

# How to Tame Your Online Services

Qingwei Lin, Jian-Guang Lou, Hongyu Zhang, Dongmei Zhang

Microsoft Research, Beijing, China

Email: {qlin, jlou, honzhang, dongmeiz}@microsoft.com

## Background

Online service systems, such as online banking systems and e-commerce systems, have been increasingly popular and important in our society. During operation of an online service, there can be a live-site service incident: an unplanned interruption, outage, or degradation in the quality of the service. Such incidents can lead to huge economic loss or other serious consequences. For example, the estimated average cost of one hour's service downtime for Amazon.com is \$180,000 [1].

Once a service incident occurs, the service provider should take actions immediately to diagnose the incident and restore the service as soon as possible. A typical procedure of incident management in practice (e.g., at Microsoft and other service-provider companies) goes as follows. When the service monitoring system detects a service violation, the system automatically sends out an alert and makes a phone call to a group of On-Call engineers to trigger an incident investigation. Given an incident, engineers need to understand what the problem is and how to resolve it. In ideal cases, engineers can identify the root cause of the incident and fix it quickly. However, in many cases, engineers are unable to identify or fix root causes within a short time, as it usually takes time to identify and fix the root causes, conduct regression testing, and re-deploy the new version to data centers. Thus, in order to recover the service as soon as possible, a common practice is to restore the service by identifying a temporary workaround solution (such as restarting a server component) to restore the service. Then after service restoration, identifying and fixing the underlying root cause for the incident can be conducted via offline postmortem analysis. Incident management has become a critical task for online services. The goal is to minimize the service downtime and to ensure high quality of the provided services. In practice, incident management of an online service heavily depends on data collected at runtime of the service, such as service-level logs, performance counters, and machine/process/service-level events. Such monitoring data typically contains information that reflects the runtime state and

behavior of the service. Based on the data collected, service incidents can be detected and mitigated in a timely way.

## Service Analysis Studio

We formulated the problem of incident management for online services as a software analytics problem [2], which can be tackled with phases of task definition, data preparation, analytic-technology development, and deployment and feedback gathering. We carried out a two-year research project, where we designed a set of incident management techniques based on the analysis of a huge amount of data collected at service runtime [3]. As a result of this project, we developed a tool called Service Analysis Studio (SAS), which targets real incident management scenarios of large-scale online services provided by Microsoft.

SAS includes a set of data-driven techniques for diagnosing service incidents. Each of these techniques targets at a specific scenario and a certain type of data. Here we just briefly introduce some of the major techniques SAS offers:

- *Identification of Incident Beacons from System Metrics:* When engineers diagnose incidents of online services, they usually start by hunting for a small subset of system metrics that are symptoms of the incidents. We call such kinds of metrics “service-incident beacons”. A service-incident beacon could provide useful information helping engineers locate the cause of an incident. For example, when a resource intensive SQL query blocks the execution of other queries accessing the same table, symptoms can be observed on monitoring data: the waiting time on the SQL-inducing lock becomes longer, and the event “SQL query time out failure” is triggered. Such metrics can be considered service-incident beacons. We developed data mining based techniques that helped engineers effectively and efficiently identify service-incident beacons from such huge number of system metrics. The technical details can be found in [4, 5].
- *Leveraging Previous Effort for Recurrent Incidents:* Engineers of an online service system may receive many similar incident reports. Therefore, leveraging the knowledge from past incidents can help improve the effectiveness and efficiency of incident management. The key here is to design a technique that automatically retrieves the past incidents similar to the new one, and then proposes a potential restoration action based on the past solutions. More details can be found in [6].

- *Mining Suspicious Execution Patterns*: Transactional logs provide rich information for diagnosing service incidents. When scanning through the logs, engineers usually look for a set of log events that appear in the log sequences of failed requests but not in the ones of the succeeded requests. Such a set of log events are named as suspicious execution patterns. A suspicious execution pattern could be an error message indicating a specific fault, or a combination of log events of several operations. For example, many normal execution paths look like {task start, user login, cookie validation success, access resource R, do the job, logout}. In contrast, a failed execution path may look like {task start, user login, cookie not found, security token rebuild, access resource R error}. The code branch reflected by {cookie not found, security token rebuild, access resource R error} indicates a suspicious execution pattern. We proposed a mining-based technique to automatically identify suspicious execution patterns. The details of our technique can be found in [6].

## Success Story

We have successfully applied SAS to one of Microsoft large-scale services (a geographically distributed, web-based service serving hundreds of millions of users). Similar to other online services, the Microsoft service is expected to provide high-quality service at all times. In the past, the service team was facing great challenges in improving the effectiveness and efficiency of their incident management in order to provide high-quality service. SAS was first deployed to the datacenters of the service in June 2011. The engineers of the service team have been using SAS for incident management since then. The actual usage experience shows that SAS helps the engineers improve the effectiveness and efficiency of incident management. According to the usage data from a 6-month empirical study, about 91% of engineers used SAS to accomplish their incident management tasks and SAS was used to diagnose about 86% of service incidents. Now SAS has been successfully deployed to many Microsoft product datacenters and widely used by on-call engineers for incident management.

Our experience shows that incident management has become a critical task for a large-scale online service. Software analytics techniques can be successfully applied to ensure high quality and reliability of the service, utilizing the data collected at service runtime.

## References

[1] D. A. Patterson. "A simple way to estimate the cost of downtime". In Proc. of LISA' 02, pp. 185-188, 2002

[2] D. Zhang, S. Han, Y. Dang, J. Lou, H. Zhang, T. Xie. "Software Analytics in Practice". IEEE Software, Special Issue on the Many Faces of Software Analytics, vol. 30 no. 5, pp. 30-37, September/October 2013.

[3] Jian-Guang Lou, Qingwei Lin, Rui Ding, Qiang Fu, Dongmei Zhang, and Tao Xie. Software Analytics for Incident Management of Online Services: An Experience Report, in Proceedings of the 28th IEEE/ACM International Conference on Automated Software Engineering (ASE 2013), Palo Alto, California, November 2013.

[4] Meng-Hui Lim, Jian-Guang Lou, Hongyu Zhang, Qiang Fu, Andrew Beng Jin Teoh, Qingwei Lin, Rui Ding and Dongmei Zhang. Identifying Recurrent and Unknown Performance Issues, in Proceedings of IEEE International Conference on Data Mining 2014 (ICDM 2014), Shenzhen, China, December 2014.

[5] Qiang Fu, Jian-Guang Lou, Qingwei Lin, Rui Ding, Dongmei Zhang, Zihao Ye, Tao Xie, Performance Issue Diagnosis for Online Service Systems, in Proceedings of 31st International Symposium on Reliable Distributed Systems (SRDS 2012), October, 2012.

[6] Rui Ding, Qiang Fu, Jian-Guang Lou, Qingwei Lin, Dongmei Zhang, Jiajun Shen, Tao Xie, Healing Online Service Systems via Mining Historical Issue Repositories, in Proceedings of the 27th IEEE/ACM International Conference on Automated Software Engineering (ASE 2012), Essen, Germany, September, 2012.