

Chapter 3.8

Model-Driven Automated Error Recovery in Cloud Computing

Yu Sun

University of Alabama at Birmingham, USA

Jules White

Virginia Tech, USA

Jeff Gray

University of Alabama, USA

Aniruddha Gokhale

Vanderbilt University, USA

ABSTRACT

Cloud computing provides a platform that enables users to utilize computation, storage, and other computing resources on-demand. As the number of running nodes in the cloud increases, the potential points of failure and the complexity of recovering from error states grows correspondingly. Using the traditional cloud administrative interface to manually detect and recover from errors is tedious, time-consuming, and error prone. This chapter presents an innovative approach to automate cloud error detection and recovery based on a run-time model that monitors and manages the running nodes in a cloud. When administrators identify and correct errors in the model, an inference engine is used to identify the specific state pattern in the model to which they were reacting, and to record their recovery actions. An error detection and recovery pattern can be generated from the inference and applied automatically whenever the same error occurs again.

INTRODUCTION

With the increasing complexity of software and systems, domain analysis and modeling are becoming more important for software development and

system applications. Applying domain-specific modeling languages and transformation engines is an effective approach to address platform complexity and the inability of third-generation languages to express domain concepts clearly (Schmidt, 2006). Building correct models for a specific domain can often simplify many complex

DOI: 10.4018/978-1-4666-0879-5.ch3.8

tasks, particularly for distributed applications based on cloud computing (Hayes, 2008) that offer several opportunities for customization and variability.

Cloud computing shifts the computation from local, individual devices to distributed, virtual, and scalable resources, thereby enabling end-users to utilize the computation, storage, and other application resources (which forms the “cloud”) on-demand (Hayes, 2008). Amazon EC2 (Elastic Compute Cloud) (<http://aws.amazon.com/ec2/>, 2009) is an example cloud computing platform that allows users to deploy different customized applications in the cloud. A user can create, execute, and terminate the application instances as needed, and pay for the cost of time and storage that the active instances use based on a utility cost model (Rappa, 2004).

In the cloud computing paradigm, the large number of running nodes increases the number of potential points of failure and the complexity of recovering from error states. For instance, if an application terminates unexpectedly, it is necessary to search quickly through the large number of running nodes to locate the problematic nodes and states. Moreover, to avoid costly downtime, administrators must rapidly remedy the problematic node states to avoid further spread of errors.

Just like standard enterprise applications, cloud computing applications can suffer from a wide range of problems stemming from hard-

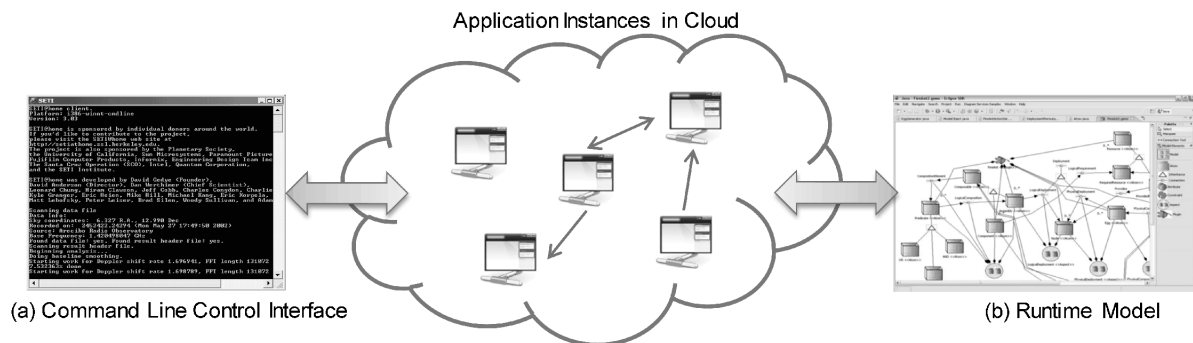
ware failure to operator error (Oppenheimer et al., 2003). For example, Amazon EC2 provides limited guarantees about availability or reliability of hardware or virtual machine (VM) instances. Operators must be prepared to re-launch VM instances when failures occur, transfer critical data to newly provisioned VM images, start critical services on new VM instances, join new nodes to virtual LANs or security contexts, or update load balancers and elastic IP addresses to reference newly provisioned infrastructure.

Although Amazon EC2 provides a user-friendly and simple interface to manage and control the application instances (Figure 1a), administrators must still be experienced with the administrative commands, the configuration of each application, as well as some domain knowledge about each running instance. Administrators must therefore be highly trained to effectively and efficiently handle error detection and error recovery. The complexity of managing a large cloud of nodes can increase maintenance costs, especially when personnel are replaced due to turnover or downsizing.

Even with experienced administrators, the process of error recovery involves the following challenges:

- It is hard to locate errors accurately with a large number of running application instances.

Figure 1. Two options to control application instances



19 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage:
www.igi-global.com/chapter/model-driven-automated-error-recovery/64509

Related Content

A Collaborative Model of Information Technology Strategic Plan for the Government Sector

Wagner N. Silva, Marco A. Vazand Jano M. Souza (2013). *International Journal of Distributed Systems and Technologies* (pp. 29-38).

www.irma-international.org/article/a-collaborative-model-of-information-technology-strategic-plan-for-the-government-sector/104716

Security Aspects in Utility Computing

Mayank Swarnkarand Robin Singh Bhadoria (2016). *Emerging Research Surrounding Power Consumption and Performance Issues in Utility Computing* (pp. 262-275).

www.irma-international.org/chapter/security-aspects-in-utility-computing/139847

Power Optimization Using Clock Gating and Power Gating: A Review

Arsalan Shahid, Saad Arif, Muhammad Yasir Qadriand Saba Munawar (2016). *Innovative Research and Applications in Next-Generation High Performance Computing* (pp. 1-20).

www.irma-international.org/chapter/power-optimization-using-clock-gating-and-power-gating/159037

A High Performance Model for Task Allocation in Distributed Computing System Using K-Means Clustering Technique

Harendra Kumar, Nutan Kumari Chauhanand Pradeep Kumar Yadav (2018). *International Journal of Distributed Systems and Technologies* (pp. 1-23).

www.irma-international.org/article/a-high-performance-model-for-task-allocation-in-distributed-computing-system-using-k-means-clustering-technique/207689

A Scalable Approach to Real-Time System Timing Analysis

Alan Griggand Lin Guan (2012). *Grid and Cloud Computing: Concepts, Methodologies, Tools and Applications* (pp. 637-668).

www.irma-international.org/chapter/scalable-approach-real-time-system/64507