Atmosphere and the TCP/IP Stack for Maximizing Security

Undergraduate Thesis

Submitted in partial fulfillment of the requirements of BITS F421T Thesis

Ву

Tirth Jain ID No. 2019A7TS0120P

Under the supervision of:

Dr. Anton Burtsev

&

Dr. K Hari Babu



BIRLA INSTITUTE OF TECHNOLOGY AND SCIENCE PILANI, PILANI CAMPUS November 2022

Declaration of Authorship

I, Tirth Jain, declare that this Undergraduate Thesis titled, 'Atmosphere and the TCP/IP Stack for Maximizing Security' and the work presented in it are my own. I confirm that:

- This work was done wholly or mainly while in candidature for a research degree at this University.
- Where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated.
- Where I have consulted the published work of others, this is always clearly attributed.
- Where I have quoted from the work of others, the source is always given. With the exception of such quotations, this thesis is entirely my own work.
- I have acknowledged all main sources of help.
- Where the thesis is based on work done by myself jointly with others, I have made clear exactly what was done by others and what I have contributed myself.

Signed:			
Date:			

Certificate

This is to certify that the thesis entitled, "Atmosphere and the TCP/IP Stack for Maximizing Security" and submitted by <u>Tirth Jain</u> ID No. <u>2019A7TS0120P</u> in partial fulfillment of the requirements of BITS F421T Thesis embodies the work done by him under my supervision.

Supervisor
Dr. Anton Burtsev
Associate Professor,
University of Utah

Date:

Co-Supervisor

Dr. K Hari Babu

Asst. Professor,

BITS Pilani Pilani (

BITS Pilani Pilani Campus

Date:

BIRLA INSTITUTE OF TECHNOLOGY AND SCIENCE PILANI, PILANI CAMPUS

Abstract

Bachelor of Engineering (Hons.)

Atmosphere and the TCP/IP Stack for Maximizing Security

by Tirth Jain

Even after decades of work to make monolithic kernels more secure, serious vulnerabilities in them are still reported every year. Because the entire monolithic kernel is in one address space, an attacker is just one vulnerability away from owning the entire machine. We argue that it is time to decompose monolithic kernels like Linux into smaller parts that run in isolated compartments and communicate using secure interfaces. We think this is timely due to recent trends in hardware that make it easier and efficient to isolate kernel components.

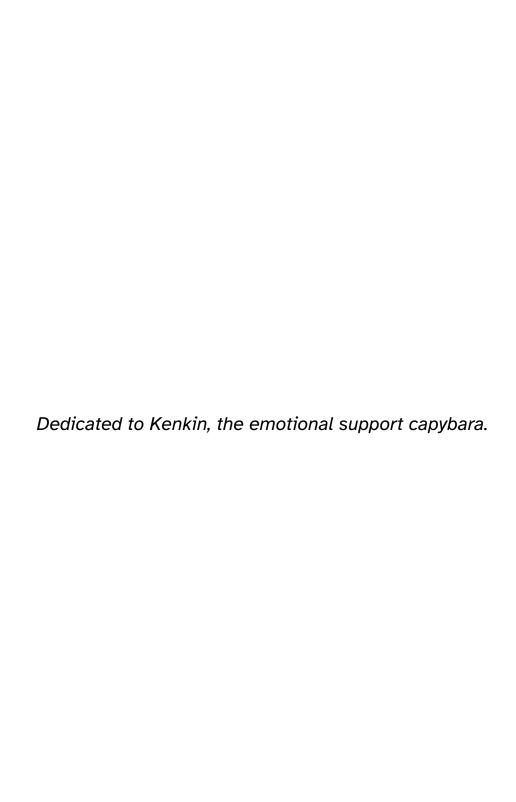
In this Thesis, I discuss Atmosphere, our approach to microkernelization and fault isolation in Operating Systems. Specifically, I will focus on the network stack we built for Atmosphere. We argue that the network stack is a source for bugs and that isolation is the way forward to minimize the impact of these bugs.

Acknowledgements

Thanks to Anton for calling me to Utah and giving me this opportunity. Thanks to Prof. Hari Babu for teaching me NetProg. Thanks to Xiangdong for dragging me to the gym everyday and all of his compilers. Thanks to Zhaofeng for all the Nix shells and gatekeeping all our code. Without all of them, this Thesis wouldn't have been possible.

Contents

De	eclaration of Authorship	i			
Ce	ertificate	ii			
Abstract					
Ac	cknowledgements	iv			
Co	ontents	V			
1	Introduction	1			
2	Background and Related Work 2.1 Theseus	3			
	2.2 Singularity OS	4			
	2.3 sel4	4			
	2.4 RedLeaf	4			
3	Atmosphere Design	5			
	3.1 Safe IPC	5			
4	The Network Stack	6			
5	Evaluation	7			
Ri	bliography	8			



Introduction

Operating system reliability is a topic that has been studied for decades but remains a major concern till date. Even thought its been 30 years since Linux kernel was developed, critical vulnerabilities are being found till date. On top of that, the monolithic design (ie everything in the kernel running in a single address space) makes it such that the attacker is only one exploit away from taking control of the entire operating system. This is concerning because device drivers, often developed by third parties are a major source of vulnerabilities. Redundancy in hardware and software protection mechanisms can protect against transient faults but persistent faults because of undetected logic flaws still remain a major concern. In recent times, with the advent of cloud services, this becomes even more significant because clouds often have to run untrusted code that needs to be isolated from other parts of the operating system.

Logic errors and faults can propogate from within a component and propogate to other parts of the system. Microkernelization aims to solve this problem by isolating various parts of an operating system and minimizing the effect of failures in a single subsystem. In a microkernel, only the bare minimum that is required to boot an operating system is included in the kernel. All other components such as the filesystem, the device drivers, and the network stack are run in separate isolated processes that interact with each other using some form of IPC. Inter component fault propogation can thus be reduced by introducing safe IPC mechanisms.

Atmosphere is a microkernel based operating system. In this thesis, I discuss the design of Atmosphere mainly focusing on the the network stack and how its design can be used to model other services of the operating system. The network stack was built two main principles in mind: maximizing fault isolation and minimizing its cost on performance. For the first principle, we take microkernelization to its extremes and build a per connection (or per socket for UDP) network stack. This means, every new connection has its own TCP/IP stack that can process and dispatch packets. This ensures that if one connection is corrupted, the rest of the connections on the system can keep operating. Although this is not completely possible since there is some inherent

shared state on any host machine on a network. In later sections we discuss how we minimize this shared state. For the second principle, we implement a zero copy model and explain how using Rust's guarantees help us implement this with confidence. On failure, a component can be safely restarted and meanwhile, all IPC calls to it can return an error.

The rest of this Thesis is organised as follows:

- In Background and Related Work we take a look at the contemporary approaches to microkernelization.
- Atmosphere Design discusses the design principles of Atmosphere.
- The Network Stack talks about the API design and the internals of the network stack we built.
- Evaluation evaluates the performance of the network stack as compared to the current standards.

Background and Related Work

Tirth some intro here? idk

2.1 Theseus

Theseus[2] is a microkernel based operating system written in Rust. Rust's safety guarantees enable Theseus to run all software written in Rust, including userspace applications, to run in a Single Address Space system with a Single Privilege level. Thus, eliminating the need for virtual memory management and protection rings. Theseus introduces the idea of cells that are described as a software-defined unit of modularity that serves as the core building block of the OS. A cell can be modeled as a crate in Rust. On booting, only the microkernel, ie, the nano_core is loaded which bootstraps the system. All other cells are dynamically loaded on demand. Theseus piggybacks on Rust's safety guarantees to enable reliable IPC between cells. For example, memory mappings are a 1-1 mapping to a physical frame. They can only be shared behind a read-only &MappedPages reference eliminating the double-free and the use-after-free problem.

The most important contribution by Theseus is the idea of eliminating (or minimizing) shared state between components. All OS services can be modelled as servers and applications requesting these services are modelled as clients. Theseus eliminates state-spill[1] across cells by eliminating all state from the servers. Everything that is needed to service a client's request is stored in the client itself and thus the servers can be completely stateless. This allows a server cell to be replaced by a new cell without any state loss. This, however, is not always possible. In case a state cannot be eliminated, as in the case of descriptor tables, the state is stored in a state_db. The state_db is a key-value database that stores states with a static lifetime that the server can request a (weak) reference to. In case a server needs to be restarted, its state can be recovered from the statedb. The statelessness of cells also allows for live updates. A server

can be replaced, ie, a patch can be applied to a component without having to restart the entire operating system. Applying a patch is as simple as swapping a cell with a new cell.

Tirth Maybe add limitations and performance evaluations of Theseus here?

2.2 Singularity OS

2.3 sel4

2.4 RedLeaf

Atmosphere Design

3.1 Safe IPC

Tirth Describe RedLeaf RRef interface.

The Network Stack

Evaluation

Bibliography

- [1] Kevin Boos, Emilio Del Vecchio, and Lin Zhong. "A Characterization of State Spill in Modern Operating Systems". In: *Proceedings of the Twelfth European Conference on Computer Systems*. EuroSys '17. Belgrade, Serbia: Association for Computing Machinery, 2017, 389–404. isbn: 9781450349383. doi: 10.1145/3064176.3064205. url: https://doi.org/10.1145/3064176.3064205.
- [2] Kevin Boos et al. "Theseus: an Experiment in Operating System Structure and State Management". In: 14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20). USENIX Association, Nov. 2020, pp. 1–19. isbn: 978-1-939133-19-9. url: https://www.usenix.org/conference/osdi20/presentation/boos.