# Availability modeling and analysis of a data center for disaster tolerance

Tuan Anh Nguyen [a,*], Dong Seong Kim [b], Jong Sou Park [a]

[a] *Department of Computer Engineering, Korea Aerospace University, South Korea*
[b] *Department of Computer Science and Software Engineering, University of Canterbury, New Zealand*

## HIGHLIGHTS

- We present a comprehensive availability modeling and analysis of a data center system for disaster tolerance.
- We assess availability characteristics of a data center regarding disaster occurrence, unexpected failure of network connection and complicated dependencies.
- The study reflects significance of the incorporation of disaster and fault tolerant techniques into geographically distributed data centers for high availability of cloud based businesses.
- The study provides a selection basis of designing for disasters considering the trade-off between system availability and downtime cost with infrastructure construction cost.

## ARTICLE INFO

## ABSTRACT

Availability assessment of a data center with disaster tolerance (DT) is demanding for cloud computing based businesses. Previous work attempted to model and analyze the computing systems without a good consideration on disaster occurrence, unexpected failure of network connection and proper dependencies between subsystems in a data center. This paper presents a comprehensive availability model of a data center for DT using stochastic reward nets (SRN). The model incorporates (i) a typical two-level high availability (HA) configuration (i.e., active/active between sites and active/passive within a site), (ii) various fault and disaster tolerant techniques; (iii) dependencies between subsystems (e.g. between a host and virtual machines (VMs), between a network area storage (NAS) and VMs) and dependency between disastrous events and physical subsystems; and (iv) unexpected failures during data transmission between data centers. The constructed SRN model is analyzed on the basis of steady state analysis, downtime cost analysis, and sensitivity analysis with regard to major impacting parameters. The analysis results show the availability improvement of the disaster tolerant data center (DTDC) and featured system responses corresponding to the selected variables. The modeling and analysis of the DTDC in this paper provide a selection basis of designing for disasters in consideration of the trade-off between system availability and downtime cost with infrastructure construction cost.

© 2015 Elsevier B.V. All rights reserved.

## 1. Introduction

Businesses of all sizes have been more and more dependent on their escalating IT infrastructures to achieve automation in management and analyses of business strategies and operations. Thus, business continuity [1] becomes a critical concern of business operations as well as physical computing systems in enterprises. To assure business continuity, computing systems must tolerate different single points of failure to mitigate downtime and improve system availability during long-term operations. The adverse conditions could be a failure with technology such as hardware failures, software defects; and/or a natural/man-made disaster like a storm, a fire, power outages or terrorist acts for instances. The increasing of catastrophic threats and disasters in the past all over the world emphasizes the demand of development in disaster tolerant computing systems. The impacts of a disaster or a severe system fault always highlight the interdependencies within a computing system in the way that a failure of a component might be propagated to other dependent components and escalate into a severe system failure, thus damaging business transactions. In practice, enterprises commonly employ Disaster Recovery (DR) [1–3] as a subset of Business Continuity Planning (BCP) to repair, recover and restore the business operations and resources. How-

---

\* Corresponding author.
*E-mail addresses:* anhnt@kau.ac.kr, anhnt2407@gmail.com (T.A. Nguyen).

ever, the continuance of business transactions in an uninterrupted manner is so critical and demanding nowadays that computing systems must be able to tolerate any type of unexpected disaster at any system level in order to provide continuous operations and connections even around the time of disaster occurrence. The ICT computing systems and resources, therefore, must handle this demand by adopting a variety of fault tolerant and disaster tolerant techniques under some amount of autonomy. Among modern IT computing infrastructures, a data center is the squeezed and consolidated information facility which is placed at the core of information business operations of most enterprises [4]. The quick paces in development of social networking services, e-commerce and cloud computing based businesses have pushed a vast evolution on data center systems in recent years. Additionally, the essential of maintaining users confidence and avoiding revenue losses requires High Availability (HA) of cloud services hosted on data centers. The difference in minutes of unavailability, while seemingly nominal, may cost the owner a massive amount of money due to interruption of business transactions [5]. However, the higher availability a data center is designed to gain, the more capital and operational costs the enterprises need to spend for building and managing it [6]. Thus, assuring HA and assessing the availability characteristics of such disaster tolerant computing systems are of paramount importance. The ICT planners and managers of enterprises critically demand to conduct quantitative evaluation of the availability of a pre-designed data center with Fault Tolerance (FT) and Disaster Tolerance (DT) in advance of system construction due to the conflicting requirements of the HA levels demanded against cost values. Consequently, the approaches on the evaluation and optimization of these requirements are in sharp focus of data center designers. We will be focusing on modeling and analysis of a typical Disaster Tolerant Data Center (DTDC) in this paper.

The fundamental definitions and understanding of disaster tolerant computing and communication systems was first introduced in the works [7,8]. There have been also several papers on the perception of high available and disaster tolerant ICT infrastructure [9–12]. Some others presented different methodologies for implementation of DT in large-scale computer and communication systems [13,14]. Nevertheless, none of these works provided a proper quantitative evaluation of a particular system with DT. Only a few recent work [15,16] tried to model and analyze simplified cluster and cloud computing systems with DT. This motivates us to model and analyze a typical disaster tolerant computing system using a stochastic model. The approach in the previous work [16] showed: (i) a system configuration with one operating data center and the other in cold-standby mode, which is active/cold-standby HA configuration at data center scale according to the work [17]; the configuration is further extensible to other types of HA configurations in the work [17] defined by the initial number of running VMs in each data center; (ii) a fixed number of VMs in a certain site are in running state at most of the time but all of the remaining VMs in both sites stay in standby state for the sake of operational failover within/between data centers without significant contribution to system throughput; (iii) VMs (depicted by tokens in system model) still retain in physical servers (represented by tokens in ready state) in the cases of disaster occurrence or multiple hosts simultaneous failure; however in their extended work [18], this drawback was resolved by transferring all VMs in operational states to the remaining operational data center; (iv) typical VM migration mechanisms are implicitly incorporated both in the works [16,18]; (v) dependencies between different parts of the system architecture were captured incompletely in the work [16] but were taken into account in their extended work [18]; (vi) at last, the transmission of large-sized VM image file (at gigabytes) throughout such a long distance (at hundreds or thousands kilometers) between data centers or between backup server and the data centers were carried out smoothly without proper consideration of the probability of failure occurrence, however, this is only an ideal case because it is not practical to be free from any failure during data transmission over such a long distance in real ICT systems. Enumerating the key features of the previous work, we find that modeling and analyzing such a sophisticated DTDC using stochastic models [19] is still a preliminary endeavor. We attempt to incorporate and capture more detailed system behaviors in a complete manner. We summarize the main contributions of our work as follows:

- modeled two configurations (active/active configuration between sites and active/standby one within a site) for a DTDC system with two geographically distributed data centers and two servers in a data center under a two-level HA configuration [17],
- incorporated different fault tolerant and disaster tolerant techniques as well as systems detailed behaviors such as VM failover and VM live-migration between hosts, geographically distributed site redundancy, and geographical VM migration and failover,
- captured the dependencies between components in the system architecture in detail: (a) between hosts and VMs, (b) between network area storage (NAS) and VMs, and (c) between the events of disaster occurrence and the above mentioned physical components (hosts and NAS),
- incorporated unexpected failures during VM transmission (VM failover and migration) between geographically distant systems using an imperfect coverage factor,
- constructed a DTDC system model using SRN which enables one to model a detailed system behaviors and dependencies and to improve the clarification and intuition of system model presentation,
- analyzed the constructed SRN model in terms of steady-state availability, downtime cost, and sensitivity with respect to several major parameters.

*Through the modeling and analysis, we have found the followings:*

- A data center with FT and DT located in a place with a few disastrous events (safer area) has higher availability even with less quality network connectivity; the steady state availability of such data center can achieve tier-4 (four nines) of the renowned uptime standards for industrial data centers [5,20,21].
- Failures in long-haul VM transmission between backup server and data centers do affect significantly while the one between data centers has less effect on the system steady state availability.
- A correlation between network speed and VM image size in consideration of their impact on the system availability:
  – In a system with high-speed network connection between geographically distributed components, the image size slightly affects the system availability
  – In the case of slow-speed of network connection, the image size with bigger value drastically pull down the system availability.
- This work in relation to the previous work [16] and its extension [18] brings about a proper consideration of failures in data transmission process and provides an evaluation and trade-offs basis of infrastructure costs in designing for disaster tolerance of data centers.

The rest of this paper is organized as follows. Related work is presented in Section 2. Section 3 introduces a DTDC. Section 4 presents SRN models for the DTDC. The numerical analysis and discussion are presented in Section 5. Finally, Section 6 concludes the paper.