

# ALI RAZA

[alirazabhutta.10@gmail.com](mailto:alirazabhutta.10@gmail.com) - 857-350-5959

## OVERVIEW

---

Interests: AI/ML System Level Optimizations, HPC, Linux, Operating Systems, Low-level Debugging, Computer Networks, Cloud Computing, Distributed Systems, Performance Analysis and Application Tuning

Skills: Linux Kernel Programming and Debugging, Systems and Network Programming, C, Python, Bash

## EXPERIENCE

---

Mar 2023 - Present

**Senior Engineer, AI and HPC System Architecture**  
SAMSUNG SEMICONDUCTOR INC., USA

- Explored optimizations for distributed AI training frameworks like PyTorch by enhancing the MPI library to support shared memory operations, enabling efficient utilization of emerging memory-coupled compute and global shared memory architectures.
- Reduced inter-process communication overhead and developed system software for novel shared memory architectures, ensuring seamless integration and readiness for scalable AI workloads, including transformer-based architectures like GPT.
- Addressed noise issues in Linux-based environments to enhance stability and predictability for AI/ML and high-throughput compute environments.
- Authored a patent for an innovative threading model to offload computations to accelerators, potentially enhancing performance in AI training and inference workflows.

## EDUCATION

---

Sep 2017 - Mar 2023

**Ph.D. Computer Science**  
BOSTON UNIVERSITY  
Advisor: [Prof. Orran Krieger](#)

Sep 2014 - Jun 2016

**M.S. Computer Science**  
LAHORE UNIVERSITY OF MANAGEMENT SCIENCES, PAKISTAN  
Advisor: [Prof. Ihsan Ayyub Qazi](#)

Sep 2008 - Jun 2012

**B.S. Electrical Engineering**  
LAHORE UNIVERSITY OF MANAGEMENT SCIENCES, PAKISTAN  
Senior Year Project Supervisor: [Prof. Sohaib Khan](#)

## SELECTED PROJECTS

---

### PhD Thesis - Unikernel Linux (UKL)

This research, which I lead, explores if unikernel techniques (e.g., highly optimized transitions between kernel and application code, customized code paths, run-to-completion,

link-time optimization of kernel and application code, etc.) can be integrated into Linux while preserving its battle-tested code, development community, and ecosystem of tools, applications, and supported hardware. Our prototype, called Unikernel Linux (UKL) demonstrates that it is possible, and can support a large class of unmodified applications and hardware, and result in significant performance improvements. Complex unmodified applications (e.g., Redis, Memcached) can improve by over 10% in 99th tail latency and throughput, and small modifications to the applications can result in more than 20% improvement. The core UKL 550 line patch was shared with the community as an RFC (link below); the full set of optimizations is just over 1200 LoC.

**Github Link:** <https://github.com/unikernelLinux/ukl>

**UKL patch:** <https://lore.kernel.org/lkml/20221003222133.20948-1-aliraza@bu.edu/>

### **Elastic Secure Datacenter**

This research, which I participated in at the start of my PhD, explored how we can securely and rapidly multiplex physical servers between many different tenants while minimizing trust in the provider. The Bolted prototype we developed enabled tenants to use existing unmodified provisioning systems for HPC, enterprise and cloud services. Security sensitive tenants can deploy their own attestation services to ensure that previous tenants did not compromise the firmware. Keylime attestation service used in Bolted has developed a large open source community and is being used in both Red Hat products and IBM's public cloud. The fundamental mechanisms from this work have now been adapted by the OpenStack community, and are currently being integrated and tested at the MGHPC data center.

### **Context Aware WiFi Bitrate Adaptation**

Indoor wireless channel conditions suffer from signal attenuation, multi-path interference and shadowing effects, decreasing Wi-Fi throughput, especially when high data rates are used. On the contrary, in outdoor settings, better channel conditions allow higher bit rates, thus improving throughput. For mobile devices, especially those which frequently transition between indoor and outdoor settings, WiFi bitrates needs to quickly adapt to channel conditions. But the WiFi protocol (IEEE 802.11 family) estimates channel conditions over timescales of hundreds of packets or needs to be trained to understand the SNR-BER relationship, and is better suited to stationary devices. As part of my MS thesis, I designed an indoor-outdoor detection system which uses sensors present in mobile devices such as accelerometer, light sensor and GPS etc., and use this information to quickly adapt WiFi bitrate. The results show faster WiFi bitrate adaptation and better average throughput for a mobile device.

## **PUBLICATIONS**

---

### **Unikernel Linux (UKL) - EuroSys '23 [PDF]**

This work explores if unikernel techniques can be integrated into a general-purpose OS while preserving its battle-tested code, development community, and ecosystem of tools, applications, and hardware support. Our prototype demonstrates both a path to integrate unikernel techniques in Linux and that such techniques can result in performance advantages for unmodified applications. Expert developers can modify the application to call internal kernel functionality and optimize across the application/kernel boundary for more significant gains. The changes to the Linux kernel are modest (1250 LOC).

### **Unikernels: The Next Stage of Linux's Dominance - HotOS '19 [PDF]**

Unikernels have demonstrated advantages over Linux in many important domains, but have lagged in adoption and ease of use. They lack the battle-tested code base

and large open source community which Linux has. In this paper, we explore if Linux can be built into a unikernel, and would that give us any performance benefits. We build an early Linux unikernel prototype and demonstrate that some simple changes can bring dramatic performance advantages.

#### **A Secure Cloud with Minimal Provider Trust - HotCloud 18 [\[PDF\]](#)**

In this paper we propose Bolted which is a new architecture for a bare metal cloud with the goal of providing security-sensitive customers of a cloud the same level of security and control that they can obtain in their own private data centers. It allows tenants, rather than the provider, to control the trade-offs between security, price, and performance. A prototype demonstrates scalable end-to-end security with a small overhead compared to a less secure alternative.

#### **It's All in the Name: Why Some URLs are More Vulnerable to Typosquatting - IEEE INFOCOM 2018 [\[PDF\]](#)**

Typosquatting is a blackhat practice that relies on human error and low-cost domain registrations to hijack legitimate traffic from well-established websites. We explore the relationship between human hand anatomy, keyboard layouts and typing mistakes to understand why some URLs are more vulnerable to typing mistakes than others. Furthermore, we predict the most likely typos for each URL which can then be defensively registered.

#### **An Anomaly Detection Fabric for Clouds Based on Collaborative VM Communities - CCGrid 2017 [\[PDF\]](#)**

To mitigate security threats in the cloud, we propose a hypervisor layer anomaly detection system based on system call monitoring, which compresses the stream of system calls at their source making the system scalable and near real-time. Not requiring modifications to the guest OS or the application make it ideal for the datacenter setting. Additionally, for robust and early detection of threats, we share information about attacks in their early stages to provide immunity to the entire deployment.

## **TALKS AND PRESENTATIONS**

---

- DevConf.us 2021 [video](#)
- DevConf.us 2020 [video](#)
- Red Hat Research Day 2020 [video](#)
- DevConf.us 2019 [video](#)
- DevConf.us 2018 [video](#)