

Can *Silhouette Execution* solve the VM Bootstorm Problem?

by

Syed Aunn Hasan Raza

S.B., Computer Science and Engineering, M.I.T., May 2010

Submitted to the Department of Electrical Engineering and Computer
Science

in partial fulfillment of the requirements for the degree of

Master of Engineering in Electrical Engineering and Computer Science

at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

September 2011

© Massachusetts Institute of Technology 2011. All rights reserved.

Author
Department of Electrical Engineering and Computer Science
September 1, 2011

Certified by
Dr. Saman P. Amarasinghe
Professor
Thesis Supervisor

Accepted by
Dr. Christopher J. Terman
Chairman, Masters of Engineering Thesis Committee

Can *Silhouette Execution* solve the VM Bootstorm Problem?

by

Syed Aunn Hasan Raza

Submitted to the Department of Electrical Engineering and Computer Science
on September 1, 2011, in partial fulfillment of the
requirements for the degree of
Master of Engineering in Electrical Engineering and Computer Science

Abstract

Both server and desktop virtualization rely on high VM density per physical host to reduce costs and improve consolidation. In the case of *boot-storms*, such high VM density per host can be a problem....

(To be filled in)

Thesis Supervisor: Dr. Saman P. Amarasinghe

Title: Professor

Acknowledgments

I would like to thank Professor Saman Amarasinghe for his huge role in both this project and my wonderful undergraduate experience at MIT. Saman was my professor for 6.005 (Spring 2008), 6.197/6.172 (Fall 2009) and 6.035 (Spring 2009). These three exciting semesters not only convinced me of his unparalleled genius, but also ignited my interest in computer systems. Over the past year, as I have experienced the highs and lows of research, I have really benefited from Saman's infinite insight, encouragement and patience.

As an M-Eng student, I have been blessed to work with two truly inspirational and gifted people from the COMMIT group: Marek Olszewski and Qin Zhao. Their expertise and brilliance is probably only eclipsed by their humility and helpfulness. I have learned more from Marek and Qin than I probably realize, and their knowledge of operating systems and dynamic instrumentation is invaluable and, frankly, immensely intimidating. I hope to emulate (or even approximate) their excellence some day.

This past year, I have also had the opportunity to work with Professor Srinivas Devadas, Professor Fraans Kaashoek, and Professor Dina Katabi as a TA for 6.005 and 6.033. It has been an extraordinarily rewarding experience, and I have learned tremendously from simply interacting with these peerless individuals. Professor Dina Katabi was especially kind to me for letting me work in G916 over the past few months.

I would like to thank Abdul Munir, whom I have known since my first day at MIT; I simply don't deserve the unflinchingly loyal and supportive friend I have found in him. I am also indebted to Osama Badar, Usman Masood, Brian Joseph, and Nabeel Ahmed for their unrelenting support and encouragement; this past year would have been especially boring without our never-ending arguments and unproductive 'all-nighters'. I also owe a debt of gratitude to my partners-in-crime Prannay Budhraj, Ankit Gordhandas, Daniel Firestone and Maciej Pacula, who have been great friends and collaborators over the past few years.

I am humbled by the countless sacrifices made by my family in order for me to be where I am today. My father has been the single biggest inspiration and support in my life since childhood. He epitomizes, for me, the meaning of selflessness and resilience. This thesis, my work and few achievements were enabled by – and dedicated to – him, my mother and my two siblings Ali and Zahra. Ali has been a calming influence during my years at MIT; the strangest (and most unproductive) obsessions unite us, ranging from Chinese *Wuxia* fiction to, more recently, *The Game of Thrones*. Zahra’s high-school math problems have been a welcome distraction over the past year; they have also allowed me to appear smarter than I truly am.

Finally, I would like to thank my wife Amina for her unwavering love and support throughout my stay at MIT, for improving and enriching my life every single day since I have known her, and for knowing me better than even I know myself. Through her, I have also met two exceptional individuals, Drs. Fatima and Anwer Basha, whom I have already learnt a lot from.

“It is impossible to live without failing at something, unless you live so cautiously that you might as well not have lived at all – in which case, you fail by default.”

J.K. Rowling, Harvard Commencement Speech 2008

Contents

1	Introduction	13
1.1	Motivation	13
1.2	Goal of Thesis	15
1.3	Contributions	16
1.4	Importance of Deterministic Execution	17
1.5	Thesis Organization	18

List of Figures

List of Tables

Chapter 1

Introduction

1.1 Motivation

Large organizations increasingly use virtualization to consolidate server applications in data centers. This reduces operating costs, simplifies administrative tasks and improves performance scalability. One reason for the success of server virtualization is that it resolves the tension between typically conflicting goals of high isolation and effective resource utilization. Ideally, organizations prefer to assign each server application its own dedicated machine to achieve high isolation. However, each application typically utilizes only a modest fraction of a machine's hardware resources, so this can waste hardware resources. With the development of virtualization technology, applications can be assigned dedicated virtual machines (VMs) for isolation, while many such VMs can be hosted by the same physical host for high resource utilization.

The ability to consolidate numerous server applications on fewer physical hosts is so important to the success of server virtualization and is measured by *VM density per host*.

organizations can isolate different server applications by assigning dedicated virtual machines (VMs) to them.

For isolation, each server application is typically assigned a dedicated virtual machine (VM), but this can

Because a single VM would typically be desirable in most cases for effective resource

utilization.

Data centers increasingly use server virtualization to reduce operating costs, simplify administrative tasks and improve performance scalability. Through virtualization, it is possible to achieve high resource utilization and isolation at the same time: each application is typically assigned a dedicated server virtual machine (VM), while many VMs are consolidated on powerful host computers to reduce wasted cycles. The use of techniques such as memory overcommitment (including *transparent page sharing*, *ballooning* and *hypervisor swapping*) [18] has further improved the consolidation ratios and cost-effectiveness of server virtualization, and augurs well for the future of the technology.

Given the success of server virtualization, many companies are extending the use of virtualization to their desktop computers. In a Virtual Desktop Infrastructure [17] (VDI), desktop operating systems and applications are hosted in virtual machines that reside in a data center; users access virtual desktops from desktop PCs or thin clients via a remote display protocol. A VDI provides simplicity in administration and management: applications can be centrally added, deleted, upgraded and patched. VDI deployments also promise even higher consolidation ratios than those achieved via server virtualization because desktop virtual machines typically require less resources than server virtual machines.

Consolidation ratios (measured by VM density per host) in data centers are expected to increase in the future, not only because of improvements in virtualization technology, but also because newer generations of processors support more cores and memory [4].

However, correlated spikes in the CPU/memory usage of many VMs can suddenly cripple host machines. For instance, a *boot storm* [4, 6, 9, 14, 16] can occur after some software is installed or updated, requiring hundreds or thousands of identical VMs to reboot at the same time. Bootstorms can be particularly frequent in VDIs because users typically show up to work at roughly the same time in the morning each day.

Concurrently booting VMs create unusually high I/O traffic, generate numerous disk and memory allocation requests, and can saturate host CPUs. To avoid the

prohibitively high boot latencies that result from boot storms, data centers usually either boot machines in a staggered fashion, or invest in specialized, expensive and/or extra-provisioned hardware for network/storage [5, 6]. There is also anecdotal evidence that VDI users sometimes leave their desktop computers running overnight to avoid morning boot storms; this practice represents an unnecessary addition to already exorbitant data center energy bills [13]. Data deduplication [3], through which hosts reclaim/reuse disk blocks common to several VMs, has been proven to reduce the memory footprint of concurrently booting machines. However, while data deduplication can mitigate the stress on the memory subsystem in a boot storm, lowered memory latency can in turn overwhelm the CPU, fibre channel, bus infrastructure or controller resources and simply turn them into bottlenecks instead [10].

With the spread of virtualization, it is important to address the bootstorm problem in a way that does not involve simply skirting around the issue. Data deduplication is partly effective because identical VMs load the same data from disk when they boot up. In this thesis, we pose the following question: is it possible to generalize deduplication of data to deduplication of *execution*? If many identical VMs are concurrently booting up in a data center, do they execute the same set of instructions? Even if there are some differences in the instructions executed, are they caused by controllable sources of non-determinism? Ultimately, if there is a way to ensure that concurrently booting VMs execute mostly the same set of instructions and perform the same I/O requests, one way to solve the boot storm problem may be remarkable simple in essence: instead of booting N identical VMs concurrently, we can boot one VM as a leader; the remaining $(N - 1)$ VMs follow the leader by executing a tiny subset of the instructions they would otherwise execute; we split execution into N different instances as late as possible into the boot process. This approach could potentially reduce pressure on the underlying host hardware, and thereby enable data centers to handle boot storms effectively.

1.2 Goal of Thesis

This thesis aims to address the following questions:

1. When identical VMs boot up concurrently, how similar are the sets of instructions executed? What is the statistical profile of any differences in the executed instructions?
2. What are the source(s) of any differences in the instruction streams of concurrently booting VMs? Are there ways to minimize the non-determinism in booting VMs?

The answers to these questions are clearly crucial in determining the feasibility of *deduplication of execution* as a possible solution to the boot storm problem.

1.3 Contributions

For this work, we used dynamic instrumentation frameworks such as DynamoRio [2] and Pin [7] to study user-level instruction streams from a few representative Linux services at boot-time.

In this document, we:

1. quantify nondeterminism in Linux services, and show that it is bursty and rare;
2. document the sources of nondeterminism in Linux services – both obvious and obscure – and specify strategies for overcoming them in the boot storm scenario;
3. use simple dynamic instrumentation techniques to show that *fully* deterministic execution is achievable without *any* modifications to Linux or an executing service.

Strategies to achieve deterministic execution have been studied at the operating system layer [1] before, but they require modifications to Linux. Deterministic execution can be achieved in multi-threaded programs using record-and-replay approaches

[12] or deterministic logical clocks [11]. Our study of non-determinism has different goals from both approaches: we wish to avoid changing existing software (to ease adoption); we also wish to make several distinct – and potentially different – executions *overlap* as much as possible, rather than replay one execution over and over. In our case, we do not know *a priori* whether two executions will behave identically or not. That the behavior of system calls or signals in Linux can lead to different results or side-effects across multiple executions of an application is well known: what is not documented is the application *context* in which these sources of nondeterminism originate. To the best of our knowledge, this is the first attempt to study the statistical profile and context of nondeterminism in Linux services in such detail. While we hope this work ultimately proves the basis for an implementation of our proposed solution to the boot storm problem, we also note that deterministic execution can immediately improve the effectiveness of existing virtualization technologies such as transparent page sharing and data deduplication.

1.4 Importance of Deterministic Execution

While our study of nondeterminism is driven by a specific application, deterministic execution is desirable in a variety of scenarios. The motivations for deterministic multithreading listed in [11, 12] apply to our work as well.

Mainstream Computing, Security and Performance: If distinct executions of the same program can be expected to execute the same set of instructions, then any significant deviations can be used to detect security attacks. Runtime detection of security attacks through the identification of anomalous executions is the focus of *mainstream computing* [15], and deterministic execution obviously helps in reducing false positives. Anomalous executions can also be flagged for performance debugging.

Testing: Deterministic execution in general facilitates testing, because outputs and internal state can be checked at certain points with respect to expected values. Our

version of determinism allows for a particularly strong kind of test case that may be necessary for safety-critical systems: a program must execute the exact same instructions across different executions (for the same inputs).

Debugging: Erroneous behavior can be more easily reproduced via deterministic execution, which helps with debugging. Deterministic execution has much lower storage overhead than traditional record-and-replay approaches.

1.5 Thesis Organization

In what follows, Chapter ?? presents an overview of the Linux boot process, along with the dynamic instrumentation techniques we used to profile non-determinism in Linux services. Chapter ?? presents a summary of the sources of nondeterminism discovered in this work and the strategies we used to eliminate them. Chapter ?? presents a detailed case study of three Linux services to identify the common context in which non-determinism arises. Chapter ?? presents design ideas for an implementation of deduplication of execution. Finally, Chapter ?? concludes this thesis and discusses future work.

Bibliography

- [1] T. Bergan, N. Hunt, L. Ceze, and S.D. Gribble. Deterministic process groups in dos. *9th OSDI*, 2010.
- [2] D.L. Bruening. *Efficient, transparent, and comprehensive runtime code manipulation*. PhD thesis, Citeseer, 2004.
- [3] A.T. Clements, I. Ahmad, M. Vilayannur, and J. Li. Decentralized deduplication in san cluster file systems. In *Proceedings of the 2009 conference on USENIX Annual technical conference*, pages 8–8. USENIX Association, 2009.
- [4] J.G. Hansen and E. Jul. Lithium: virtual machine storage for the cloud. In *Proceedings of the 1st ACM symposium on Cloud computing*, pages 15–26. ACM, 2010.
- [5] Solving Boot Storms With High Performance NAS. http://www.storage-switzerland.com/Articles/Entries/2011/1/3_Solving_Boot_Storms_With_High_Performance_NAS.html, 2011. [Accessed 1-August-2011].
- [6] X.F. Liao, H. Li, H. Jin, H.X. Hou, Y. Jiang, and H.K. Liu. Vmstore: Distributed storage system for multiple virtual machines. *SCIENCE CHINA Information Sciences*, 54(6):1104–1118, 2011.
- [7] C.K. Luk, R. Cohn, R. Muth, H. Patil, A. Klauser, G. Lowney, S. Wallace, V.J. Reddi, and K. Hazelwood. Pin: building customized program analysis tools with dynamic instrumentation. In *ACM SIGPLAN Notices*, volume 40, pages 190–200. ACM, 2005.
- [8] Bootchart: Boot Process Performance Visualization. <http://www.bootchart.org>, 2011. [Accessed 29-July-2011].
- [9] S. Meng, L. Liu, and V. Soundararajan. Tide: achieving self-scaling in virtualized datacenter management middleware. In *Proceedings of the 11th International Middleware Conference Industrial track*, pages 17–22. ACM, 2010.
- [10] VMware Bootstorm on NetApp. <http://ctistrategy.com/2009/11/01/vmware-boot-storm-netapp/>, 2009. [Accessed 29-July-2011].
- [11] M. Olszewski, J. Ansel, and S. Amarasinghe. Scaling deterministic multithreading. *2nd WoDet*, 2011.

- [12] H. Patil, C. Pereira, M. Stallcup, G. Lueck, and J. Cownie. Pinplay: a framework for deterministic replay and reproducible analysis of parallel programs. In *Proceedings of the 8th annual IEEE/ACM international symposium on Code generation and optimization*, pages 2–11. ACM, 2010.
- [13] A. Qureshi, R. Weber, H. Balakrishnan, J. Guttag, and B. Maggs. Cutting the electric bill for internet-scale systems. In *Proceedings of the ACM SIGCOMM 2009 conference on Data communication*, pages 123–134. ACM, 2009.
- [14] Vijayaraghavan Soundararajan and Jennifer M. Anderson. The impact of management operations on the virtualized datacenter. *SIGARCH Comput. Archit. News*, 38:326–337, June 2010.
- [15] M.W. Stephenson, R. Rangan, E. Yashchin, and E. Van Hensbergen. Statistically regulating program behavior via mainstream computing. In *Proceedings of the 8th annual IEEE/ACM international symposium on Code generation and optimization*, pages 238–247. ACM, 2010.
- [16] S.B. Vaghani. Virtual machine file system. *ACM SIGOPS Operating Systems Review*, 44(4):57–70, 2010.
- [17] VMware Virtual Desktop Infrastructure. http://www.vmware.com/pdf/virtual_desktop_infrastructure_wp.pdf, 2011. [Accessed 29-July-2011].
- [18] C.A. Waldspurger. Memory resource management in vmware esx server. *ACM SIGOPS Operating Systems Review*, 36(SI):181–194, 2002.