

Modeling Variation of Performance Metric of Distributed Memory Heterogeneous Parallel Computer System Using Analytic and Recursive Models

Osondu E. Oguike¹, Monica N. Agu² and Stephenson C.Echezona³

Department of Computer Science
University of Nigeria
Nsukka, Enugu State
Nigeria

e-mail:

¹ osondu.oguike@unn.edu.ng

² monica.agu@unn.edu.ng

³ stephenson.echezona@unn.edu.ng

Abstract— In a heterogeneous parallel computer system, the computational power of each of the processors differs from one another. Furthermore, with distributed memory, the capacity of the memory, which is distributed to each of the processors, differs from one another. Using queuing system to describe a distributed memory heterogeneous parallel computer system, each of the heterogeneous processors will have its own heterogeneous queue. The variation of a performance metric of heterogeneous parallel computer system with distributed memory needs to be modeled because it will help designers of parallel computer system to determine the extent of variation of a performance metric. It will also help users to know when to realize minimum variation of a performance metric. This paper models the variation of a performance metric of distributed memory heterogeneous parallel computer system using analytic and recursive models.

Keywords - heterogeneous parallel computer, distributed memory, parallel computer system, queuing network, variation, recursive model, analytic models

I. INTRODUCTION

A heterogeneous parallel computer system is one in which the computational power of each of the processors differs from one another. With distributed memory, it means that each of the heterogeneous processors has its own memory. Describing the system using queuing network, each of the processors has its own queue. With a round robin scheduling algorithm, processes can be scheduled to the various parallel processors, whenever a process needs to perform an I/O operation, it joins the appropriate I/O queue. Therefore, the queuing network of a heterogeneous parallel computer system consists of parallel processors, parallel processor queues, I/O processors and I/O queues. Suppose there are n different parallel processor queuing systems and k different I/O queuing systems. A queuing system in this context is defined as a processor, together with its own queue. We also assume that the various queues are finite [1, 2, 3, 4] i.e. there is a limit to the number of jobs that can be admitted into the queues, and negligible communication overhead. Suppose $X_1, X_2, X_3, \dots, X_n, X_{n+1}, X_{n+2}, X_{n+3}, \dots, X_{n+k}$ are the maximum number of processes that can be admitted into the respective queues. We assume that processes arrive at the various queues according to Poisson distribution, and they are serviced according to exponential distribution [5, 6]. Figure 1 illustrates a model of the queuing network of a heterogeneous parallel computer system with distributed memory.

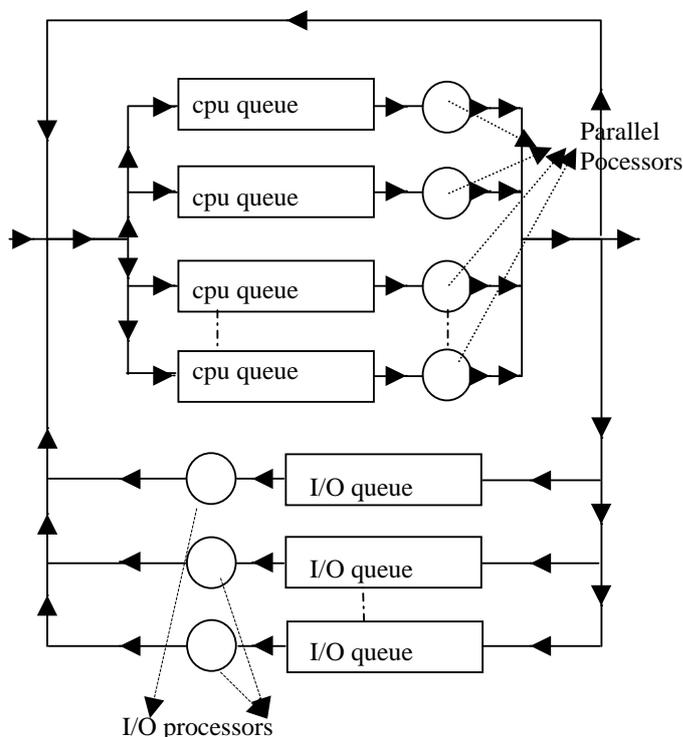


Figure 1: Queuing network of a heterogeneous parallel computer system with distributed memory.

II. STATEMENT OF THE PROBLEM

Though analytic queuing method can be used to model the performance of a heterogeneous parallel computer system with distributed memory, however, it may not be possible to model some performance metrics of heterogeneous parallel computer, like variation of a performance metric. The reason is because the analytic method cannot determine the exact convergence of some mathematical series that are used in modeling the variation of a performance metric of a heterogeneous parallel computer system. Therefore, there is the need for the use of another model, rather than analytic model. The use of efficient linear recursive model [9] can efficiently model the variation of a performance metric of a distributed memory heterogeneous parallel computer system, because recursive models can be used to determine the exact convergence of any series used in modeling the variation of a performance metric of distributed memory parallel computer system.

III. LITERATURE REVIEW

Queuing approach has been used extensively in the literature to model the performance of computer systems. However, this has been done in different ways and for different models of computer systems. In [20], the authors used a recursive computation approach to solve the steady state equations, thereby leading to the modeling of the various performance metrics of a multi-terminal system that is subject to breakdown. Furthermore, the author in [24] used a rigorous approach to model the performance of heterogeneous parallel computer system without introducing any constraint on the kind of interconnection between the heterogeneous nodes. Furthermore, in [24], systems with the same interconnection speed were considered when modeling the performance of heterogeneous parallel computer system. The authors in [25] looked at alternative ways of measuring the performance of heterogeneous parallel computer system, by modeling linear speed and linear efficiency using simulation-modeling techniques. In [26], the author showed that Little's formulae could be universally applicable, if properly interpreted to take account of state-varying entrance rates, batch arrivals, and multiple customer classes. In [27], the author confirmed that Little's formula could be applied to very general queuing systems (not just M/M/1), even whole networks! The authors in [28] considered a new performance metric, variation of the computing power as a unique performance metric that is ideal for a heterogeneous network of workstations, though an approach different from queuing approach was used to do this. In [29], analytic models were used to model the performance of computer intensive applications of parallel computers, while [30] used recursive models only to evaluate the performance of compute intensive application of a parallel computer system. In [31], recursive models were used to evaluate the performance of heterogeneous parallel computer system with distributed memory.

IV. METHODOLOGY

This paper models the variation of a performance metric of distributed memory heterogeneous parallel computer system. A queuing approach, with finite queues has been used to achieve the above aim, with parallel processors depicting parallel servers. The statistical method of probability density function and other probability theory concepts have used [15, 23]. A novel method of deriving the recursive model that determines the xth terms and the convergence of important mathematical series has been used to develop the recursive models. The simulation of the models on the computer has been done using Java programming language and the statistical regression/trend line analysis has been used to analyze the results of the simulation [11].

V. DEVELOPING THE MODELS

As a result of the use of the above methodologies, the following models have been developed for one queuing system and for all the queuing systems of the queuing network.

A. Models Based on a Queuing System

The following models have been developed for one queuing system

- Probability Density Function of the Number of Processes in a Queue.

Let X_i denotes the maximum number of processes that can be in the i th finite queuing system at any time [12, 13, 14]. Suppose the arrival rate, λ_{x_i} when x_i processes are in the i th queuing system of the queuing network be described as:

$$\lambda_{x_i} = \begin{cases} \lambda_i, & x_i = 0,1,2,3,\dots,X_i - 1 \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

Since the various processors are heterogeneous, therefore, it implies that the departure rate will vary, which can be described as:

$$\mu_{x_i} = \begin{cases} \mu_i, & x = 1,2,3,4,\dots, X_i \\ 0, & \text{otherwise} \end{cases} \quad (2)$$

Using the steady state probability as stated in [7, 16] the probability that x_i processes will be in the i th queuing system is

$$P_{x_i} = \begin{cases} \rho_i^{x_i} P_{0_i}, & x \leq X_i \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

The utilization factor for the i th queuing system, ρ_i is defined as:

$\frac{\lambda_i}{\mu_i}$. To obtain the value of P_{0_i} in equation (3), we sum all μ_i the probabilities for the i th queuing system and equate it to 1. This implies that:

$$\sum_{x_i=0}^{X_i} P_{x_i} = 1. \tag{4}$$

From equations (3) and (4), it implies that:

$$P_{0i} + \rho_i P_{0i} + \rho_i^2 P_{0i} + \rho_i^3 P_{0i} + \rho_i^4 P_{0i} + \dots + \rho_i^{X_i} P_{0i} = 1. \tag{5}$$

Factorizing equation (5), it implies that

$$P_{0i} (1 + \rho_i + \rho_i^2 + \rho_i^3 + \dots + \rho_i^{X_i}) = 1. \tag{6}$$

Analytically, the series in equation (6) converges to:

$$\frac{1 - \rho_i^{X_i+1}}{1 - \rho_i}, \text{ provided that } \rho_i \neq 1, \text{ otherwise it converges to } X_i + 1.$$

It implies that $P_{0i} \left(\frac{1 - \rho_i^{X_i+1}}{1 - \rho_i} \right) = 1$, when $\rho_i \neq 1$ and

$P_{0i}(X_i + 1) = 1$, when $\rho_i = 1$. Solving for P_{0i} , we have that:

$$P_{0i} = \begin{cases} \frac{1 - \rho_i}{1 - \rho_i^{X_i+1}}, & \rho_i \neq 1 \\ \frac{1}{X_i + 1}, & \rho_i = 1 \end{cases} \tag{7}$$

Using equation (7) in equation (3), we have the following:

$$P_{x_i} = \begin{cases} \frac{\rho_i^{x_i} (1 - \rho_i)}{1 - \rho_i^{X_i+1}}, & \rho_i \neq 1, \quad x_i = 0, 1, 2, 3, \dots, X_i \\ \frac{1}{X_i + 1}, & \rho_i = 1 \end{cases} \tag{8}$$

Equation (8) is the probability density function that models the probability that x_i processes will be admitted in the i th queuing system.

- Average Number of Processes in One Queuing System.

Furthermore, the average number of processes in the i th queuing system (i.e the queue and the processor) can be described statistically as expectation of x_i , where x_i is the random variable that denotes the number of processes in the i th queuing system. This can be written as

$$E(x_i) = \sum_{x_i=0}^{X_i} x_i P_{i x_i}. \tag{9}$$

Using equation (8) in equation (9), we obtain the following:

$$E(x_i) = \sum_{x_i=1}^{X_i} \left(\frac{x_i \rho_i^{x_i} (1 - \rho_i)}{1 - \rho_i^{X_i+1}} \right), \rho_i \neq 1, \quad x_i = 0, 1, 2, 3, \dots, X_i \tag{10}$$

Equation (10) can be simplified as:

$$E(x_i) = \left(\frac{1 - \rho_i}{1 - \rho_i^{X_i+1}} \right) \rho_i (1 + 2\rho_i + 3\rho_i^2 + 4\rho_i^3 + \dots + X_i \rho_i^{X_i-1}) \tag{11}$$

A recursive model has been used in [30,31] to determine the convergence of the series in equation (11). The recursive model is called $\text{Sum2}_i(X_i, \rho_i)$, and it is given as:

$$\begin{cases} 1, & X_i = 1 \\ \text{Term2}_i(X_i) * \text{term1}_i(X_i-1, \rho_i) + \text{Sum2}_i(X_i-1, \rho_i), & X_i \neq 1 \end{cases} \tag{12}$$

$\text{Term2}_i(X_i)$ and are given as:

$$\text{Term2}_i(X_i) = \begin{cases} 1, & X_i = 1 \\ 1 + \text{term2}_i(X_i-1), & X_i \neq 1 \end{cases} \tag{13}$$

and $\text{term1}_i(X_i, \rho_i)$ is defined below as:

$$\text{Term1}_i(X_i, \rho_i) = \begin{cases} 1, & X_i = 0 \\ \rho_i * \text{Term1}_i(X_i-1, \rho_i), & X_i \neq 0 \end{cases} \tag{14}$$

Therefore, using equation (12) in equation (11), we obtain:

$$E(x_i) = \left(\frac{1 - \rho_i}{1 - \rho_i^{X_i+1}} \right) \rho_i \text{Sum2}_i(X_i, \rho_i) \tag{15}$$

B. Models Based on The Whole Queuing Network.

Having developed the models for the performance metrics of one queuing system, these models can be extended to the whole queuing systems of the queuing network of a heterogeneous parallel computer system. It is necessary to define δ_i as the probability that a process will join the i th queue after each cpu burst, and δ_0 as the probability that the execution of a process has been completed. Arrival of processes into the various parallel processor queues can come from the outside world or from the various I/O queues or from the particular parallel processor, at the expiration of the time quantum for that process. Let λ_i be the rate of arrival of processes into the i th queuing system, and λ , the rate of arrival of processes from the outside world.. Under the steady state, when we

consider the queuing network, the overall utilization factor has been defined in [31] as:

$$\rho_i = \begin{cases} \frac{\lambda}{\delta_0 \mu_i}, & i = 0 \\ \frac{\lambda \delta_i}{\delta_0 \mu_i}, & i = 1, 2, 3, \dots, n + k \end{cases} \quad (16)$$

- Variation of Average Number of Processes in all the Queuing Systems of the Queuing Network

Suppose x_i is the random variable that denotes the number of processes in the i th queuing system. Therefore, another random variable, Y , can denote the average number of processes in the queuing systems of the queuing network, as:

$$Y = \frac{\sum_{i=1}^{n+k} x_i}{n+k} \quad (25)$$

Therefore, the variation of the average number of processes in all the queuing systems of the queuing network can be defined statistically as:

$$VAR(Y) = VAR\left(\frac{\sum_{i=1}^{n+k} x_i}{n+k}\right) \quad (26)$$

Using one of probability theory laws in [23], we obtain:

$$VAR(Y) = \frac{1}{(n+k)^2} \sum_{i=1}^{n+k} VAR(x_i) \quad (27)$$

From [23], the variance can be defined statistically as:

$$VAR(x_i) = E(x_i^2) - (E(x_i))^2 \quad (28)$$

Simplifying equation (28) further, we obtain:

$$E(x_i^2) = \sum x_i^2 P_{i,x_i} \quad (29)$$

Using equation (8) in equation (29), we obtain:

$$E(x_i^2) = \sum_{x_i=1}^{x_i} \left(\frac{x_i^2 \rho_i^{x_i} (1-\rho_i)}{1-\rho_i^{x_i+1}} \right) \quad (30)$$

Simplifying equation (30), we obtain:

$$E(x_i^2) = \left(\frac{1-\rho_i}{1-\rho_i^{x_i+1}} \right) \left(\rho_i^2 + 2\rho_i^3 + 3\rho_i^4 + 4\rho_i^5 + \dots + X_i^2 \rho_i^{X_i} \right) \quad (31)$$

Simplifying equation (31) further, we obtain:

$$E(x_i^2) = \left(\frac{1-\rho_i}{1-\rho_i^{x_i+1}} \right) \left(\rho_i + 4\rho_i^2 + 9\rho_i^3 + 16\rho_i^4 + \dots + X_i^2 \rho_i^{X_i} \right) \quad (32)$$

Factorising equation (32), we obtain:

$$E(x_i^2) = \left(\frac{1-\rho_i}{1-\rho_i^{x_i+1}} \right) \rho_i \left(1 + 4\rho_i + 9\rho_i^2 + 16\rho_i^3 + \dots + X_i^2 \rho_i^{X_i-1} \right) \quad (33)$$

The convergence of the series may be difficult or impossible to obtain analytically; therefore we seek for its convergence using recursive models. The same approach used earlier can be used to determine the convergence of the series,

$(1 + 4\rho_i + 9\rho_i^2 + 16\rho_i^3 + \dots + X_i^2 \rho_i^{X_i-1})$. The series can be considered as two sequences, which are: sequence1 = 1, 4, 9, 16, ..., X^2 , while the other sequence is: sequence2 = $(1, \rho_i, \rho_i^2, \rho_i^3, \dots, \rho_i^{X_i-1})$.

The recursive model that can be used to determine the x th terms of sequence1 can be obtained by adding $2X-1$, which is the common difference between the x th term and the $(x-1)$ th term to the $(x-1)$ th term of the sequence. The recursive model can be represented as shown below in equation (34), as:

$$\text{Term3}_i(X_i) = \begin{cases} 1, X_i = 1 \\ (2 * X_i - 1) + \text{Term3}_i(X_i - 1), X_i \neq 1 \end{cases} \quad (34)$$

The recursive model that determines the x th terms of sequence2 has been developed in equation (14). Therefore, combining equation (30) and equation (14), the series in equation (33) converges to this recursive model, called $\text{Sum3}_i(X_i, \rho_i)$, which is shown below as:

$$\begin{cases} 1, X_i = 1 \\ \text{Term1}_i(X_i - 1, \rho_i) * \text{Term3}_i(X_i) + \text{Sum3}_i(X_i - 1, \rho_i), X_i \neq 1 \end{cases} \quad (35)$$

Therefore, using equation (35) in equation (33), we obtain:

$$E(x_i^2) = \left(\frac{1-\rho_i}{1-\rho_i^{x_i+1}} \right) \rho_i \text{Sum3}(X_i, \rho_i) \quad (36)$$

Using equations (15) and (36) in equation (28), we obtain:

$$VAR(x_i) = \left(\frac{\rho_i(1-\rho_i)(\text{Sum3}(X_i, \rho_i))}{(1-\rho_i^{x_i+1})} \right) - \left(\frac{\rho_i(1-\rho_i)(\text{Sum2}(X_i, \rho_i))}{(1-\rho_i^{x_i+1})} \right)^2 \quad (37)$$

Therefore, using equation (37) in equation (27), we obtain:

$$VAR(Y) = \frac{1}{(n+k)^2} \left(\sum_{i=1}^{n+k} \left(\frac{\rho_i(1-\rho_i)(\text{Sum3}(X_i, \rho_i))}{(1-\rho_i^{x_i+1})} \right) - \left(\frac{\rho_i(1-\rho_i)(\text{Sum2}(X_i, \rho_i))}{(1-\rho_i^{x_i+1})} \right)^2 \right) \quad (38)$$

Equation (38) models the variation of the average number of processes in all the queuing systems of the queuing network by combining analytic and recursive models.

VI. RESULTS OF THE SIMULATION

The result of the simulation was analyzed to determine how variation of the performance metric under consideration changes as a particular parameter varies, while other parameters remain constant [10]. Figure 2 shows the result of the simulation, if the probability of a process leaving the system is known to be 0.2 and the probabilities that a process will join the first and second queues are 0.775 and 0.025, respectively. Suppose the first processor is a high-speed processor with high departure rate of 30, while the second processor is a low speed processor with a low departure rate of 10. Suppose the arrival rate in the system from the outside world is 4 and maximum number of processes to be allowed into first queue is 20, while maximum number of processes to be allowed into the second queue is 5. The experimental trials were carried out several times, in each trial, the arrival rate was changed, and the corresponding variation was obtained as the result of the simulation.

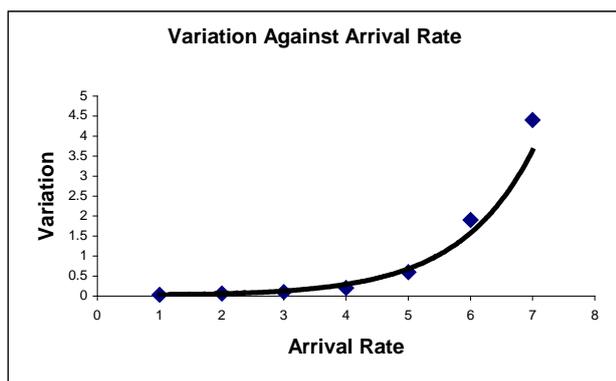


Figure 2: Variation Against Arrival Rate

The result shows that for non-compute intensive application, where the overall utilization factor of each of the processors is less than 1, increasing the departure rate will lead to a corresponding increase in the variation of the performance metric of the heterogeneous parallel computer.

Furthermore, figure 3 shows the simulation result when as we keep the following input parameters constant, the probability that a process will leave the network is 0.2, the probabilities that a process will join queue 1 and 2 are 0.775 and 0.025, respectively. The departure rates for processor 1 and 2 are 30 and 10, respectively, while the maximum number of processes in queue 1 and 2 (degree of multiprogramming for the two queues) are 20 and 5, respectively, and the arrival rate from the outside world is 4. By varying the degree of multiprogramming (maximum number of processes in the system) for the two queues of a two-processor parallel computer system, we obtain the corresponding variation shown in figure 3.

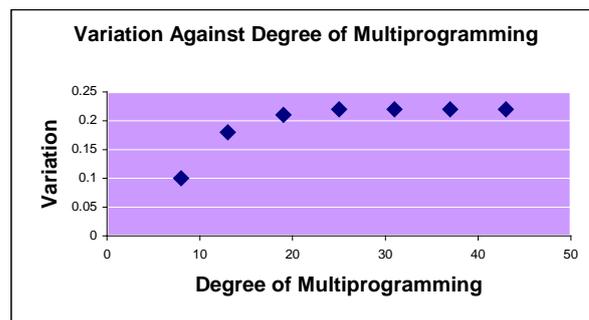


Figure 3: Variation Against the Degree of Multiprogramming

The interpretation of the result is that for non-compute intensive applications, where the overall utilization factor for each of the queues is less than 1, increasing the maximum number of processes that can be in each of the queues (the degree of multiprogramming) will lead to corresponding increase in the variation of the performance metric.

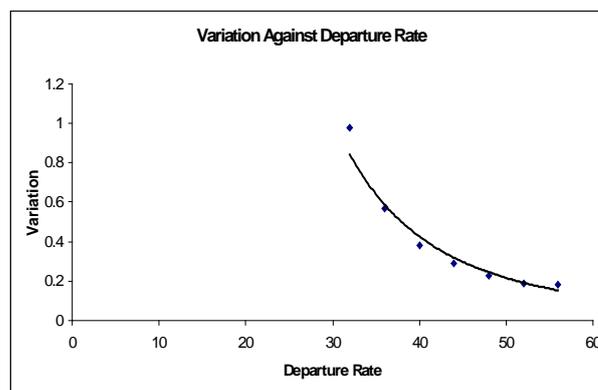


Figure 4 Variation Against Total Departure Rate

In a similar manner, as we keep the following input parameters constant, probability of a process leaving the network is 0.2, while the probability of a process going to queue 1 and 2 is 0.4. The maximum number of processes that can be in queue 1 and 2 are 15 and 14, respectively. By changing the departure rates of the two processors, we obtain the corresponding variation of the performance metric, as shown in figure 4. The result shows that for non-compute intensive application, increasing the speed of the various processors will lead to a corresponding decrease in the variation of the performance metric under consideration.

VII. SUMMARY AND CONCLUSION

This paper has been able to model the variation of a performance metric of heterogeneous parallel computer system with distributed memory by combining analytic and recursive models, using queuing approach. The models have been simulated on the computer and the results of the simulation have been analyzed in order to determine when to realize minimum variation.

REFERENCES

- [1] Henry H. Liu and Pat V. Crain, An Analytic Model for Predicting the Performance of SOA-Based Enterprise Software Applications, Proc. International Conference of Computer Measurement Group, (2004).
- [2] S. Balsamo et al, A Review of Queueing Network Models with Finite Capacity Queues for Software Architecture Performance Prediction, (2002).
- [3] Catalina M. Liado et al, A Performance Model Web Service, Proc. International Conference of Computer Measurement Group, (2005).
- [4] Rosselio, J et al, A Web Service for Solving Queueing Network Models Using PMIF. www.perfeng.com/papermdx.htm, (2005).
- [5] Cathy H. Xia, Zhen Liu., Queueing systems with long-range dependent input process and subexponential service time. Proc. ACM SIGMETRICS international conference on Measurement and modeling of computer systems,(2003).
- [6] Shanti Subramanyam, Performance Modelling of a J2EE Application to meet Service Level s, Agreement, Proc. International Conference of Computer Measurement Group, (2005)
- [7] Hamdy A. T.,. Operation Research: An Introduction, Prentice-Hall of India, (1999).
- [9] Ivan Stojmenovic; Recursive Algorithms in Computer Science Courses : Fibonacci Numbers and Binomial Coefficients; IEEE Transactions on Education; Vol. 48, No. 3
- [10] Arjan J.C. van Gemund; Performance Modelling of Parallel Systems: An Introduction.
- [11] Justyna Berlinska, The Statistical models of parallel applications, Annales UMCS Informatica, (2005).
- [12] Arranchenkov, K.E., Vilchersky, N.O., Shevlyakov, G.L Priority queueing with finite buffer size and randomized push-out; mechanism. Proc. of ACM SIGMETRICS international conference on measurement and modeling of computer systems.; (2003).
- [13] Abunday, B.D., and Khorram, E. The finite source queueing model for multiprogrammed computer systems with different CPU times and different I/O times. Acta Cybern. 8, 4 , (1998)
- [14] J. Sztrik; Finite-Source Queueing Systems and their Applications: A Bibliography;
- [15] Trivedi K. Shridharbhai, Probability and Statistics with Reliability, Queuing and Computer Science Applications, John Wiley & Sons Inc., (2002).
- [16] Per Brinch Hansen. Operating System Principles. Prentice-Hall of India Private Limited, (1990).
- [20] J. Sztrik^a and T. Gál A recursive solution of a queueing model for a multi-terminal system subject to breakdowns; Performance Evaluation Volume 11, Issue 1, Published by Elsevier, (1990).
- [23] Robert V. Hogg and Allen T. Craig; Introduction to Mathematical Statistics; Macmillan Publishing Co. Inc.; (1978).
- [24] Andrea Clemantis, Angelo Corana; Modelling Performance of Heterogeneous Parallel Computer System; Journal of Parallel Computing, Volume 12, Issue 9, Elsevier; pages 1131-1145; (1999).
- [25] E. Post, H.E. Goosen; Evaluating the Parallel Performance of a Heterogeneous System
- [26] Beutler, F; Mean sojourn times in markov queueing network: Little's formula revisited; IEEE Transaction on Information Theory; Volume 29, Issue 2, page 233-241; (2003).
- [27] Ken Vastola;
<http://networks.ecse.rpi.edu/~vastola/pslinks/perf/node46.html>
- [28] Xiaodong Zhang, Yong Yan; Modeling and Characterizing Parallel Computing Performance on Heterogeneous Network of workstations; Proceedings of the 7th IEEE Symposium on Parallel and Distributed Processing (SPDP '95) 1063-6374/95 \$10.00 © 1995 IEEE
- [29] O.E. Oguike et al; Modelling the Performance of Computer Intensive Applications of Parallel Computer System; Proc. Of IEEE 2nd International Conference on Computational Intelligence, Modeling and Simulation; (2010).
- [30] O.E. Oguike et al; Evaluating the Performance of Parallel Computer System Using Recursive Models; Proc. Of IEEE 4th UKSim European Modeling Symposium; (2010).
- [31] O.E. Oguike et al; Evaluating the Performance of Heterogeneous Distributed Memory Parallel Computer System Using Recursive Models; 2nd IEEE International Conference on Intelligent Systems, Modeling and Simulation; (2011).