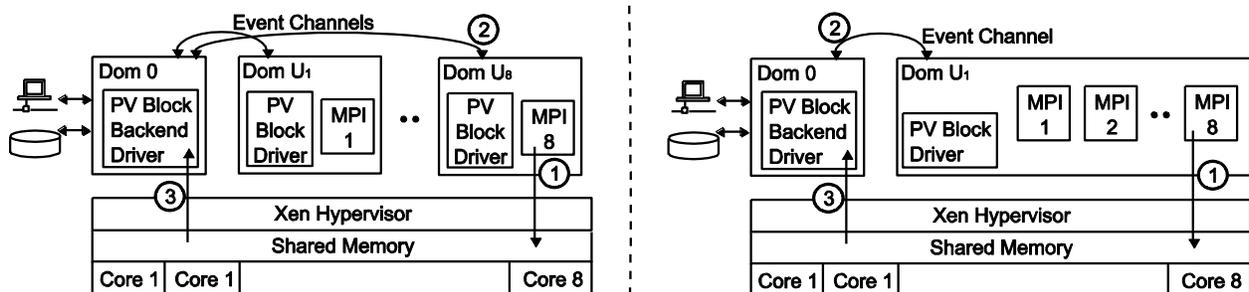


Figure 18. Communication between dom0 and domU when 1-VM per node is deployed (left). Communication between dom0 and domUs when 8-VMs per node were deployed (right).

In all the above parallel applications we tested, the timing figures measured correspond to the

time for computation and communication inside the applications. Therefore, all the I/O operations performed by the applications are network-dependent. From figure 19 (right), it is clear that Dom0 needs to handle 8 event channels when there are 8-VM instances deployed on a single bare-metal node. Although the 8 MPI processes run on a single bare-metal node, since they are in different virtualized resources, each of them can only communicate via Dom0. This explains the higher overhead in our results for 8-VMs per node configuration. The architecture reveals another important feature as well - that is, in the case of the 1-VM per node configuration, when multiple processes (MPI or other) that run in the same VM communicate with each other via the network, all the communications must be scheduled by the dom0. This results in higher latencies. We could verify this by running the above tests with LAM MPI (a predecessor of OpenMPI, which does not have improved support for in-node communications for multi-core nodes). Our results indicate that, with LAM MPI, the worst performance for all the test occurred when 1-VM per node is used. For example, figure 19 shows the performance of Kmeans clustering under bare-metal, 1-VM, and 8-VMs per node configurations. This observation suggests that, when using VMs with multiple CPUs allocated to each of them for parallel processing, it is better to utilize parallel runtimes, which have better support for in-node communication.

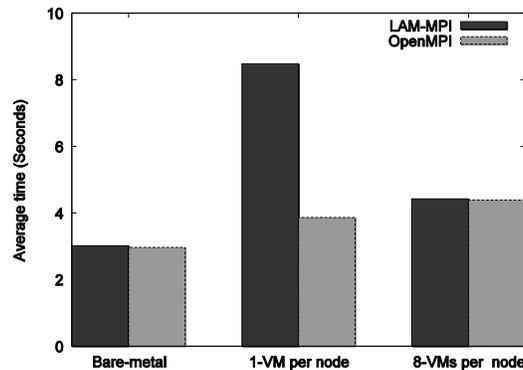


Figure 19. LAM vs. OpenMPI (OMPI) under different VM configurations

Several others have also performed relevant research on the performance implications of the virtualized resources. Youseff, Wolski, et al. (2006) presents an evaluation of the performance impact of Xen on MPI. According to their evaluations, the Xen does not impose considerable overheads for HPC applications. However, our results indicate that the applications that are more sensitive to latencies (smaller messages, lower communication to computation ratios) also experience higher overheads under virtualized resources, and this overhead increases as more and more VMs are deployed per hardware node. From their evaluations, it is not clear how many VMs they deployed on the hardware nodes, or how many MPI processes were used in each VM. According to our results, these factors cause significant changes in possible results. Running 1-VM per hardware node produces a VM instance with a similar number of CPU cores, such as in a bare-metal node. However, our results indicate that, even in this approach, if the parallel processes inside the node communicate via the network, the virtualization may produce higher overheads under the current VM architectures.

Evangelinos and Hill (2008) discuss the details of their analysis of the performance of HPC benchmarks on EC2 cloud infrastructure. One of the key observations noted in their paper is that both the OpenMPI and the MPICH2-nemesis show extremely large latencies, while the LAM MPI, the GridMPI, and the MPICH2-scok show smaller smoother latencies. This observation is similar to what we observed with the LAM-MPI in our tests, and the same explanation holds valid for their observation as well.

Walker (2008) presents benchmark results of the performance of HPC applications using “high CPU extra large” instances provided by EC2, and on a similar set of local hardware nodes. The local nodes are connected using infiniband switches; whereas Amazon EC2 network technology is unknown. The results indicate about a 40%-1000% performance degradation on

the EC2 resources, compared to the local cluster. Since the differences in operating systems and the compiler versions between the VMs and bare-metal nodes may cause variations in results, for our analysis, we used a cloud infrastructure over which we have complete control. In addition, we used similar software environments in both the VMs and the bare-metal nodes. In our results, we noticed that applications that are more susceptible to latencies experience higher performance degradation (around 40%) under virtualized resources. Bandwidth does not seem to be a consideration in private cloud infrastructures.

Gavrilvska, Kumar et al. (2007) discuss several improvements over the current virtualization architectures to support HPC applications such as HPC hypervisors and self-virtualized I/O devices. We notice the importance of such improvements and research. In our experimental results, we used hardware nodes with 8 cores, and we deployed and tested up to 8VMs per node in these systems. Our results show that the virtualization overhead increases with the number of VMs deployed on a hardware node. These characteristics will have a larger impact on systems having more CPU cores per node. A node with 32 cores running 32 VM instances may produce very large overheads under the current VM architectures.

7. Conclusions and Future Work

We have described several different studies of clouds and cloud technologies on both real applications and standard benchmark. These address different aspects of parallel computing using either traditional (MPI) or the new cloud inspired approaches. We find that cloud technologies work well for most pleasingly-parallel problems (“Map-only and “Map-reduce” classes of applications). In addition, their support for handling large data sets, the concept of moving computation to data, and the better quality of services provided such as fault tolerance and monitoring, all serve to simplify the implementation details of such problems. Applications

with complex communication patterns observe higher overheads when implemented using cloud technologies, and even with large data sets, these overheads limit the usage of cloud technologies for such applications. Enhanced MapReduce runtimes such as CGL-MapReduce allows iterative style applications to utilize the MapReduce programming model, while incurring minimal overheads, as compared to the other runtimes such as Hadoop and Dryad.

Handling large data sets using cloud technologies on cloud resources is an area that needs more research. Most cloud technologies support the concept of moving computation to data where the parallel tasks access data stored in local disks. Currently, it is not clear to us how this approach would work well with the VM instances that are leased only for the duration of use. A possible approach is to stage the original data in high performance parallel file systems or Amazon S3 type storage services, and then move data to the VMs each time they are leased to perform computations.

MPI applications that are sensitive to latencies experience moderate-to-higher overheads when performed on cloud resources, and these overheads increase as the number of VMs per bare-hardware node increases. For example, in Kmeans clustering, 1-VM per node shows a minimum of an 8% total overhead, while 8-VMs per node shows at least a 22% overhead. In the case of the Concurrent Wave Equation Solver, both these overheads are around 50%. Therefore, we expect the CPU core assignment strategies, such as $\frac{1}{2}$ of a core per VM, to produce very high overheads for applications that are sensitive to latencies. Applications that are not susceptible to latencies, such as applications that perform large data transfers and/or higher Communication/Computation ratios, show minimal total overheads in both bare-metal and VM configurations. Therefore, we expect that the applications developed using cloud technologies will work fine with cloud resources, because the milliseconds-to-seconds latencies that they

already have under the MapReduce model will not be affected by the additional overheads introduced by the virtualization. This is also an area we are currently investigating. We are also building applications (biological DNA sequencing) whose end-to-end implementation from data processing to filtering (data-mining), involves an integration of MapReduce and MPI (Fox, Bae et al. 2008).

Acknowledgements

We would like to thank Joe Rinkovsky and Jenett Tillotson from IU UITs for their dedicated support in setting up a private cloud infrastructure and helping us with various configurations associated with our evaluations. We would also want to thank ARTS team at Microsoft Research for their support on hardware and software infrastructures. We are grateful for Mina Rho and Haixu Tang from IU School of Informatics and Computing for their help in understanding Alu sequence clustering and providing human and chimpanzee gene sequence data.

References

- Amazon.com, Inc. 2009. Simple Storage Service (S3). <http://aws.amazon.com/s3>.
- ASF. 2009. Apache Hadoop Core. <http://hadoop.apache.org/core>.
- ASF. 2009. Apache Hadoop Pig. <http://hadoop.apache.org/pig/>.
- Barham, P., B. Dragovic, et al. 2003. Xen and the art of virtualization. *Proceedings of the nineteenth ACM symposium on Operating systems principles*. Bolton Landing, NY, USA.
- Batzer M.A., and P.L. Deininger, 2002, ALU repeats and human genomic diversity. *Nat. Rev. Genet.* 3 (5): 370-379.
- Condor Team. 2009. Condor DAGMan, <http://www.cs.wisc.edu/condor/dagman/>.
- Dean, J. and S. Ghemawat. 2008. MapReduce: simplified data processing on large clusters. *Commun. ACM* 51(1): 107-113.
- Dongarra, J., C. A. Geist, et al. 1993. Integrated PVM Framework Supports Heterogeneous Network Computing. *Computers in Physics*: 166-175.

- Ekanayake, J., S. Pallickara, et al. (2008). MapReduce for Data Intensive Scientific Analyses. eScience, 2008. IEEE Fourth International Conference on eScience '08
- ElasticHosts Ltd. 2009. Cloud Hosting. <http://www.elastichosts.com/>.
- Ericsson 2009. Erlang programming language. <http://www.erlang.org/>.
- Evangelinos, C. and C. Hill. 2008. Cloud Computing for parallel Scientific HPC Applications: Feasibility of Running Coupled Atmosphere-Ocean Climate Models on Amazon's EC2, *The First Workshop on Cloud Computing and its Applications (CCA '08)*. Chicago, IL.
- Forum, MPI. n.d.. MPI (Message Passing Interface). <http://www.mcs.anl.gov/research/projects/mpi/>.
- Foster, I. 2001. The anatomy of the grid: enabling scalable virtual organizations. *International Journal of Supercomputer Applications* 15: 2001.
- Fox, G. C., Hey, A. and Otto, S., Matrix Algorithms on the Hypercube I: Matrix Multiplication, *Parallel Computing*, 4, 17 (1987)
- Fox, G., S. Bae, et al. 2008. Parallel Data Mining from Multicore to Cloudy Grids. Proc. of *International Advanced Research Workshop on High Performance Computing and Grids*, HPC2008, Cetraro, Italy.
- Gabriel, E., G.E. Fagg, et al. 2004. Open MPI: Goals, concept, and design of a next generation MPI implementation. In Proc of, *11th European PVM/MPI Users' Group Meeting*, Budapest, Hungary.
- Gavrilovska, A., S. Kumar, et al. 2007. Abstract High-Performance Hypervisor Architectures: Virtualization in HPC Systems. Proc. of *HPCVirt 2007*, Portugal, Mar 2007.
- Ghemawat, S., H. Gobiuff, et al. 2003. The Google file system. *SIGOPS Oper. Syst. Rev.* **37**(5): 29-43.
- Gotoh, O. 1982. An improved algorithm for matching biological sequences. *Journal of Molecular Biology* 162:705-708
- Gu, Y. and R. L. Grossman. 2009. Sector and Sphere: the design and implementation of a high-performance data cloud. *Philos Transact A Math Phys Eng Sci* **367**(1897): 2429-45.
- Huang, X. and A. Madan. 1999. CAP3: A DNA sequence assembly program. *Genome Res* **9**(9): 868-77.
- Hull, D., K. Wolstencroft, et al. 2006. Taverna: a tool for building and running workflows of services. *Nucleic Acids Res* **34**(Web Server issue): W729-32.
- Isard, M., M. Budiu, et al. 2007. Dryad: distributed data-parallel programs from sequential building blocks. *SIGOPS Oper. Syst. Rev.* **41**(3): 59-72.
- Vermorel, J. 2005. NAligner (Smith Waterman software with Gotoh enhancement). <http://jaligner.sourceforge.net/naligner/>.
- Johnsson, S. L., T. Harris, et al. 1989. Matrix multiplication on the connection machine. Proc of the *1989 ACM/IEEE conference on Supercomputing*. Reno, Nevada, United States, ACM.

- Jurka, J. 2000. Repbase Update: a database and an electronic journal of repetitive elements. *Trends Genet.* 9:418-420 (2000).
- Keahey, K., I. Foster, et al. 2005. Virtual workspaces: Achieving quality of service and quality of life in the Grid. *Sci. Program.* **13**(4): 265-275.
- Ludscher, B., I. Altintas, et al. 2006. Scientific workflow management and the Kepler system: Research Articles. *Concurr. Comput. : Pract. Exper.* **18**(10): 1039-1065.
- Macqueen, J. (1967). Some methods for classification and analysis of multivariate observations. Proc. of *Fifth Berkeley Symp. on Math. Statist. and Prob.*
- Moretti, C., H. Bui, et al. 2009. All-Pairs: An Abstraction for Data Intensive Computing on Campus Grids. *IEEE Transactions on Parallel and Distributed Systems.* **99**(1)
- Nokia. 2009. Disco project. <http://discoproject.org/>.
- Nurmi, D., R. Wolski, et al. 2009. The Eucalyptus Open-Source Cloud-Computing System. Proc. of *9th IEEE/ACM International Symposium on Cluster Computing and the Grid, 2009*.
- Pallickara, S. and G. Fox. 2003. NaradaBrokering: a distributed middleware framework and architecture for enabling durable peer-to-peer grids. Proc of the *ACM/IFIP/USENIX 2003 International Conference on Middleware*. Rio de Janeiro, Brazil, Springer-Verlag New York, Inc.
- Pallickara, S. L. and M. Pierce. 2008. SWARM: Scheduling Large-Scale Jobs over the Loosely-Coupled HPC Clusters. Proc of *IEEE Fourth International Conference on eScience '08(eScience, 2008)*. Indianapolis, USA
- Pike, R., S. Dorward, et al. 2005. Interpreting the data: Parallel analysis with Sawzall. *Sci. Program.* **13**(4): 277-298.
- Raicu, I., Y. Zhao, et al. 2007. Falkon: a Fast and Light-weight task executiON framework. *Proceedings of the 2007 ACM/IEEE conference on Supercomputing*. Reno, Nevada, ACM.
- ServePath. 2009. GoGrid Cloud Hosting. <http://www.gogrid.com/>.
- Smit, A. F. A., R. Hubley, and P. Green. 2004. Repeatmasker. <http://www.repeatmasker.org>
- Smith, T.F. and Waterman, M.S. 1981. Identification of common molecular subsequences. *Journal of Molecular Biology* 147:195-197
- Smith Waterman Software. <http://jaligner.sourceforge.net/naligner/>.
- Walker, E. 2008. Benchmarking Amazon EC2 for high-performance scientific computing. <http://www.usenix.org/publications/login/2008-10/openpdfs/walker.pdf>.
- Y.Yu, M. Isard, et al. 2008. DryadLINQ: A System for General-Purpose Distributed Data-Parallel Computing Using a High-Level Language. Proc of *Symposium on Operating System Design and Implementation (OSDI)*. San Diego, CA.

Youseff, L., R. Wolski, et al. 2006. Evaluating the Performance Impact of Xen on MPI and Process Execution For HPC Systems. Proc of *First International Workshop on Virtualization Technology in Distributed Computing*, 2006.