# Dependability Prediction of High Availability OSCAR Cluster Server

Chokchai (Box) Leangsuksun, Lixin Shen,
Tong Liu, Hertong Song
*Department of Computer Science*
*Louisiana Tech University*
*Ruston, LA, U.S.A.*

Stephen L. Scott
*Computer Science and Mathematics Division*
*Oak Ridge National Laboratory*
*Oak Ridge, TN, U.S.A.*

## Abstract

*High availability (HA) computing has recently gained much attention, especially in enterprise and mission critical systems. The HA is now a necessity that is no longer regarded as a luxury feature. Thus, we, conjunctively with the open source community, are in process of enhancing the HA feature to Open Source Cluster Application Resources (OSCAR), a widely adopted Linux PC cluster system. Server redundancy will be our initial key aspect of the next generation HA OSCAR cluster system. In this paper, we introduce a HA server for OSCAR cluster system. Its architecture and mechanism is discussed, and then we model and predict the dependability of the system by a Petri net-based model, Stochastic Reword Net (SRN). The reliability and instantaneous availability of the system are presented as a result.*

*Keywords: cluster computing, Petri net modeling, Stochastic Reword Net, high availability, OSCAR, heartbeat*

## 1. Introduction

Cluster computing is becoming increasingly practical for high-performance computing (HPC) research and development. Because they are affordable and easily scaleable, PC clusters have taken the key role on HPC infrastructure. In the past few years, an increasing number of Beowulf clusters have been set up for the purpose of scientific research or simply for exploration of the frontier of supercomputer building [1]. As a fully integrated software package for easy to build Beowulf clusters, a significant number of HPC users, has adopted Open Source Cluster Application Resources (OSCAR) for their platform of choice [2]. High availability (HA) cluster computing has gained much attention, especially for mission critical and enterprise applications. As a widely used Linux cluster computing system, OSCAR urgently needs some HA features. The server is the most important part of the cluster system. However, it is subjected to the single point of failure. Once the failure of server occurs,

the whole system will be out of service. To eliminate the single point of failure, a failover server is introduced as a redundancy—a common technique to improve availability and reliability of a system [3].

Performance analysis has been relatively matured enough that one can estimate the total system performance (e.g. in gigaflops, tera-flops). However, system dependability analysis is the total opposite and is typically ignored during design phase. In this paper, we discuss an OSCAR cluster system with HA enhanced architecture. We then model and predict the dependability of the system by a Petri net-based model, Stochastic Reword Net (SRN). With the dependability analysis, we demonstrate the system reliability and availability modeling results and then decide whether the dependability requirement can be met in advance. SRN has been successfully used in the dependability evaluation and prediction for complicated system, especially when the time-dependent behavior is of great interest: instantaneous availability and reliability for fault-tolerant system [4]. The Stochastic Petri Net Package (SPNP) [5] allows the specification of SRN models, the computation of steady and transient state.

In the following sections, we discuss the architecture and mechanism of the dual server system in section 2, build the SRN model, and analyze dependability of the system in section 3. As a result, the reliability and instantaneous availability of the system is presented. The conclusion and future works are discussed in section 4.

## 2. System architecture and mechanism

In order to eliminate the single point of failure in an OSCAR cluster system, the dual server system will provide a failover server. Figure 1 shows architecture of a dual server system [6]. The reliable communication feature with a heartbeat mechanism is applied to the cluster. The periodical transmission of heartbeat messages traverses across a dedicated serial link, as well as an IP channel, and works as health detection of the working server. When a failure of the primary server occurs, the system can automatically transfer the control of services

to the failover server, allowing services to remain available with minimal interruption. The failed server can then be repaired while the application continues on the failover server. The failover server is idle when the primary server is in service. After the repair, the primary server takes over the control from the failover server, and the failover server becomes a backup server which continuously checks the health of the primary server by the heartbeat [7].
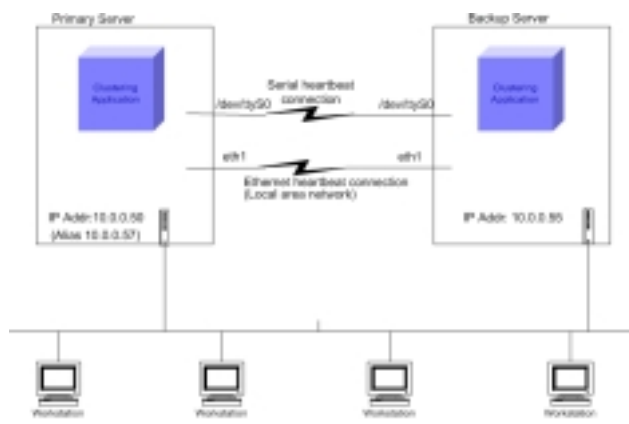


**Figure 1. The architecture of a dual server cluster system**

In our dual server system, both of the servers are running on Red Hat *7.3*. The workstations visit the primary server through alias IP address *10.0.0.57*, and the backup server is standby. Heartbeat is established between the two servers via channel *ttyS0* and *eth1*, respectively. Two channels are established for improving the accuracy of the heartbeat. If the backup server does not receive any heartbeat message packets from the primary server through neither of the two channels, it will assume that the primary server is down. In this case, the failover server takes over the control from the primary server by the alias IP address *10.0.0.57* so that the workstations are able to continue functioning seamlessly. When the primary server is ready back to work, the failover server releases the control of service back to the primary server, and works as an idle server, checking the heartbeat message from the primary server [8].

# 3. Stochastic reward nets model

We consider using SRN model to our cluster dependability analysis. The method is detailed in [4]. We assume that the two servers are dissimilar. The time to failures between two servers $S_1$ and $S_2$ is assumed to be a random variable with the corresponding distributions being exponential with rates $\lambda_1$ and $\lambda_2$, respectively. We consider the system to be functioning as long as one of the two servers is functioning [9]. We will specify our

model to SPNP and provide its solution in the next sections.

## 3.1 Reliability for dual server system

We consider the system reliability in two cases: (1) failures of the two servers are independent, and more accurately, when one server is suffering a failure, the other one can not fail. (2) there is also a failure mode where both servers can fail simultaneously, called common-mode failure (CMF). The SRN in Figure 2 models the failure behavior of this system [10].
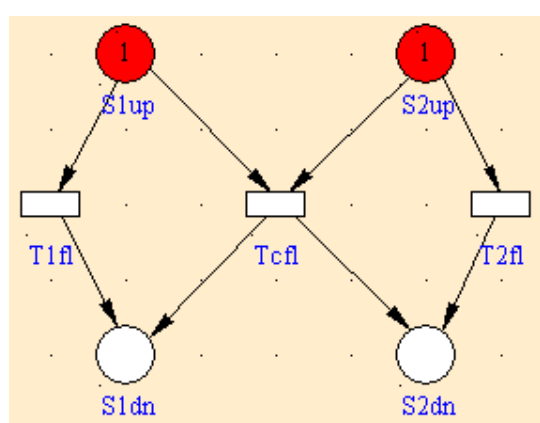


**Figure 2: The SRN model for system reliability**

We compute the reliability for this system assuming that $\lambda_1 = 1 \times 10^{-3}\ hr^{-1}$, $\lambda_2 = 2 \times 10^{-3}\ hr^{-1}$ and $\lambda_3 = 5 \times 10^{-4}\ hr^{-1}$. The reward rate assignment is specified as follows:

$$r_i = \begin{cases} 1 & if \quad (\#(S1up) = 1) \quad or \quad (\#(S2up) = 1) \quad in \quad marking \quad i \\ 0 & otherwise \end{cases}$$

where $r_i$ represents the reward rate assigned to state $i$ of the SRN, and $\#(s)$ represents the number of tokens in place $s$. The reliability of the system at time $t$ is computed as the expected instantaneous reward rate $E[X(t)]$ at time $t$ and its general expression is

$$E[X(t)] = \sum_{k \in \tau} r_k \pi_k(t)$$

where $\tau$ is the set of tangible marking, $\pi_k(t)$ is the probability of being in marking $k$ at time $t$ [11].

The system reliability as a function of time is plotted in Figure 3. The graph shows the system reliability with and without the CMF. We observe the contribution of CMF to the degradation in the system reliability. The mean time to failure (MTTF) of the system with CMF is *1000* hours, while the MTTF of the system without CMF is *1167* hours.

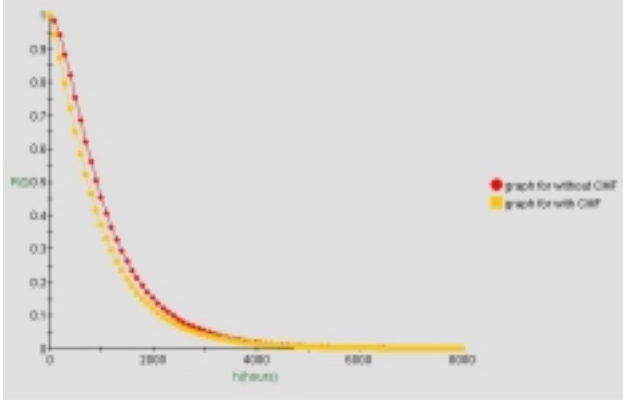## 3.2 Availability for the dual server system

**Figure 3: System reliability with and without the common-mode failure (CMF).**

We further consider a repair of the servers where the time to repair the servers is also exponentially distributed with rates $\mu_1$ and $\mu_2$, respectively. We also assume that when both the servers are waiting for repair, $S_1$ has priority for repair over $S_2$. The SRN in Figure 4 illustrates a model for the failure behavior of this system [10].
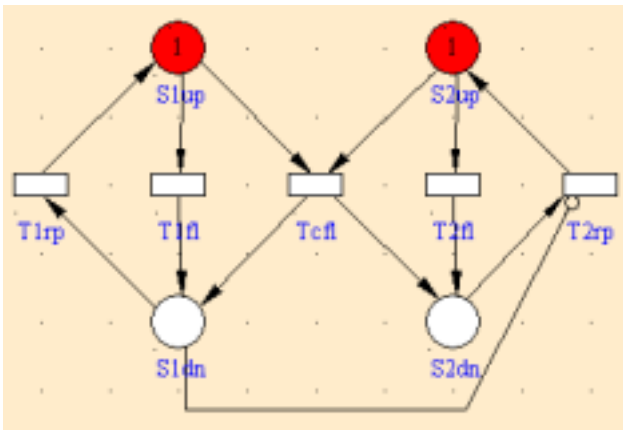


**Figure 4: The SRN model for system availability**

We compute the instantaneous availability of the system (probability that the system is availably at time $t$) assuming that $\lambda_1 = 1 \times 10^{-3}\ hr^{-1}$, $\lambda_2 = 2 \times 10^{-3}\ hr^{-1}$, $\lambda_3 = 5 \times 10^{-4}\ hr^{-1}$, and $\mu_1 = 0.5\ hr^{-1}$, $\mu_2 = 1\ hr^{-1}$. The same reward rate assignment as used for reliability is used for the expected instantaneous reward rate of the system.

The instantaneous availabilities for single server and dual server systems are plotted in Figure 5. We observed that an instantaneous availability decreases with time and reaches a steady-state value, which is equal to the steady-state availability of the system. From the figure, we can see that the system availability is greatly improved by the server redundancy.
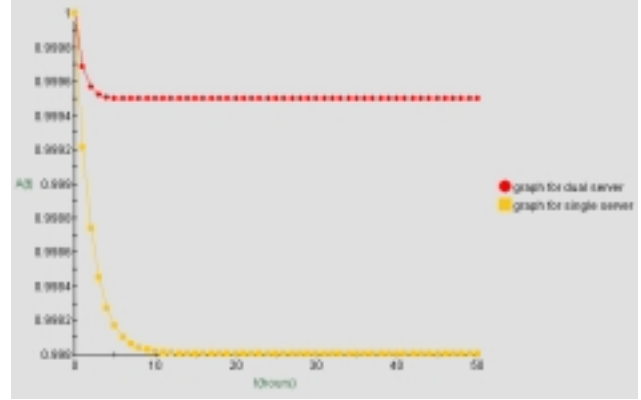


**Figure 5: System availability comparison between a single and dual server.**

## 4. Conclusions and future work

With dependability analysis, we can evaluate the HA features and compare them with dependability requirement of the system. The dual server system demonstrates the system reliability and availability improvement. We also show that SRN and SPNP methods are efficient for building and solving this model.

Besides the redundant server, we are considering adding more HA features to the next generation OSCAR cluster system, for example, the redundant disk storage, redundant network component, providing redundant network connection, etc [12]. For the dependability analysis, SRN and SPNP will be further studied in system dependability evaluation of the complex system.

## 5. Acknowledgements

## 6. References

[1] http://www.linuxjournal.com/article.php?sid=5862

[2] M.J. Brim, T.G. Mattson, and S.L. Scott, "OSCAR: Open Source Cluster Application Resources", *Ottawa Linux Symposium 2001*, Ottawa, Canada, 2001

[3] Hewlett Packard, *Managing MC/ServiceGuard,* Hewlett-Packard Company, October, 1998

[4] J. Muppala, G. Ciardo, and K. Trivedi, "Stochastic Reward Nets for Reliability Prediction", *Communications*

*in Reliability, Maintainability and Serviceability,* SAE, July 1994, Vol. 1, No. 2, pp. 9-20.

[5] G. Ciardo, J. Muppala, and K. Trivedi, "SPNP: Stochastic Petri net package", *Proc. Int. Workshop on Petri Nets and Performance Models,* IEEE Computer Society Press, Los Alamitos, CA, Dec. 1989. pp 142-150.

[6] P.S. Weygant, X. Bui, and W. Sawyer, *Clusters for High Availability: A Primer of HP Solutions (2nd Edition),* Prentice Hall PTR, May, 2001

[7] F. Piedad, and M. Hawkins, *High Availability: Design, Techniques and Processes,* Prentice Hall PTR, December, 2000

[8] O. Kolesnikov and B. Hatch, *Building Linux Virtual Private Networks.* New Riders Publishing, February, 2002

[9] R.A. Sahner, K.S. Trivedi, and A. Puliafito, *Performance and Reliability Analysis of Computer Systems: An Example-Based Approach Using the SHARPE Software Package,* Kluwer Academic Publishers, 1996.

[10] K.S. Trivedi, *SPNP User's Manual Version 6.0*, Duke University, September, 1999.

[11] M. Malhotra, K.S. Trivedi, *Dependability Modeling Using Petri-Nets*, IEEE Transactions on Reliability, Sept., 1995, Vol. 44, No. 3, pp. 428-440.

[12] http://linux-ha.org/