

Availability Modeling and Analysis of Cloud Systems Involving HA Techniques via Stochastic Reward Nets

Luyao Zou, and Choong Seon Hong

Department of Computer Science and Engineering, Kyung Hee University, South Korea.

Email: {zouluyao, cshong}@khu.ac.kr

Abstract

Nowadays, cloud computing is giving significant advantages to both enterprises and users. However, serious server downtime is possible to occur, which can make substantial economic losses to cloud providers and reduce user experience. Therefore, it is essential to create a system with high availability. Accordingly, this paper proposed stochastic models, which contain specific interactions and system behaviors for availability analysis of a cloud system. Besides, physical server/VM recovery, redundancy, and VM live migration techniques are incorporated to enhance availability. Furthermore, various metrics of interest: SLA-based steady state availability (SSA), SLA-based SSA sensitivity analysis, SLA violation downtime, and cost are analyzed in this paper.

1. Introduction

Cloud services, run on the centralized physical systems, require huge demands of high availability and less downtime costs [1]. Therefore, enterprises are seeking researches on availability evaluation and analysis with small-scale system before deploying and implementing real scale system. However, analyzing availability in real system not realistic due to it may cause problems like security. Consequently, it is important to make availability model before deploying a real system. Stochastic reward nets (SRN) which can capture system dynamics and system dependable behaviors more deeply, compared with Fault Tree, Reliability Block Diagram and Markov chains [2], is used in this study. Many researches such as [3–5] used SRN to make system model. However, they neglect VM level which has significant impact for system availability. Hence, not only physical level, but VM level are considered in this paper.

The main contributions of this paper are stated below. (1) The proposed high availability model with redundant (hot/cold pool) cloud system by SRN, in which all transitions follow exponential distributions. (2) The model includes failure/repair behavior, and interactions between hot/cold host and VM live migration to enhance the system availability. (3) We have compared various techniques (e.g., redundancy) involved in different SRN models. (4) We have analyzed SLA-based SSA, SLA violation downtime and cost on the default parameters, and sensitivity analysis regarding main impacting parameters: MTTF (mean time to failure) and MTTR (mean time to repair) of the

hot host, VM, respectively. The following of this paper is organized as follows. Section 2 illustrates the proposed SRN models. Section 3 shows the comparison between proposed SRN models and others. Sensitivity analysis is also showed in this section. Lastly, section 4 makes the conclusion of this paper.

2. System Description and SRN model

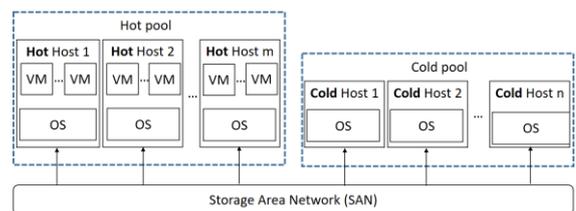


Figure 1 Redundancy Cloud Architecture

Figure 1 is the architecture of cloud system which contains two pools: hot pool for VM placement and cold pool as a redundancy of hot host. SAN (Storage area network) is equipped in this system as intermediate shared storage to save the VM data and the image files in case VM failed [1]. A case of SRN model of proposed system with two state (UP and DOWN), which contains two hot host, a cold host and two VMs on each hot host is studied in this section, showed as Fig.2 and Fig.3.

Figure 2 shows SRN model of physical server (hot host H1, H2 and cold host C). When H1 goes down, the transition T_{H1f} is fired and token in the place P_{H1up} is removed and be deposited into P_{H1f} . If P_{Cup} is not empty, token of P_{Cup} moves to P_{H1up} through firing transition $T_{C1start}$. Likewise, if H2 fails,

T_H2f will be fired, and the token moves from P_H2up to P_H2f. At this time, T_C2start is fired to transfer token from the place P_Cup to P_H2up. The failure event of C is the same as host but with different transition time. In regarding to repair behavior of each server, the tokens in the DOWN state are removed and be deposited back to its UP state by firing each corresponding transition: T_H1r, T_H2r or T_Cr.

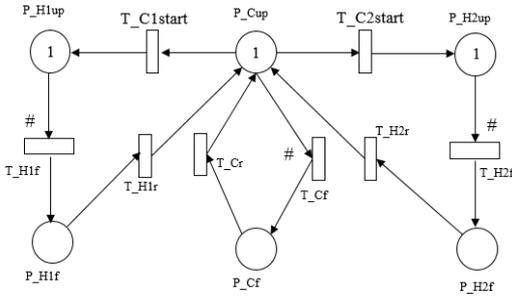


Figure 2. Host model with redundancy

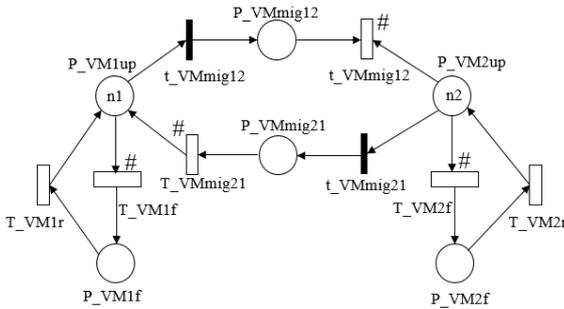


Figure 3. VM model with LM

Figure 3 depicts the Virtual Machine subsystem incorporating with VM live migration (LM) technique. VM failure and recovery process are the same like host. In regarding to LM, SAN is modeled as P_VMmig12 and P_VMmig21. When hot host goes down, tokens of VM will transfer to P_VMmig12 or P_VMmig21 immediately. Then, t_VMmig12 or t_VMmig21 is fired to transfer tokens to another hot host.

The failure and LM event are dependent on current number of tokens, which belong to P_H1up, P_H2up, P_Cup, P_VM1up, P_VM2up, P_VMmig12 and P_VMmig21, respectively. This is also named marking dependence and is marked with # notation near relative transitions. Moreover, guard function of SRN models is used to help the specify relationship between VMs and hosts and reward function is used to help analyze SLA-based availability, showed as table 2 and 3, respectively.

3. Numerical analysis and results

In this section, we use SPNP [6], a powerful and popular package for modeling complex system

behaviors of SRN models, to develop all the aforementioned SRN models. Physical system without incorporate the redundancy technique is also modeled to compare with above models. Table 1 shows the default input parameters of proposed SRN models which are referenced from different previous studies. As for the downtime (mins/(operational period hours) and its cost, the calculations are defined as follows:

$$DT = (1 - ssa)L * 60 \quad [7]$$

$$DTC = DT * 16,000 \text{ (USD)} \quad [1]$$

Table 1. Default parameters of the SRN models

Name	Transition	Description	Value
1/λ _h	T_H1f, T_H2f	Mean time to failure of a hot host	30 days
1/μ _h	T_H1r, T_H2r	Mean time to repair of a hot host	5 hours
1/λ _c	T_Cf	Mean time to failure of a cold host	60 days
1/μ _c	T_Cr	Mean time to repair of a cold host	4 hours
1/λ _{VM}	T_VM1f, T_VM2f	Mean time to failure of a VM	7 days
1/μ _{VM}	T_VM1r, T_VM2r	Mean time to repair of a VM	20 mins
1/ω _{VMmig}	T_VMmig12, T_VMmig21	Mean time to LM of a VM to a hot host	10 sec
1/ω _{Cstart}	T_C1start, T_C2start	Mean time to starting a cold host to take place a hot host	5 mins

Table 2. Reward function

Reward Function for SLA-based SSA analysis	
H2 and H2C1VM0:	
1	#P_H1up+#P_H2up>=SLA*nH
0	otherwise
H2C0VM4LM and H2C1VM4LM:	
1	#P_VM1up+#P_VM2up>=SLA*(n1+n2)
0	otherwise

Table 3. Guard function

Function	Transition	Definition
H2C0VM4LM		
gt_VMmig12	t_VMmig12	1 #P_H1up==0 0 otherwise
gt_VMmig21	t_VMmig21	1 #P_H2up==0 0 otherwise
H2C1VM0		
gT_C1start	T_C1start	1 #P_H1up==0 0 otherwise
gT_C2start	T_C2start	1 #P_H2up==0 0 otherwise
H2C1VM4LM		
gT_C1start	T_C1start	1 #P_H1up==0 0 otherwise
gT_C2start	T_C2start	1 #P_H2up==0 0 otherwise
gt_VMmig12	t_VMmig12	1 #P_H1up==0 0 otherwise
gt_VMmig21	t_VMmig21	1 #P_H2up==0 0 otherwise

Table 4. SSA and Downtime cost

Policy	Quantity of Host		VM	SLA-based Availability	Downtime (mins/year)	Downtime Cost(USD/year)
	Hot	Cold				
Non-Redundancy&VM LM (H2)	2	0	0	0.986207552554	7249.3103776176	11598896.60418816
Redundancy+Non-VM LM (H2C1VM0)	2	1	0	0.999610126508	204.9175073952	327868.01183232
Non-redundancy+VM LM (H2C0VM4LM)	2	0	2 per hot host	0.999901188778	51.9351782832	83096.28525312
Redundancy+VM LM (H2C1VM4LM)	2	1	2 per hot host	0.999947643201	27.5187335544	44029.97368704

Table 4 shows the SSA, downtime and relative cost per year of four different systems with various techniques. This table depicts that proposed system model (H2C1VM4LM) of this paper has the highest SLA-based SSA and lowest SLA violation downtime and cost.

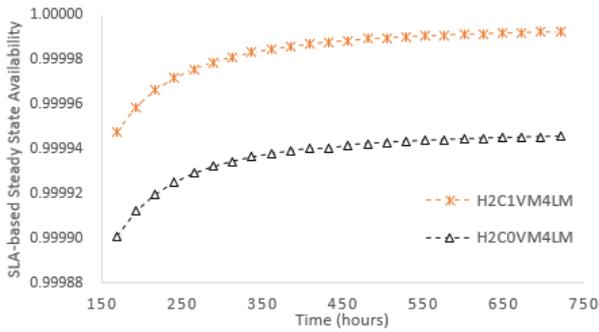


Figure 4 MTTF of VM

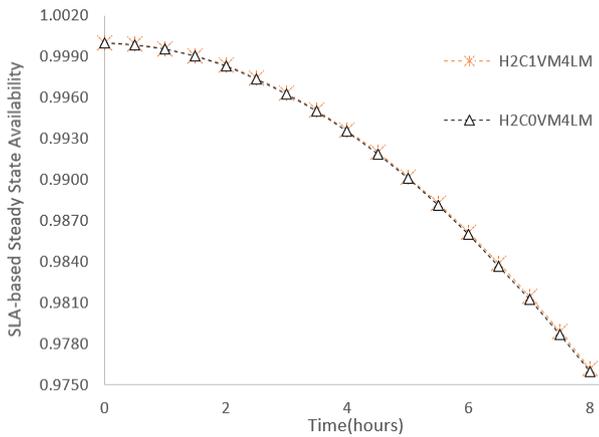


Figure 5 MTTR of VM

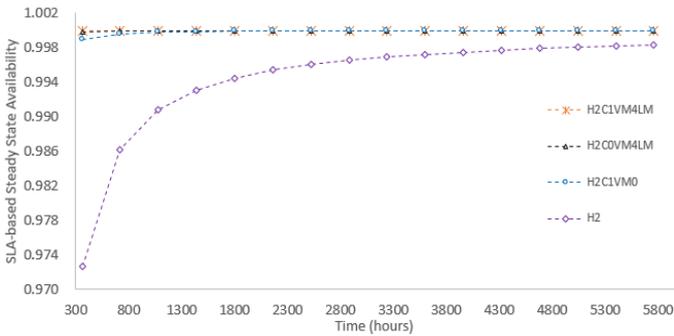


Figure 6. MTTF of Hot Host

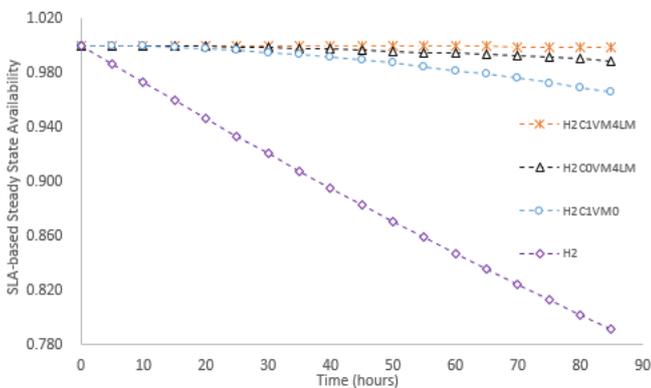


Figure 7. MTTR of Hot Host

Figure 4–7 illustrate the dependence of SSA with MTTF/MTTR of hot host and VM. From these figures, our proposed system always has the highest availability regardless of increasing MTTF or MTTR. In

addition, when increasing MTTF, the SSA gradually increased, and SSA experienced a downward trend when MTTR decreased.

4. Conclusion

In this paper, through modeling with SRN, it can be seen that the proposed cloud systems have high availability and low downtime cost. Furthermore, availability analysis with SRN models can solve the problems such as big financial losses, security, etc., which are caused by analyzing availability in real machines. Therefore, this method plays an important role for cloud providers before deploy the real systems.

Acknowledgement

This research was partially supported by the MSIT (Ministry of Science and ICT), Korea, under the Grand Information Technology Research Center support program (IITP–2018–2015–0–00742) supervised by the IITP(Institute for Information & communications Technology Promotion) and by Basic Science Research Program through the National Research Foundation of Korea (NRF) funded by the Ministry of Education (NRF–2016R1D1A1B01015320). *Dr. CS Hong is the corresponding author

References

- [1] Tuan Anh Nguyen, Dugki Min & Eunmi Choi. A comprehensive evaluation of availability and operational cost for a virtualized server system using stochastic reward nets. The Journal of Supercomputing. 2017.
- [2] Jogesh K. Muppala, Gianfranco Ciardo, Kishor S. Trivedi*. Stochastic Reward Nets for Reliability Prediction. CiteSeerX. 1994.
- [3] Pagadala Nagendra Babu, Seelam Naresh, Metta Lakshmi Prasanna. IMPLEMENTATION OF IAAS IN CLOUD USING MONOLITHIC MODEL. International Journal of Technical Research and Applications. 2016
- [4] Rahul Ghosh, Francesco Longo, Flavio Frattini, Stefano Russo, and Kishor S. Trivedi. Scalable Analytics for IaaS Cloud Availability. IEEE. 2014.
- [5] Dario Bruneo. A Stochastic Model to Investigate Data Center Performance and QoS in IaaS Cloud Computing Systems. IEEE. 2014.
- [6] Gianfranco Ciardo, Jogesh Muppala, Kishor Trivedi. SPNP: Stochastic Petri Net Package. IEEE. 1989.
- [7] Ghosh, Rahul (2012). Scalable Stochastic Models for Cloud Services. Dissertation, Duke University. Retrieved from <http://hdl.handle.net/10161/6110>.