

Reto Achermann

Faithful Hardware Representation and Least-Privilege Memory Management in Operating Systems

University of Texas at Austin, December 6, 2019

PhD Student, Systems Group, Department of Computer Science, ETH Zurich



My Operating Systems Coding Experiences



- Platform support: Cavium ThunderX / Xeon Phi coprocessor / ARM FastModels
- Drivers: USB Stack / DMA Engines / Xeon Phi / CPU / IOMMU
- Runtimes: OpenMP / libnuma / message-passing
- Extensions to the capability system
- Domain Specific Languages

- Memory management
- Message-passing system
- Memory allocation policies
- Page-table replication
- Kernel modules



A deep **dissatisfaction** about the way operating systems abstract and represent hardware.

Computer Architecture 101

Figure 1.4 Hardware organization of a typical system. CPU: Central Processing Unit, ALU: Arithmetic/Logic Unit, PC: Program counter, USB: Universal Serial Bus.



systems, but all systems have a similar look and feel. Don't worry about the complexity of this figure just now. We will get to its various details in stages throughout the course of the book.

The OS just runs on a new platform

- A single, flat physical address space
 - physical address as a unique identifier
- Physical address means the same for all CPUs and devices

Reality: Hardware Violates the Assumptions Made by the OS

6+ heterogeneous cores A9 A9 DSP Next, we zoom in here 5+ Interconnects **Firewalls Devices attached** to different interconnects A heterogeneous network of cores, interconnects, devices and memory.

TexasInstruments OMAP4460, Q4 2011

EHzürich

Ambiguous Physical Addresses and Non-Uniform Views

OMAP 4460 SoC

A9 Like an alarm clock, set a time There are multiple physical addresses for the same device! 0x40138000/12 0x49038000/12 0x01D38000/12 0x00 0x01D38000/12 0x00 0x001D38000/12 0x00 0x001D38000/12 0x00 0x000 0x00 0x

General Purpose Timer

Problem: There is no right address.



EHzürich

Complicated Memory Topologies are a Universal Problem



Secure and non-secure Co-Processors Segmentation Segmentation (ARM TrustZone / Intel SGX) Addressetranslations are (configurable.

Direct Segments



Mismatch: Hardware Abstraction in OS Real Hardware

This mismatch is a **problem** – the OS does not seem to get it right

Bugs and vulnerabilities in systems software:

- CVE-2014-3601: Miscalculation of affected pages
- CVE-2016-5349: Not enough memory address information provided
- CVE-2017-8061: Wrong DMA addresses
- CVE-2017-16994: Ignoring holes in huge-pages
- CVE-2014-9888: wrong access rights for data pages
- CVE-2019-2250: authorization bug allows writing to memory locations
- CVE-2018-11994: SMMU misconfiguration allows access to memory
- 30% of patches to Linux memory manager are bugfixes

My Research Focus

My goal: A sound basis for reasoning about reliable operating system on any hardware platform.



Design and Implementation of Operating Systems



Domain Specific Languages for Systems Engineering



Applying Formal Methods to Operating Systems and Hardware Specification





Faithful Hardware Representation and Least-Privilege Memory Management in Operating Systems





A new model to express memory addressing on modern machines

Generate OS code and hardware configuration



Express changes of the configuration and the required authority



Efficient Implementation in an Operating System





Faithful Hardware Representation and Least-Privilege Memory Management in Operating Systems



A new model to express memory addressing on modern machines



Generate OS code and hardware configuration



Express changes of the configuration and the required authority



Efficient Implementation in an Operating System

A More Realistic View of the Platform

Figure 1.4 Hardware organization of a typical system. CPU: Central Processing Unit, ALU: Arithmetic/Logic Unit, PC: Program counter, USB: Universal Serial Bus.





The OS needs to be ported:

- Requires engineering effort
- Manual process: source of bugs (more engineering effort)

Configurable I/O bus:

- PCI bridge programming
- PCI hot-plug
- I/O MMUs / System MMUs
- Virtualization

Interconnects:

- Configurable, multi-stage translations
- Configurable Firewalls
- Private access ports

Processors

- Core-Local resources
- Cache Hierarchy, operation modes
- Configurable Multi-Stage Translation
- Virtualization

Devices

- Memory access restrictions
- Configurable Translations
- Self-virtualization

The OS must correctly configure and manage these

Specification and Modeling Hardware

- Industry Standards: DeviceTrees / UEFI / ACPI / USB / PCI Express
 - Limited topology information, not available everywhere
 - Assumptions: a uniform view & unique physical addresses hardware representation
- Memory & Processor Models: ARM's ASL and Sewell et al.
 - Model the behavior of instructions and memory requests
 - Stop at processor boundaries
- Verified Operating Systems: seL4, CertiKOS
 - Proofs based on a linear flat array to physical memory

Complimentary problem orthogonal to address translation

Proofs need an accurate

hardware representation



The Address Space Abstraction – A More Faithful Representation

- Memory management needs an unambiguous reference to physical resources.
- Address Space (Naming Problem)
 - Context for resolving addresses
 - Range of address values e.g. [0, 2^b)
 - Regions of address mappings
 - Regions of local resources



ETH zürich NUMA Node 0

NUMA Node 1





Block Diagram of a heterogeneous, two-socket server

2x Intel Xeon E5 processor 2x 128GB RAM

Intel Xeon Phi Co-Processor attached to PCI Express

Decomposing a Heterogeneous Server into Address Spaces



Applying Formal Methods to the Address Space Model

- Formalization of the model in Isabelle/HOL [Decoding Net, MARS'17 / ITP'18]
- Well-defined semantics of address resolution termination proofs, ...
- Verification of algorithms on top of the model
- Capture the current, static state of the system



Sound foundation to express address translation of real hardware (e.g. TLB models)

EHzürich

"What is reachable from a core and at which address?"



EHzürich

pera

Syst

Using the Correct Output to write Platform Specific OS Code

New Platform Repeat the Process

Raspberry Pi 4

Your tiny, dual-display, desktop computer ...and robot brains, smart home hub, media centre, networked Al core, factory controller, and much more

https://www.raspberrypi.org/products/raspberry-pi-4-model-b/

Many new SoC-Platforms are released every year

System-on-a-Chip Platform Vendors

Actions Semiconductor, Advanced Micro Devices (AMD), Advanced Semiconductor Engineering (ASE), Aeroflex Gaisler, Agate Logic, Alchip, Allwinner Technology, Altera, Amkor Technology, Amlogic, Analog Devices, Anyka, Apple Inc., Applied Micro Circuits Corporation (AMCC), ARM Holdings, ASIX Electronics, Atheros, Atmel, Axis Communications, Broadcom, Cambridge Silicon Radio, Cavium Networks, CEVA, Inc., Cirrus Logic, Conexant, Core Logic, Coronis (Wavenis Technology), Cortina Systems, CPU Tech, Cypress Semiconductor, FameG (Fulhua Microelectronics Corp.), Freescale Semiconductor, Frontier Silicon Ltd, Fujifilm, HiSilicon, Horizon Semiconductors, Imagination Technologies, Infineon Technologies, Innova Card, Intel Corporation, InvenSense, Lattice Semiconductor, Leadcore Technology, LSI Corporation, Marvell Technologies, Mistletoe Technologies, Maxim Integrated Products, Milkymist, MIPS Technologies, Mistletoe Technologies, Nokia, NuCORE Technology, Nufront, NVIDIA, NXP Semiconductors (formerly

SoC Released 2018 Apple A12, S4, W3

Samsung

Exynos 9 Series, Exynos 7 Octa, Exynos 5 Hexa

QualComm

SDM439, SDM429, SDM632, SDM670, SDM710, SDM845, QCC5120

Philips Semiconductors), Corporation, PMC-Sien Scorpion, Redpine Sign Sequence Design, Shar Silicon Motion, Skywor SolidRun, Spreadtrum, Sunplus Technology, Sy Tensilica, Teridian Semi TranSwitch, Vimicro, RDA Microelectronics



Each platform is different. Operating system needs to be adapted every time. Manual porting: - time effort

- source of bugs...

T6755S, Helio Helio P70





Faithful Hardware Representation and Least-Privilege Memory Management in Operating Systems



A new model to express memory addressing on modern machines



Generate OS code and hardware configuration



Express changes of the configuration and the required authority



Efficient Implementation in an Operating System

EHzürich

Automatic Generation of OS Code From the Model

DATASHEET



Convright 2019 Raspherry Pi (Trading) Ltd. All rights reserve

Vendor supplied data (e.g. Hardware Manual)



Machine readable description of the platform



Executable representation of the model + algorithms

Domain Specific Language for Sockeye Hardware Descriptions

Generation of correct low-level OS code



ETH zürich

Toolchain



Use Case Example: Correct-by-Construction Page-Table Generation

Specify the observing core Flatten the decoding net using CPU Core Local view-preserving operations CPU Socket **CPU Socke** RAM IOMMU System Interconnect RAM SMPT PCI Root Window PCI Root Window IOMMU Co-Processor Socket PCI Bridge PCI Bridge x57 CPU Core CPU GDDR Memory Device Device Device Core Core Registers Local Registers Topology as generated ARMv7, "Local Address Space" by Sockeye ARMv8, x86 64, Device DRAM Registers K10M Core Local GDDR DRAM Resources Page Table Receipt Page-Table Binary Architecture = [ARMv8]Mappings = [DRAM0 @ 0x8000000, Generate page-tables based on DevRegs @ 0x1000000 the flattened representation

Validation: Custom Simulation Platforms for ARM FastModels







Faithful Hardware Representation and Least-Privilege Memory Management in Operating Systems



A new model to express memory addressing on modern machines







Express changes of the configuration and the required authority



Efficient Implementation in an Operating System

Dynamic Configuration of the Decoding Net



Semantics of **dynamic** configuration of the **translate** function



Rights and authority required to **change** the translate function



Principle of Least Privilege

From the Static Model to a Dynamic Implementation in an OS

1. Identification of all relevant **objects** and **subjects** and their relationship

2. Development of an **executable specification** for rapid prototyping

3. Executable specification as a guide for an efficient **OS implementation**

Fine-Grained Decomposition of Rights and Operations



Expressing Fine-Grained Authority in a System



Objects Subjects	DRAM Region	Core Address Space	•••
•••			
Process	Grant	Мар	
•••			

The access control matrix defines what address space configurations and transitions are valid.



From the Static Model to a Dynamic Implementation in an OS

- ✓ Identification of all relevant **objects** and **subjects** and their relationship
- Development of an **executable specification** for rapid prototyping
- Executable specification as a guide for an efficient **OS implementation**
 - Barrelfish with Multiple Address Spaces + Least-Privilege
 - Access control with Capabilities as a natural match for least privilege
 - Support for heterogeneous platforms





Faithful Hardware Representation and Least-Privilege Memory Management in Operating Systems



A new model to express memory addressing on modern machines

Generate OS code and hardware configuration



Express changes of the configuration and the required authority



Efficient Implementation in an Operating System

Evaluation

Multiple address spaces & least-privilege access control. What's the cost of implementing this in an operating system?

- 1) What is the cost of memory management operations?
- 2) What is the overhead for dynamic address space configuration following the least-privilege principle?

E *H* zürich

Evaluation of virtual memory management operations



The Cost of Dynamic Reconfiguration with Least-Privilege

Task: Setup a shared buffer between the host CPU and the co-processor

- 1) Invoke model to obtain
 - memory resources to allocate
 - list of address spaces to configure

Allocation & mapping of memory
Configuration of address translation



Two-Socket server with Xeon Phi Co-Processor

The Cost of Dynamic Reconfiguration with Least-Privilege







Faithful Hardware Representation and Least-Privilege Memory Management in Operating Systems



A new model to express memory addressing on modern machines

Generate OS code and hardware configuration



Express changes of the configuration and the required authority



Efficient Implementation in an Operating System

Future Directions







Vision

Apply the same approach to other areas of the operating system to obtain correct and reliable system software running on any platform.

Combining OS design & implementation with programming languages, code synthesis and Formal Methods.

ETH zürich

Thanks to my collaborators

Timothy Roscoe	Lukas Humbel	Nora Hossle	David Cock	Daniel Schwyn
Simon Gerber	Kornilios Kourtis	Dejan Milojicic	Stefan Kaestle	Tim Harris
Gerd Zellweger	Roni Haecki	Moritz Hoffmann	Sabela Ramos	Jayneel Gandhi
Izzat El Hajj	Alexander Merritt	Contributors to the Barrelfish Operating System		

List of Related Publications

- R. Achermann, A. Panwar', J. Gandhi, A. Bhattacharjee, T. Roscoe. Mitosis: Transparently Self-Replicating Page-Tables for Large-Memory Machines (ASPLOS20)
- R. Achermann, N. Hossle, L. Humbel, D. Schwyn, D. Cock, T. Roscoe. A Least-Privilege Memory Protection Model for Modern Hardware. (ArXiv)
- L. Azriel, L. Humbel, **R. Achermann**, A. Richardson, M. Hoffmann, A. Mendelson, T. Roscoe, RN. Watson, P. Faraboschi, D. Milojicic D. *Memory-side protection with a capability enforcement co-processor*. (TACO).
- R. Achermann, L. Humbel, D. Cock, T. Roscoe. *Physical addressing on real hardware in Isabelle/HOL*. (ITP'18).
- L. Humbel, R. Achermann, D. Cock, T. Roscoe. Towards Correct-by-Construction Interrupt Routing on Real Hardware. (PLOS'17).
- R. Achermann, C. Dalton, P. Faraboschi, M. Hoffmann, D. Milojicic, G. Ndu, A. Richardson, T. Roscoe, A. L. Shaw; R. N. M. Watson. Separating Translation from Protection in Address Spaces with Dynamic Remapping. (HOTOS'XVI).
- R. Achermann, L. Humbel, D. Cock and T. Roscoe. Formalizing Memory Accesses and Interrupts. (MARS 2017).
- S. Kaestle, R. Achermann, R. Haecki, M. Hoffmann, S. Ramos, and T. Roscoe. Machine-Aware Atomic Broadcast Trees for Multicores. (OSDI'16).
- I. El Hajj, A. Merritt, G. Zellweger, D. Milojicic, **R. Achermann**, W. Hwu, K. Schwan, T. Roscoe, P. Faraboschi. *SpaceJMP: Programming with Multiple Virtual Address Spaces*. (ASPLOS XXI).
- S. Kaestle, R. Achermann, T. Roscoe, T. Harris. Shoal: Smart Allocation and Replication of Memory For Parallel Programs (ATC'15)
- S. Gerber, G. Zellweger, R. Achermann, K. Kourtis, T. Roscoe, D. Milojicic. Not Your Parents' Physical Address Space. (HotOS XV).