
On the Scalability of Multi-Criteria Protein Structure Comparison in the Grid

GIANLUIGI FOLINO^{1*}, AZHAR ALI SHAH^{**}, AND NATALIO KRASNOGOR^{***}

RECEIVED ON 10.06.2011 ACCEPTED ON 01.10.2011

ABSTRACT

MC-PSC (Multi-Criteria Protein Structure Comparison) is one of the GCAs (Grand Challenge Applications) in the field of structural proteomics. The solution of the MC-PSC grand challenge requires the use of distributed algorithms, architectures and environments. This paper is aimed at the analysis of the scalability of our newly developed distributed algorithm for MC-PSC in the grid environment. The scalability in the grid environment indicates the capacity of the distributed algorithm to effectively utilize an increasing number of processors across multiple sites. The results of the experiments conducted on the UK's NGS (National Grid Service) infrastructure are reported in terms of speedup, efficiency and cross-site communication overhead.

Key Words: Protein Structure Comparison, Grid, Scalability, Bioinformatics

1. INTRODUCTION

The theoretical analysis of the scalability of the 'computation-centric' parallel applications on the grid appears in [1] with a prompt to the Grid community for the demonstration of this idea in terms of real Grid computing environments. This theoretical analysis is based on the idea of HCG (Homogeneous Computational Grid) and fits well with the real Grid computing infrastructure provided by the UK NGS [2] (Please see Section 3 for the details of the NGS infrastructure). The HCG model is based on the concept of HRM (Hierarchical Resource Manager) [3] and assumes that the Grid consists of C number of identical CEs (Computing Elements) and each CE (being a HPC

system) has p number of identical processors connected using the same type of network. The workload decomposition on such a system consists of two-level hierarchy: at first the un-decomposed work (W expressed e.g. in Mflops) is equally distributed in C CE's (i.e W/C decomposition) and then within each CE the portion of the work is assigned to each of the p processors (i.e $(W/C)/p$ decomposition). Consequently, this two-level hierarchy gives rise to two sources of communication overhead, i.e. the communication overhead among C CE's $Q_1(W,C)$ and the communication overhead among p processors of each CE $Q_2(W/C,p)$. With this formalism, the execution time on HCG could be defined as:

* Institute of High Performance Computing and Networking, Italian National Research Council, Cosenza 87036, Italy.
** Institute of Information & Communication Technology, University of Sindh, Jamshoro, Pakistan.
*** School of Computer Science, University of Nottingham, NG81BB, UK.

$$T_{C,p}(W) = \frac{W_p}{pC\Delta} + Q_2\left(\frac{W}{C,p}\right) + Q_1(W,C) \quad (1)$$

Where Δ indicates the computing capacity of a processor e.g Mflops/s. Please note that if $C=1$ and if $Q_1(W,1)=0$ then the overhead of Equation (1) returns to the standard parallel case i.e $Q_2(W,p)=Q(W,p)$.

Equation (1) makes it clear that running the parallel application on more than one CEs introduces an additional communication overhead in terms of $Q_1(W,C)$ which increases the execution time. However, this increase in the execution time could be masked by the value of C , which decreases the execution time by increasing the number of processors and also by reducing the communication overhead in terms of $Q_2(W/C,p)$ as compared to $Q(W,p)$ on one CE.

In order to analyze the added value of parallelism we normally compare the parallel execution time on P processors with the sequential execution time on 1 processor. However, as suggested by [1], in a Grid environment, we need to compare the parallel execution time on C CE's with the parallel execution time on 1 CE. This comparison is named as Grid Speedup and is mathematically defined as:

$$\Gamma_p^C = \frac{T_1, p(W)}{T_{C,p}(W)} \quad (2)$$

where Γ_p^C is the 'Grid Speedup' (with p processors and C CEs), T_1 is the execution time on a single CE and T_c is the execution time on C CEs.

The Grid Speedup (Equation (2)) is one of the scalability metrics for the parallel applications on the Grid. Its value

indicates how better a parallel application performs when decomposed on C CEs as compared to its performance on a single CE in terms of execution time. From Equation (2) we could also derive the expression for the Grid efficiency as:

$$\gamma_p^C = \frac{T_1, p(W)}{CT_{C,p}(W)} \quad (3)$$

where γ_p^C is the 'Grid efficiency' and p, C, T_1 and T_c represent the same parameters as described in Equation (2).

The description of the 'Grid Efficiency' in Equation (3) follows Amdahl's popular statement that "for a given instance of a particular problem, the system efficiency decreases when the number of available processors is increased" [4]. In the case of the Grid efficiency, in addition to the number of processors, it is the value of the C (number of CEs) that affects the system efficiency. Based on these concepts of scalability, this paper performs empirical analysis of our parallel algorithm for MC-PSC as described in the following sections.

The remainder of this paper is organized as follows: Section 2 describes the background related to the MC-PSC Grand Challenge, Section 3 describes the experimental setup; Section 4 presents the results and discussions and finally Section 5 concludes the paper.

2. THE MC-PSC GRAND CHALLENGE

The problem of large scale MC-PSC could be represented as a 3D cube. The x and y axis of the cube representing the different proteins being compared, while the z axis representing different comparison methods being used such as the USM (Universal Similarity Metric), MaxCMO (Maximum

Contact Map Overlap), DaliLite (Distance Alignment Matrix), CE (Combinatorial Extension), FAST (Fast Alignment And Search Tool) and TM-Align etc. While processed, each cell of this 3D cube holds the output of each comparison method in terms of different measures and metrics. That is, each cell of the 3D cube represents both the processing as well as the storage perspective of the problem space while cell boundaries specify the communication overhead. Given the ever growing number of protein structure comparison methods as well as the number of protein structures being deposited in the PDB; the dimensions of this cube go on increasing and making its computation, in our opinion, to be one of the GCAs (Grand Challenge Applications) in the field of structural biology. GCAs are defined as "fundamental problems in science and engineering with great economic and scientific impact, whose solution is intractable without the use of state-of-the-art parallel/distributed systems" [5]. Many examples of the use of parallel/distributed systems for the solution of GCAs in the field of life sciences in general and structural proteomics in particular are available in the literature [6]. It is believed that most of the GCAs may have several parallel solutions; therefore, a methodological approach based on an exploratory nature will help in finding the best available solution [7]. An example of such approach that is widely used for the design of parallel and distributed algorithms is the PCAM (Partitioning, Communication, Agglomeration, and Mapping) approach. Introduced by Foster, [7], the beauty of this approach is that it enables the designer to consider the machine-independent issues (e.g. concurrency, scalability and communication) first and machine-specific issues (e.g. granularity and load-balancing) later in the design

process. Based on the philosophy of the PCAM approach, a high-throughput distributed framework for the solution of the grand challenge of MC-PSC using MPI (Message Passing Interface) model of parallel programming has been introduced [8]. The performance of this framework along with different load balancing strategies was evaluated on a 64-node cluster as reported in [8]. However, it was observed that for datasets having relatively large number of proteins (e.g. 1000+), even the 64-node cluster becomes a bottleneck and it takes about 11 days for the computation to complete. Hence, we tried to deploy our algorithm on the UK NGS to take advantage of greater number of cores available across multiple sites. The deployment on the NGS is reported in the next section.

3. DEPLOYMENT ON THE NGS INFRASTRUCTURE

The NGS, provides the eScience infrastructure to all the UK-based scientists free of cost [2]. For our case we used the Globus-based MPIg [9] (grid-based implementation of MPI) to spawn the jobs across two NGS sites; one at Leeds and the other at Manchester. Like its predecessors (e.g. MPICH-G and MPICH-G2), the MPIg library extends the Argonne MPICH implementation of MPI to use services provided by the Globus Toolkit for cross-site job execution using IP-based communication for inter-cluster messaging. However, being the latest implementation, the MPIg includes several performance enhancements such as in the case of inter-cluster communication it uses multiple threads as compared to the single thread communication of the previous implementations. Furthermore, besides being backward compatible with

the pre-web service Globus, the MPIg also makes use of the new web services provided by Globus version 4x. By making use of the new web services, the MPIg provides much more enhanced functionality, usability and performance. The use of the MPIg for cross-site runs requires advanced resource reservation so that jobs (processes) can run simultaneously across all the sites. To facilitate this, NGS provides the HARC (High-Available Resource Co-allocation) [10] as a command line utility to perform automatic reservation. Each of the two NGS sites (Leeds and Manchester) consists of 256 cores (AMD Opteron with 2.6GHz and 8GB of main memory) interconnected with Myrinet M3F-PCIXD-2. However, the NGS policies allow the advance reservation of maximum of 128 cores at each site for the maximum duration of 48 hours. Once the reservation is done, then the Globus-based job submission could be achieved with the RSL (Resource Specification Language) scripts and other Globus services could be used for job monitoring and control. For the MPI based jobs to run on different sites, the source code of the application needs to be compiled with MPIg libraries at each site and the executable placed in the appropriate working directory under the respective local file system. The compilation of the MPI based application with MPIg does not require any change in the source code and hence from the user's perspective the deployment is as straight forward as running the parallel application on a single site/cluster with the exception that the RSL scripts specifies the resources of the additional site to be used. Fig. 1, shows the overall architecture and setup of deploying the MC-PSC application on the Grid.

3.1 Dataset

The dataset used in these experiments is the one introduced by Kinjo, et. al. [11] consisting of 1012 non-redundant protein chains having a total of 252, 569 residues. The 1012 chains result in as many as 1,024, 144 pairwise comparisons for each method/algorithm. While using all the six methods (i.e. USM, MaxCMO, CE, DaliLite, FAST and TM-Align), the total number of pairwise comparisons becomes $1,024, 144 \times 6 = 6, 144, 864$. Given that the average time for the comparison of 1 pair using all the six methods on a single processor machine is about 8 secs, this computation requires about 569 days to complete on a single processor and it took about 10.7 days to complete on a 64-node cluster [10]. The results on the Grid infrastructure are presented in the next section.

4. RESULTS AND DISCUSSION

Both the single-site and cross-site experiments for MC-PSC were conducted with varying number of processors using the Kinjo, et. al. [11] dataset. The Grid speedup and efficiency (Equations (2-3) respectively) were calculated based on the results of these experiments and are shown in Figs. 2-3 . Fig. 2 shows that initially (for less number of processors), running the MC-PSC experiments across two sites almost doubles the performance to that of the single-site. However, as the number of processors increases (thereby decreasing the level of granularity and increasing the communication overhead), the speedup decreases slightly and finally reaches to about 1.65. There is also same trend in the Grid efficiency as shown in Fig. 3.

Figs. 4-5 provide the comparison of the algorithmic speedup on a single-site (S_1 , having 128 processors) and the speedup obtained while running the experiments on the two sites (S_2 , having a total of 256 processors). The speedup in this case is taken as the ratio of the execution time on single-machine (single processor) (T_1) to the execution time on p processors (T_p) (i.e $S_1=S_2=T_1T_p$). As indicated by Fig. 4 though initially, the cross-site speedup is slightly low as

compared to the single-site speedup; however, given the large number of processors available on the later, the overall speedup increases by almost a factor of 2. The total time for the computation of the given dataset on 256 cores (2.4GHz each) was reduced to 38.6 hours. Comparing this with the 569 days on the single-machine and 10.7 days required on a 64-node (though having less processor power i.e 1.4GHz each) cluster we observe a good scalability and performance of our

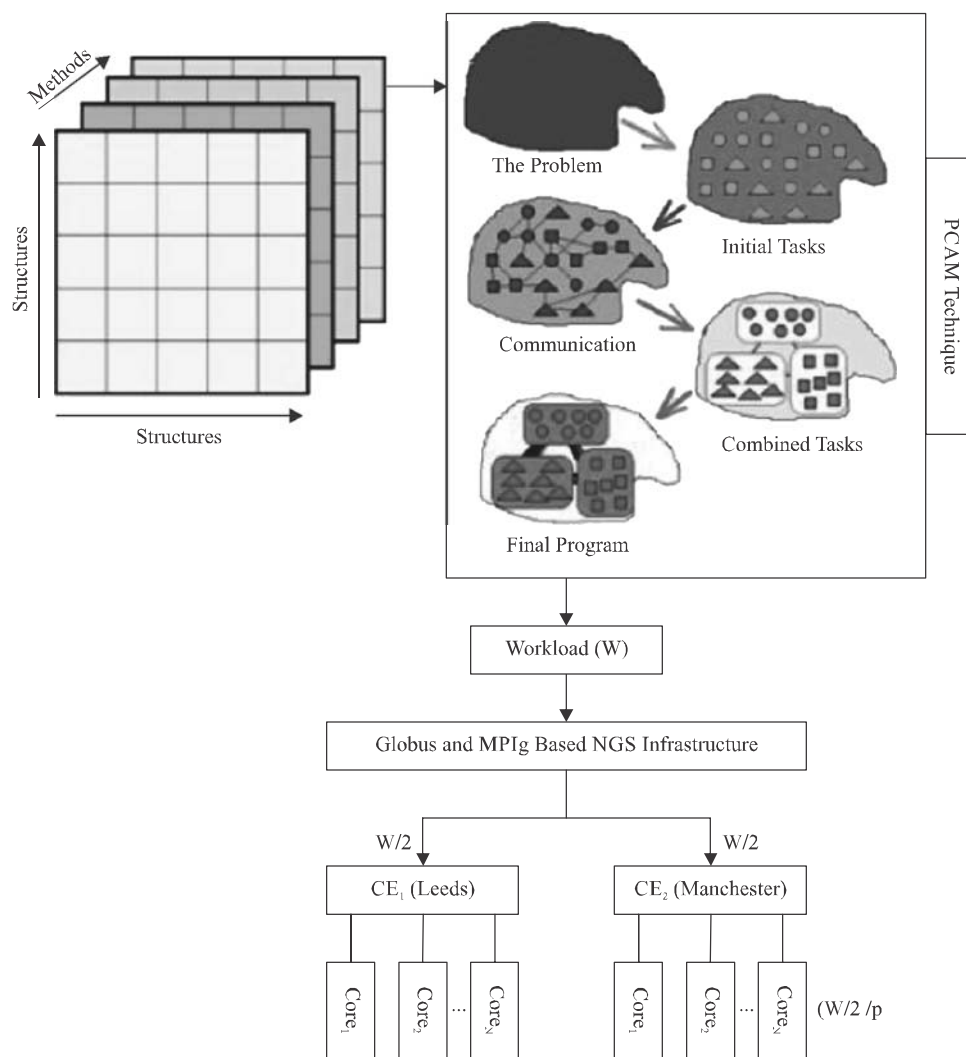


FIG. 1. DEPLOYMENT OF THE MC-PSC APPLICATION ON THE GRID: THE APPLICATION TAKES PROTEIN 3D STRUCTURES AS INPUT AND PREPARES THE BALANCED WORKLOAD W TO BE DISTRIBUTED ON THE GRID. HALF OF THE TOTAL WORKLOAD ($W/2$) IS ASSIGNED TO EACH SITE (CE). EACH SITE FURTHER DISTRIBUTES THE $W/2$ INTO p NUMBER OF CORES

algorithm on the Grid. The boost in the speedup and performance is two folds i.e the large number of processors (physical speedup) coupled with high speed of each individual processor (power scalability). Fig. 5 shows the corresponding efficiency of the algorithm on single-site and cross-site architecture. The efficiency, in this case measures the effective use

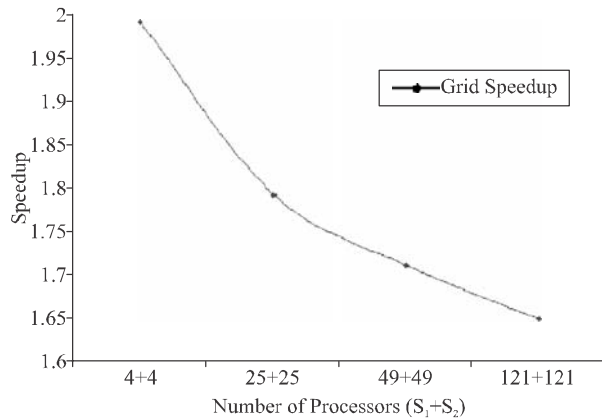


FIG. 2. PERFORMANCE OF THE MC-PSC ON THE GRID: GRID SPEEDUP; INITIALLY THE SPEEDUP IS ALMOST IDEAL FOR LESS NUMBER OF NODES BUT AS THE NUMBER OF NODES INCREASES ON EACH SITE THE CORRESPONDING LEVEL OF GRANULARITY DECREASES WHILE THE LEVEL OF COMMUNICATION OVERHEAD INCREASES AND HENCE IT CAUSES THE SPEEDUP TO DEGRADE SLIGHTLY. NEVERTHELESS, THE OVERALL SPEEDUP IS MUCH GREATER (1.6) AS COMPARED TO SPEEDUP ON THE SINGLE SITE (<1)

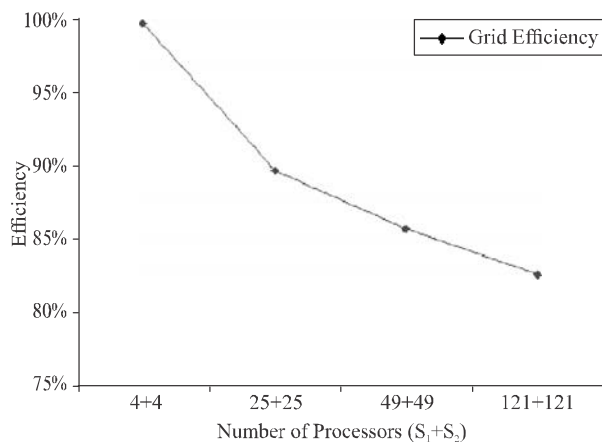


FIG. 3. GRID EFFICIENCY; AS EXPECTED THE SLIGHT DEGRADATION OF SPEEDUP CAUSES THE DEGRADATION IN THE EFFICIENCY OF THE SYSTEM

of the hardware and is equal to the ratio of the speedup on p processors to p (i.e $E=S_p/p$). Fig. 6 shows the cross-site communication overhead in terms of running the MC-PSC application in the Grid. It shows that, when a few processors are used, the load of the processors and the amount of data to be exchanged is high and consequently there is a considerable communication overhead. However, when we use a larger number of processors, the overhead is negligible in comparison with the computation time.

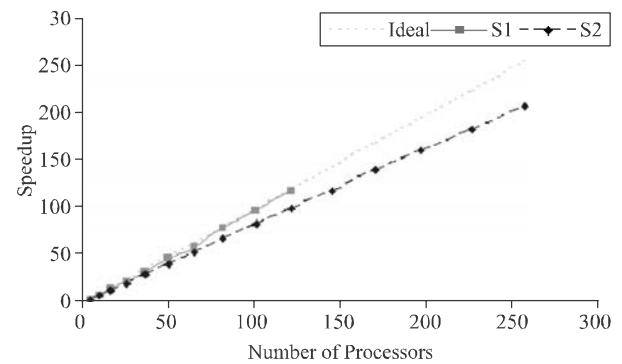


FIG. 4. SINGLE-SITE AND CROSS-SITE: SPEEDUP; THE GRAPH SHOWS THAT THOUGH INITIALLY, THE CROSS-SITE SPEEDUP (S2) IS SLIGHTLY LOW AS COMPARED TO THE SINGLE-SITE SPEEDUP (S1); HOWEVER, GIVEN THE LARGE NUMBER OF PROCESSORS AVAILABLE ON THE LATER, THE OVERALL SPEEDUP (S2) INCREASES BY ALMOST A FACTOR OF 2

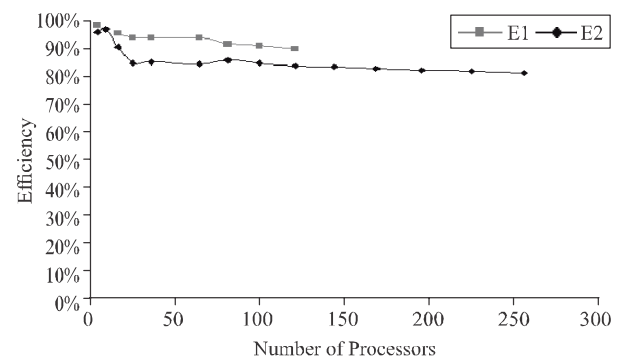
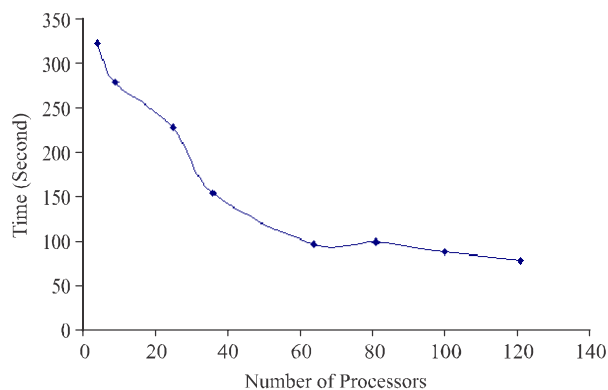


FIG. 5. EFFICIENCY; AS EXPECTED THE CROSS-SITE EFFICIENCY (E2 IS SLIGHTLY LESS AS COMPARED TO THE SINGLE-SITE EFFICIENCY (E1 DUE TO EXTRA COMMUNICATION OVERHEAD

FIG. 6. CROSS-SITE COMMUNICATION OVERHEAD. THE



GRAPH SHOWS THAT WHEN A FEW PRO-CESSORS ARE USED THE LOAD OF THE PROCESSORS AND CONSEQUENTLY THE AMOUNT OF DATA TO BE EXCHANGED IS HIGH AND CONSEQUENTLY THERE IS CONSIDERABLE COMMUNICATION OVERHEAD. HOWEVER, WHEN WE USE A LARGER NUMBER OF PROCESSORS, THE OVERHEAD IS NEGLIGIBLE IN COMPARISON WITH THE COMPUTATION TIME

5. CONCLUSION

The quality of our parallel algorithm for MC-PSC has been measured in terms of Grid Speedup and efficiency. The results of the single-site and cross-site experiments indicate that by making use of the Grid resources, the algorithm scales well and that the cross-site communication overhead is not much significant. The current cross-site experiments were conducted on only two sites based on the HCG model of the NGS, UK. As the NGS is still in the process of adding more sites, in future we would like to extend this study by increasing the number of sites as well as incorporating the heterogeneous architecture of the Grid. Because, at present the maximum time allocated for continuous execution of a job/process at NGS is limited to 48 hours and hence does not allow evaluating the performance of the application with very larger datasets, hence the software developed so far could be upgraded by adding the fault tolerance mechanism in the form of checkpoint/restart. The checkpoint/restart mechanism could be added without changing the code of the application by using some libraries such as the BLCR (Berkeley Lab Checkpoint/Restart). With these

improvements, it would be possible for the MC-PSC to perform real time computation with even large datasets and to develop a database of pre-computed results.

Acknowledgments

the authors would like to acknowledge the use of the National Grid Service, UK, in carrying out this work. Second Author acknowledges the University of Sindh, Jamshoro, Pakistan, for the scholarship SU/PLAN/F.SCH/794.

REFERENCES

- [1] Hoekstra, A.G., and Sloot, P.M.A., "Introducing Grid Speedup G: A Scalability Metric for Parallel Applications on the Grid", EGC, pp. 245-254, 2005.
- [2] Richards, A., and Sinclair, G. M., "UK National Grid Service", CRC Press, 2009.
- [3] Halderen, A.W., Overeinder, B.J., Sloot, P.M.A., Van Dantzig, R., Epema, D.H.J., and Livny, M., "Hierarchical Resource Management in the Polder Metacomputing Initiative", Parallel Computing, Volume 24, pp.12-13, 1998.
- [4] Amdahl, G., "Validity of the Single Processor Approach to Achieving Large-Scale Computing Capabilities", Proceedings of AFIPS Conference, Volume 30, pp. 483-485, 1967.
- [5] Silva, L., and Buyya, R., "Parallel Programming Models and Paradigms", Prentice Hall PTR, NJ, USA, 1999.
- [6] Shah, A.A., Barthel, D., Lukasiak, P., Blacewicz, J., and Krasnogor, N., "Web and Grid Technologies in Bioinformatics, Computational Biology and Systems Biology: A Review", Current Bioinformatics, Volume 3, No. 1, pp. 10-23, 2008.
- [7] Foster, I., "Parallel Computers and Computation", Designing and Building Parallel Programs: Concepts and Tools for Parallel Software Engineering, 1995.

- [8] Shah, A.A., Folino, G., and Krasnogor, N., "Towards High-Throughput, Multi-Criteria Protein Structure Comparison and Analysis", *IEEE Transactions on Nano Bioscience*, Volume 9, pp. 1-12, 2010.
- [9] Manos, S., Mazzeo, M., Kenway, O., Coveney, P.V., Karonis, N.T., and Toonen, B., "Distributed MPI Cross-Site Run Performance Using MPIG", *Proceedings of the 17th International Symposium on High Performance Distributed Computing*, New York, NY, USA, ACM, 2008.
- [10] Maclaren, J., Keown, M.M., and Pickles, S., "Co-Allocation Fault Tolerance and Grid Computing", *e-Science AHM*, UK, 2006.
- [11] Kinjo, A.R., Horimoto, K., and Nishikawa, K., "Predicting Absolute Contact Numbers of Native Protein Structure from Amino Acid Sequence", *Proteins Struct Funct Bioinf*, Volume 58, pp. 158-165, 2005.