



## Capability Hardware Enhanced RISC Instructions: CHERI Instruction-set architecture

Robert N. M. Watson, Peter G. Neumann,  
Jonathan Woodruff, Jonathan Anderson,  
David Chisnall, Brooks Davis, Ben Laurie,  
Simon W. Moore, Steven J. Murdoch,  
Michael Roe

December 2014

15 JJ Thomson Avenue  
Cambridge CB3 0FD  
United Kingdom  
phone +44 1223 763500  
<http://www.cl.cam.ac.uk/>

© 2014 Robert N. M. Watson, Peter G. Neumann,  
Jonathan Woodruff, Jonathan Anderson, David Chisnall,  
Brooks Davis, Ben Laurie, Simon W. Moore,  
Steven J. Murdoch, Michael Roe, SRI International

Approved for public release; distribution is unlimited.  
Sponsored by the Defense Advanced Research Projects  
Agency (DARPA) and the Air Force Research Laboratory  
(AFRL), under contract FA8750-10-C-0237 (“CTSRD”) as  
part of the DARPA CRASH research program. The views,  
opinions, and/or findings contained in this report are those of  
the authors and should not be interpreted as representing the  
official views or policies, either expressed or implied, of the  
Department of Defense or the U.S. Government. Portions of  
this work were sponsored by the RCUK’s Horizon Digital  
Economy Research Hub grant, EP/G065802/1. Portions of  
this work were sponsored by Google, Inc.

Technical reports published by the University of Cambridge  
Computer Laboratory are freely available via the Internet:

*<http://www.cl.cam.ac.uk/techreports/>*

ISSN 1476-2986

## Abstract

This technical report describes CHERI ISAv3, the third version of the Capability Hardware Enhanced RISC Instructions (CHERI) Instruction-Set Architecture (ISA) being developed by SRI International and the University of Cambridge. This design captures four years of research, development, refinement, formal analysis, and testing, and is a substantial enhancement to the ISA version described in UCAM-CL-TR-850. Key improvements lie in tighter C-language integration, and more mature support for software object-capability models; these changes result from experience gained in adapting substantial software stacks to run on prototype hardware.

The CHERI instruction set is based on a *hybrid capability-system architecture* that adds new capability-system primitives to a commodity 64-bit RISC ISA enabling software to efficiently implement *fine-grained memory protection* and a hardware-software *object-capability security model*. These extensions support incrementally adoptable, high-performance, formally based, programmer-friendly underpinnings for fine-grained software decomposition and compartmentalization, motivated by and capable of enforcing the principle of least privilege.

The CHERI system architecture purposefully addresses known performance and robustness gaps in commodity ISAs that hinder the adoption of more secure programming models centered around the principle of least privilege. To this end, CHERI blends traditional paged virtual memory with a per-address-space capability model that includes capability registers, capability instructions, and tagged memory that have been added to the 64-bit MIPS ISA via a new *capability coprocessor*. CHERI also learns from the C-language fat-pointer literature: CHERI capabilities can describe not only regions of memory, but can also capture C pointer semantics allowing capabilities to be substituted for pointers in generated code.

CHERI's hybrid system approach, inspired by the Capsicum security model, allows incremental adoption of capability-oriented software design: software implementations that are more robust and resilient can be deployed where they are most needed, while leaving less critical software largely unmodified, but nevertheless suitably constrained to be incapable of having adverse effects. For example, we are focusing conversion efforts on low-level TCB components of the system: separation kernels, hypervisors, operating system kernels, language runtimes, and userspace TCBs such as web browsers. Likewise, we see early-use scenarios (such as data compression, protocol parsing, image processing, and video processing) that relate to particularly high-risk software libraries, which are concentrations of both complex and historically vulnerability-prone code combined with untrustworthy data sources, while leaving containing applications unchanged.

This report describes the CHERI *Instruction-Set Architecture* (ISA) and design, and provides reference documentation and potential memory models, along with their requirements.<sup>1</sup> It also briefly addresses the CHERI system hardware-software architecture, documenting our current thinking on integrating programming languages and operating systems with the CHERI hardware.

---

<sup>1</sup>We have attempted to avoid confusion among three rather different uses of the word 'architecture'. The ISA specifies the interface between hardware and software, rather than describing either the (micro-)architecture of a particular hardware prototype, or laying out the total-system hardware-software architecture. We expect further documentation on the latter to emerge in the final year of the project.

## Acknowledgments

The authors of this report thank other members of the CTSRD team, and our past and current research collaborators at SRI and Cambridge:

Ross J. Anderson	Ruslan Bukin	Gregory Chadwick	Nirav Dave
Paul J. Fox	Khilan Gudka	Jong Hun Han	Alex Horsman
Alexandre Joannou	Asif Khan	Myron King	Chris Kitching
Wojciech Koszek	Patrick Lincoln	Anil Madhavapeddy	Ilias Marinos
A. Theodore Marketos	Ed Maste	Andrew Moore	Will Morland
Alan Mujumdar	Prashanth Mundkur	Robert Norton	Philip Paeps
Colin Rothwell	John Rushby	Hassen Saidi	Hans Petter Selasky
Muhammad Shahbaz	Stacey Son	Richard Uhler	Munraj Vadera
Philip Withnall	Bjoern A. Zeeb		

The CHERI team also thanks past and current members of its external oversight group for significant support and contributions:

Lee Badger	Simon Cooper	Rance DeLong	Jeremy Epstein
Virgil Gligor	Li Gong	Mike Gordon	Steven Hand
Andrew Herbert	Warren A. Hunt Jr.	Doug Maughan	Greg Morrisett
Brian Randell	Kenneth F. Shottling	Joe Stoy	Tom Van Vleck
Samuel M. Weber			

We would also like to acknowledge the late David Wheeler, whose conversations with some of the authors about the CAP computer contributed to our thinking on CHERI.

Finally, we are grateful to Howie Shrobe, MIT professor and past DARPA CRASH program manager, who has offered both technical insight and support throughout this work. We are also grateful to Robert Laddaga, who has succeeded Howie in overseeing the CRASH program.

# Contents

<b>1</b>	<b>Introduction</b>	<b>8</b>
1.1	Motivation . . . . .	10
1.1.1	Trusted Computing Bases (TCBs) . . . . .	11
1.1.2	The Compartmentalization Problem . . . . .	12
1.2	The CHERI ISA Design . . . . .	13
1.2.1	A Hybrid Capability-System Architecture . . . . .	15
1.3	Threat Model . . . . .	16
1.4	Formal Methodology . . . . .	16
1.5	CHERI and CHERI2 Reference Prototypes . . . . .	17
1.6	Historical Context . . . . .	18
1.6.1	Capability Systems . . . . .	20
1.6.2	Microkernels . . . . .	20
1.6.3	Language and Runtime Approaches . . . . .	23
1.6.4	Bounds Checking and Fat Pointers . . . . .	24
1.6.5	Influences of Our Own Past Projects . . . . .	24
1.6.6	A Fresh Opportunity for Capabilities . . . . .	26
1.7	Publications . . . . .	26
1.8	Version History . . . . .	27
1.9	Document Structure . . . . .	30
<b>2</b>	<b>CHERI Capability-System Model</b>	<b>32</b>
2.1	Design Goals . . . . .	32
2.2	A Hybrid Capability-System Architecture . . . . .	34
2.3	The CHERI Software Stack . . . . .	36
2.4	Capability Model . . . . .	37
2.4.1	Capabilities are for Compilers . . . . .	37
2.4.2	Capabilities . . . . .	37
2.4.3	Capability Registers . . . . .	38
2.4.4	Memory Model . . . . .	39
2.4.5	Local Capabilities and Revocation . . . . .	40
2.4.6	Notions of Privilege . . . . .	41
2.4.7	Traps, Interrupts, and Exception Handling . . . . .	41
2.4.8	Tagged Memory . . . . .	42
2.4.9	Capability Instructions . . . . .	43
2.4.10	Object Capabilities . . . . .	43
2.4.11	Peripheral Devices . . . . .	44

<b>3</b>	<b>CHERI Instruction-Set Architecture</b>	<b>46</b>
3.1	Capability Registers . . . . .	46
3.2	Capabilities . . . . .	49
3.2.1	<b>tag</b> . . . . .	50
3.2.2	<b>sealed</b> . . . . .	50
3.2.3	<b>perms</b> . . . . .	50
3.2.4	<b>otype</b> . . . . .	51
3.2.5	<b>offset</b> . . . . .	51
3.2.6	<b>base</b> . . . . .	51
3.2.7	<b>length</b> . . . . .	51
3.3	Capability Permissions . . . . .	51
3.4	Capability Exceptions . . . . .	52
3.5	CPU Reset . . . . .	56
3.6	Changes to Standard MIPS ISA Processing . . . . .	56
3.7	Changes to the TLB . . . . .	57
3.8	Possible Extensions to the CHERI ISA . . . . .	58
<b>4</b>	<b>CHERI Instruction Reference</b>	<b>59</b>
4.1	Common Definitions . . . . .	59
4.2	Details of Individual Instructions . . . . .	61
	CAndPerm . . . . .	63
	CBTS . . . . .	64
	CBTU . . . . .	65
	CCall . . . . .	66
	CCheckPerm . . . . .	69
	CCheckType . . . . .	70
	CClearTag . . . . .	71
	CFromPtr . . . . .	72
	CGetBase . . . . .	74
	CGetCause . . . . .	75
	CGetLen . . . . .	76
	CGetOffset . . . . .	77
	CGetPCC . . . . .	78
	CGetPerm . . . . .	79
	CGetSealed . . . . .	80
	CGetTag . . . . .	81
	CGetType . . . . .	82
	CIncBase . . . . .	83
	CIncOffset . . . . .	85
	CJALR . . . . .	86
	CJR . . . . .	88
	CL[BHWD][U] . . . . .	90
	CLC . . . . .	93
	CLLD . . . . .	95
	CPtrCmp . . . . .	97
	CReturn . . . . .	100

CS[BHWD]	102
CSC	105
CSCD	107
CSeal	108
CSetCause	110
CSetLen	111
CSetOffset	113
CToPtr	114
CUnseal	115
4.3 Assembler Pseudo-Instructions	116
4.3.1 Capability Move	116
4.3.2 Get/Set Default Capability	116
4.3.3 Capability Loads and Stores of Floating-Point Values	117
<b>5 Design Rationale</b>	<b>118</b>
<b>6 Programming Languages and Operating Systems</b>	<b>127</b>
6.1 Development Plan and Status	127
6.2 Open-Source Foundations	128
6.3 Current Software Implementation	128
6.4 CheriBSD	128
6.5 Extended GNU Assembler (gas)	128
6.6 Cheri LLVM/Clang	128
6.7 Cheri Unit-Test Suite	129
6.8 Future Plans	130
<b>7 Future Directions</b>	<b>131</b>
7.1 An Open-Source Research Processor	132
7.2 Formal Methods for Bluespec	132
7.3 ABI and Compiler Development	133
7.4 Hardware Capability Support for FreeBSD	133
7.5 Evaluating Performance and Programmability	134

# Chapter 1

## Introduction

The Capability Hardware Enhanced RISC Instructions (CHERI) architecture extends the commodity 64-bit MIPS Instruction-Set Architecture (ISA) with new security primitives to allow software to efficiently implement *fine-grained memory protection* and an *object-capability security model*. CHERI's extensions are intended to support incrementally adoptable, high-performance, formally supported, and programmer-friendly underpinnings for robust and scalable *software compartmentalization* motivated by *the principle of least privilege*. CHERI is a *hybrid capability-system architecture* that combines design elements from both conventional RISC ISAs and capability systems. We hope that this will support gradual deployment of CHERI features in existing software, offering a more gentle hardware-software adoption path. CHERI has four central design goals aimed at dramatically improving the security of contemporary Trusted Computing Bases (TCBs) through processor support for fine-grained memory protection and scalable software compartmentalization, whose (at times) conflicting requirements have required careful negotiation in our design:

1. *Fine-grained memory protection* improves software resilience to escalation paths that allow software bugs to be coerced into more powerful software vulnerabilities; e.g., through remote code injection via buffer overflows and other memory-based techniques. Unlike widely deployed approaches, CHERI's memory protection is intended to be driven by the compiler in protecting programmer-described data structures and references, rather than via coarse page-granularity protections. Fine-grained protection also provides the foundation for expressing compartmentalization within application instances. We draw on, and extend, ideas from recent work in C-language *software bounds checking* by combining *fat pointers* with capabilities, allowing capabilities to be substituted for C pointers with only limited changes to program semantics.
2. *Software compartmentalization* involves the decomposition of software into isolated components to mitigate the effects of security vulnerabilities by applying sound principles of security, such as abstraction, encapsulation, and especially least privilege. Previously, it seems that the adoption of compartmentalization has been limited by a conflation of hardware primitives for virtual addressing and separation, leading to inherent performance and programmability problems when implementing fine-grained separation. Specifically, we seek to decouple the virtualization from separation to avoid scalability problems imposed by translation look-aside buffer (TLB)-based Memory Management Units (MMUs), which impose a very high performance penalty as the number of protection



domains increases, as well as complicating the writing of compartmentalized software.

3. Simultaneously, we require a realistic *technology transition path* that is applicable to current software and hardware designs. CHERI hardware must be able to run most current software without significant modification, and allow incremental deployment of security improvements starting with the most critical software components: the TCB foundations on which the remainder of the system rests, and software with the greatest exposure to risk. CHERI features should significantly improve security so that vendors of mobile and embedded devices would seek its feature set from CPU companies (such as MIPS and ARM); these CHERI features must at the same time conform to vendor expectations for performance, power use, and compatibility to compete with less secure alternatives.
4. Finally, we wish to draw on *formal methodologies* wherever feasible, to improve our confidence in the design and implementation of CHERI. This use is necessarily subject to real-world constraints of timeline, budget, design process, and prototyping, but will help ensure that we avoid creating a system that cannot meet our functional and security requirements. Formal methods can also help to avoid many of the characteristic design flaws that are common in both hardware and software. This desire requires us not only to perform research into CPU and software design, but also to develop new formal methodologies and adaptations and extensions of existing ones.

We are concerned with trustworthy systems and networks, where *trustworthiness* is a multidimensional measure of how well a system or other entity satisfies its various requirements – such as those for security, system integrity, and reliability, as well as survivability, robustness, and resilience, notably in the presence of a wide range of adversities such as hardware failures, software flaws, malware, accidental and intentional misuse, and so on. Our approach to trustworthiness encompasses hardware and software architecture, dynamic and static evaluation, formal and nonformal analyses, good software engineering practices, and much more.

Our selection of RISC as a foundation for the CHERI capability extensions is motivated by two factors. First, simple instruction set architectures are easier to reason about, extend, and implement. Second, RISC architectures (such as ARM and MIPS) are widely used in network embedded and mobile device systems such as firewalls, routers, smart phones, and tablets – markets with the perceived flexibility to adopt new CPU facilities, and also an immediate and pressing need for improved security. CHERI’s new security primitives would also be useful in workstation and server environments, which face similar security challenges.

In its current incarnation, we have prototyped CHERI as a *coprocessor* extension to the 64-bit MIPS ISA, but our approach is intended to easily support other similar ISAs, such as 64-bit ARM. The design principles would also apply to other non-RISC ISAs, such as 32-bit and 64-bit Intel and AMD, but require significantly more adaptation work, as well as careful consideration of the implications of the diverse set of CPU features found in more CISC-like architectures. We also consider the possibility that the syntax and semantics of the CHERI model might be implemented over conventional CPUs with the help of the compiler, static checking, and dynamic enforcement approaches found in software fault isolation techniques [74] or Google Native Client (NaCl) [90]. All of these future considerations are considerably enhanced by our use of the Bluespec SystemVerilog [9] Hardware Description Language (HDL) in prototyping, which has provided significant opportunities for rapid redesigns through its use of modular abstraction, encapsulation, and hierarchicalization.

## 1.1 Motivation

The CHERI ISA provides a sound and formally based processor foundation for the principled development of highly trustworthy systems. The CHERI instruction-set architecture and system architecture build on and extend decades of research into hardware and operating-system security.<sup>1</sup> However, some of the historic approaches that CHERI incorporates (especially capability architectures) have not been adopted in commodity hardware designs. In light of these past transition failures, a reasonable question is “Why now?” What has changed that would allow CHERI to succeed where so many previous efforts have failed? Several factors have motivated our decision to begin and carry out this project:

- Dramatic changes in threat models, resulting from ubiquitous connectivity and pervasive uses of computer technology in many diverse and widely used applications such as wireless mobile devices, automobiles, and critical infrastructure.
- New opportunities for research into (and possible revisions of) hardware-software interfaces, brought about by programmable hardware (especially FPGA soft cores) and complete open-source software stacks.
- An increasing trend towards exposing inherent hardware parallelism through virtual machines and explicit software multi-programming, and an increasing awareness of information flow for reasons of power and performance that may align well with the requirements of security.
- Emerging advances in programming languages, such as the ability to map language structures into protection parameters to more easily express and implement various policies.
- Reaching the tail end of a “compatibility at all costs” trend in CPU design, due to proximity to physical limits on clock rates and trends towards heterogeneous and distributed computing. While Wintel remains entrenched on the desktop, mobile systems – such as phones and tablet PCs, as well as appliances and embedded devices – are much more diverse, running on a wide variety of instruction set architectures (especially ARM and MIPS).
- Likewise, new diversity in operating systems has arisen, in which commercial products such as Apple’s iOS and Google’s Android extend open-source systems such as FreeBSD and Linux. These new platforms abandon many traditional constraints, requiring that rewritten applications conform to new security models, programming languages, hardware architectures, and user-input modalities.
- Development of *hybrid capability-system models*, such as Capsicum [81], that integrate capability-system design tenets into current operating-system and language designs. With CHERI, we are transposing this design philosophy into the instruction-set architecture. Hybrid design is a key differentiator from prior capability-system processor designs that have typically required ground-up software-architecture redesign and reimplementation.

---

<sup>1</sup>Levy’s *Capability-Based Computer Systems* provides a detailed history of segment- and capability-based designs through the early 1990s [40]. However, it leaves off just as the transition to microkernel-based capability systems such as Mach [1] and L4 [41], as well as capability-influenced virtual machines such as the Java Virtual Machine [23], begins. Section 1.6 discuss historical influences on this work in greater detail.

- Significant changes in the combination of hardware, software, and formal methods to enhance assurance (such as those noted above) now make possible the development of trustworthy system architectures that previously were simply too far ahead of their times.

In the following sections, we consider the context and motivation for CHERI, a high-level view of the CHERI design, the role of formal methods in the project, and our work-in-progress research prototype.

### 1.1.1 Trusted Computing Bases (TCBs)

Contemporary client-server and cloud computing are premised on highly distributed applications, with end-user components executing in rich execution substrates such as POSIX applications on UNIX, or AJAX in web browsers. However, even thin clients are not thin in most practical senses: as with client-server computer systems, they are built from commodity operating system kernels, hundreds of user-space libraries, window servers, language runtime environments, and web browsers, which themselves include scripting language interpreters, virtual machines, and rendering engines. Both server and embedded systems likewise depend on complex (and quite similar) software stacks. All require confluence of competing interests, representing multiple sites, tasks, and end users in unified computing environments.

Whereas higher-layer applications are able to run on top of type-safe or constrained execution environments, such as JavaScript interpreters, lower layers of the system must provide the link to actual execution on hardware. As a result, almost all such systems are written in the C programming language; collectively, this Trusted Computing Base (TCB) consists of many tens of millions of lines of trusted (but not trustworthy) C and C++ code. Coarse hardware, OS, and language security models mean that much of this code is security-sensitive: a single flaw, such as an errant NULL pointer dereference in the kernel, can expose all rights held by users of a system to an attacker or to malware.

The consequences of compromise are serious, and include loss of data, release of personal or confidential information, damage to system and data integrity, and even total subversion of a user's online presence and experience by the attacker (or even accidentally without any attacker presence!). These problems are compounded by the observation that the end-user systems are also an epicenter for multi-party security composition, where a single web browser or office suite (which manages state, user interface, and code execution for countless different security domains) must simultaneously provide strong isolation and appropriate sharing. The results present not only significant risks of compromise that lead to financial loss or disruption of critical infrastructure, but also frequent occurrences of such events.

Software vulnerabilities appear inevitable: even as the execution substrates improve in their ability to resist attacks such as buffer overflows and integer vulnerabilities, logical errors will necessarily persist. Past research has shown that compartmentalizing applications into components executed in isolated sandboxes can mitigate exploited vulnerabilities (sometimes referred to as privilege separation). Only the rights held by a compromised component are accessible to a successful attacker. This technique is effectively applied in Google's Chromium web browser, placing HTML rendering and JavaScript interpretation into sandboxes isolated from the global file system. This technique exploits the principle of least privilege: if each software element executes with only the rights required to perform its task, then attackers lose access to most all-or-nothing toeholds; vulnerabilities may be significantly or entirely mitigated, and attackers must identify many more vulnerabilities to accomplish their goals.

## 1.1.2 The Compartmentalization Problem

The *compartmentalization problem* arises from attempts to decompose security-critical applications into components running in different security domains: the practical application of the principle of least privilege to software. Historically, compartmentalization of TCB components such as operating system kernels and central system services has caused significant difficulty for software developers – which limits its applicability for large-scale, real-world applications, and leads to the abandonment of promising research such as 1990s *microkernel* projects. A recent resurgence of compartmentalization, applied in userspace to applications such as OpenSSH [60] and Chromium [62], and most recently in our own Capsicum project [81], has been motivated by a critical security need; however it has seen success only at very coarse separation granularity due to the challenges involved.

On current conventional hardware, native applications must be converted to employ message passing between address spaces (or processes) rather than using a unified address space for communication, sacrificing programmability and performance by transforming a local programming problem into a distributed systems problem. As a result, large-scale compartmentalized applications are difficult to design, write, debug, maintain, and extend; this raises serious questions about correctness, performance, and most critically, security.

These problems occur because current hardware provides strong separation only at coarse granularity via rings and virtual address spaces, making the isolation of complete applications (or even multiple operating systems) a simple task, but complicates efficient and easily expressed separation between tightly coupled software components. Three closely related problems arise:

**Performance is sacrificed.** Creating and switching between security domains is expensive due to reliance on software and hardware address-space infrastructure, such as a quickly overflowed Translation Look-aside Buffer (TLB) that can lead to massive performance degradation. Also, above an extremely low threshold, performance overhead from context switching between security domains tends to go from extremely expensive to intolerable: each TLB entry is an access-control list, with each object (page) requiring multiple TLB entries, one for each authorized security domain.

High-end server CPUs typically have TLB entries in the low hundreds, and even recent network embedded devices reach the low thousands; the TLB footprint of fine-grained, compartmentalized software increases with the product of in-flight security domains and objects due to TLB aliasing, which may easily require tens or hundreds of thousands of spheres of protection. The transition to CPU multi-threading has not only failed to relieve this burden, but actively made it worse: TLBs are implemented using ternary content-addressable memory (TCAMs) or other expensive hardware lookup functions, and are often shared between hardware threads in a single core due to their expense.

In comparison, physically indexed general-purpose CPU caches are several orders of magnitude larger than TLBs, scaling instead with the working set of code paths explored or the memory footprint of data actively being used. If the same data is accessed by multiple security domains, it shares data or code cache (but not TLB entries) with current CPU designs.

**Programmability is sacrificed.** Within a single address space, programmers can easily and efficiently share memory between application elements using pointers from a common names-

pace. The move to multiple processes frequently requires the adoption of a distributed programming model based on explicit message passing, making development, debugging, and testing more difficult. RPC systems and higher-level languages are able to mask some (although usually not all) of these limitations, but are poorly suited for use in TCBs – RPC systems and programming language runtimes are non-trivial, security-critical, and implemented using weaker lower-level facilities.<sup>2</sup>

**Security is sacrificed.** Current hardware is intended to provide robust shared memory communication only between mutually trusting parties, or at significant additional expense; granularity of delegation is limited and its primitives expensive, leading to programmer error and extremely limited use of granular separation. Poor programmability contributes directly to poor security properties.

## 1.2 The CHERI ISA Design

The CHERI ISA embodies two fundamental and closely linked technical goals to address vulnerability mitigation: first, fine-grained capability-oriented memory protection within address spaces, and second, primitives to support both scalable and programmer-friendly compartmentalization within address spaces based on the object-capability model. The CHERI model is designed to support low-level TCBs, typically implemented in C or a C-like language, in workstations, servers, mobile devices, and embedded devices. Simultaneously, it will provide reasonable assurance of correctness and a realistic technology transition path from existing hardware and software platforms.

To this end, we have prototyped CHERI as an Instruction-Set Architecture (ISA) extension to the widely used 64-bit MIPS ISA; we are also considering the implications for the RISC-V and ARM ISAs. CHERI adds the following features to a RISC ISA via a new *capability coprocessor*, or instruction-set extension, that supports granular memory protection within address spaces:

- The *capability register file* describes the rights (*protection domain*) of the executing thread to memory that it can access, and to object references that can be invoked to transition between protection domains. Capability registers supplement the general-purpose register file, allowing capabilities to displace general-purpose registers in describing data and object references. Capability registers contain a tag, sealed bit, permission mask, base, length, and offset (allowing the description of not just bounded regions, but also a pointer into that region, improving C-language compatibility). Capability registers are suitable for describing both data and code, and can hence protect both data integrity/confidentiality and control flow. Certain registers are reserved for use in exception handling; all others are available to be managed by the compiler using the same techniques used

---

<sup>2</sup>Through extreme discipline, a programming model can be constructed that maintains synchronized mappings of multiple address spaces, while granting different rights on memory between different processes. This leads to even greater TLB pressure and expensive context switch operations, as the layouts of address spaces must be managed using cross-address-space communication. Bittau has implemented this model via *threads*, an OS primitive that tightly couples UNIX processes via shared memory associated with data types – a promising separation approach constrained by the realities of current CPU design [8].

with conventional registers. A similar approach is also possible, in which capability registers extend existing general-purpose registers: such a design might reduce the hardware resources required to implement CHERI support, but could also make it more difficult to provide a superset Application Binary Interface (ABI) that avoids the need for explicit ABI wrappers when passing control between native and capability-enhanced code.

- The *capability instructions* allow executing code to create, constrain (e.g., by increasing the base, decreasing the length, or reducing permissions), manage, and inspect capability register values. Both unsealed (memory) and sealed (object) capabilities can be loaded and stored via memory capability registers (i.e., dereferencing); object capabilities can be invoked, via special instructions, allowing a transition between protection domains, but their fields are given additional integrity protections to provide encapsulation. Capability instructions implement *guarded manipulation*: invalid capability manipulations (e.g., to increase rights or length) and invalid capability dereferences (e.g., to access outside of a bounds-checked region) result in an exception that can be handled by the supervisor or language runtime. Most capability instructions are part of the usermode ISA, rather than privileged ISA, and will be generated by the compiler to describe application data structures and protection properties.
- *Tagged memory* associates a 1-bit tag with each capability-aligned and capability-sized word in physical memory, which allows capabilities to be safely loaded and stored in memory without loss of integrity. This functionality expands a thread's effective protection domain to include the transitive closure of capability values that can be loaded via capabilities via those present in its register file. For example, a capability register representing a C pointer to a data structure can be used to load further capabilities from that structure, referring to further data structures, which could not be accessed without suitable capabilities. Writes to capability values in memory that do not originate from a valid capability in the capability-register file will clear the tag bit associated with that memory, preventing accidental (or malicious) dereferencing of invalid capabilities.

In keeping with the RISC philosophy, CHERI instructions are intended for use primarily by the compiler rather than directly by the programmer, and consist of relatively simple instructions that avoid, for example, combining memory access and register value manipulation in a single instruction. In our current software prototypes, there are direct mappings from programmer-visible, C-language pointers to capabilities in much the same way that conventional code generation translates pointers into general-purpose register values; this allows CHERI to continuously enforce bounds checking, pointer integrity, and so on. There is likewise a strong synergy between the capability-system model, which espouses a separation of policy and mechanism, and RISC: CHERI's features make possible the implementation of a wide variety of OS, compiler, and application-originated policies on a common protection substrate that optimizes fast paths through hardware support.

The capability coprocessor is a coprocessor in two senses. First, the capability coprocessor occupies a portion of the existing ISA encoding dedicated to extensions (typically referred to as coprocessor instructions). Second, the capability coprocessor supplements the general-purpose register file with its own ISA-managed registers, as well as performing (and transforming) memory access, and delivering exceptions to the main pipeline, requiring hardware resources that interact with the primary processor pipeline. This behavior is comparable in many ways to

system, floating-point, vector, or cryptographic coprocessors that will similarly supplement the base ISA and processor features.

Wherever possible, CHERI systems make use of existing hardware designs: processor pipelines and register files, cache memory, system buses, commodity DRAM, and commodity peripheral devices such as NICs and display cards. We are currently focusing on enforcement of CHERI security properties on applications running on a general-purpose processor; in future work, we hope to consider the effects of implementing CHERI in peripheral processors, such as those found in Network Interface Cards (NICs) or Graphical Processing Units (GPUs).

In order to prototype this approach, we have localized our ideas about CHERI capability access to a specific instruction set: the 64-bit MIPS ISA. This has necessarily led to a set of congruent implementation decisions about register-file size, selection of specific instructions, exception handling, memory alignment requirements, and so on, that reflect that starting-point ISA. These decisions might be made differently with another starting-point ISA as they are simply surface features of an underlying approach; we anticipate that adaptations to ISAs such as ARM and RISC-V would adopt instruction-encoding conventions, and so on, more in keeping with their specific flavor and approach.

Other design decisions reflect the goal of creating a platform for prototyping and exploring the design space itself; among other choices, this includes the selection of 256-bit capabilities, which have given us greater flexibility to experiment with various bounds-checking and capability behaviors. Reducing capabilities to 128-bit is not unreasonable, and measurements suggest that cache footprint increases from capabilities can be significantly mitigated through such a change. However, this would also introduce tradeoffs in the granularity of memory (e.g., losing access to the full 64-bit space) and flexibility of the object-capability design (e.g., loss of software-defined permission bits) that would need to be reasoned about carefully.

We believe, however, that the higher-level memory protection and security models we have described would relatively easily apply to variations in ISA-level implementation. This should allow reasonable source-level software portability (leaving aside OS assembly code and compiler code generation) across the CHERI model implemented in different architectures, in much the same way that conventional OS and application C code is moderately portable across underlying ISAs.

### **1.2.1 A Hybrid Capability-System Architecture**

Unlike past research into capability systems, CHERI allows traditional address-space separation, implemented using a memory management unit (MMU), to coexist with granular decomposition of software within each address space. As a result, fine-grained memory protection and compartmentalization can be applied selectively throughout existing software stacks to provide an incremental software migration path. We envision early deployment of CHERI extensions in selected components of the TCB's software stack: separation kernels, operating system kernels, programming language runtimes, sensitive libraries such as those involved in data compression or encryption, and network applications such as web browsers and web servers.

CHERI addresses current limitations on memory protection and compartmentalization by extending virtual memory-based separation with hardware-enforced, fine-grained protection within address spaces. Granular memory protection mitigates a broad range of previously exploitable bugs by coercing common memory-related failures into exceptions that can be handled by the application or operating system, rather than yielding control to the attacker. The

CHERI approach also restores a single address-space programming model for compartmentalized (sandboxed) software, facilitating efficient, programmable, and robust separation through the capability model.

We have selected this specific composition of traditional virtual memory with an in-address-space security model to facilitate technology transition: in CHERI, existing C-based software can continue to run within processes, and even integrate with capability-enhanced software within a single process, to provide improved robustness for selected software components – and perhaps over time, all software components. For example, a sensitive library (perhaps used for image processing) might employ capability features while executing as part of a CHERI-unaware web browser. Likewise, a CHERI-enabled application can sandbox and instantiate multiple copies of unmodified libraries, to efficiently and easily gate access to the rest of application memory of the host execution environment.

### 1.3 Threat Model

CHERI protections constrain code “in execution” and allow fine-grained management of privilege within a framework for controlled separation and communication. Code in execution can represent the focus of many potentially malicious parties: subversion of legitimate code in violation of security policies, injection of malicious code via back doors, trojan horses, and malware, and also denial-of-service attacks. CHERI’s fine-grained memory protection mitigates many common attack techniques by reducing opportunities for the conflation of code and data, corruption of control flow, and also catches many common exploitable programmer bugs; compartmentalization constrains successful attacks via pervasive observance of the principle of least privilege.

Physical attacks on CHERI-based systems are explicitly excluded from our threat model, although CHERI CPUs might easily be used in the context of tamper-evident or tamper-resistant systems. Similarly, no special steps have been taken in our design to counter undesired leakage of electromagnetic emanations and certain other side channels such as acoustic inferences: we take for granted the presence of an electronic foundation on which CHERI can run. CHERI will provide a supportive framework for a broad variety of security-sensitive activities; while not itself a distributed system, CHERI could form a sound foundation for various forms of distributed trustworthiness.

Somewhat to our chagrin, we report that the CHERI design currently includes no features for resisting covert or side-channel attacks: these have proven increasingly relevant in CPU design, but the tools CHERI provides do not improve resilience against these attacks. In some sense, they increase exposure: the greater the offers of protection within a system, the greater the potential impact of unauthorized communication channels. As such, we hope side-channel attacks are a topic that we will be able to explore in future work. Overall, we believe that our threat model is realistic and will lead to systems that can be substantially more trustworthy than today’s commodity systems.

### 1.4 Formal Methodology

Throughout this project, we apply formal methodology to help avoid system vulnerabilities. An important early observation is that existing formal methodology applied to software se-



curity has significant problems with multi-address-space security models; formal approaches have relied on the usefulness of addresses (pointers) as unique names for objects. Whereas this weakness in formal methods is a significant problem for traditional CPU designs, which offer security primarily through rings and address-space translation, CHERI’s capability model is scoped within address spaces. This offers the possibility of applying existing software proof methodology in the context of hardware isolation (and other related properties) in a manner that was previously infeasible. We are more concretely (and judiciously) applying formal methodology in three areas:

1. We have developed a formal semantics for the CHERI ISA described in SRI’s Prototype Verification System (PVS) – an automated theorem-proving and model-checking toolchain – which can be used to verify the expressibility of the ISA, but also to prove properties of critical code. For example, we are interested in proving the correctness of software-based address-space management and domain transitions. We are likewise able to automatically generate ISA-level test suites from formal descriptions of instructions, which are applied directly to our hardware implementation.
2. We have also developed a more complete ISA model incorporating both MIPS and CHERI instructions using Cambridge’s L3 instruction-set description language. Although we have not yet used this for automated theorem proving, we increasingly use the L3 description as a “gold model” of the instruction set against which our test suite is validated, software implementations can be tested in order to generate traces of correct processor execution, and so on. We have used the L3 model to identify a number of bugs in multiple hardware implementations of CHERI, as well as to discover software dependences on undefined instruction-set behavior.
3. We have developed extensions to the Bluespec compiler to export an HDL description to SRI’s PVS and SAL model checker. We have also developed a new tool (Smten) for efficient SMT (Satisfiability Modulo Theories) modeling of designs (using SRI’s Yices), and another tool for automated extraction of key properties from larger Bluespec designs, both of which greatly simplify formal analysis. These tools will allow us to verify low-level properties of the hardware design and use the power of model checking and satisfiability solvers to analyze related properties. Ideally they will also help link ISA-level specifications with the CPU implementation.

A detailed description of formal methods efforts relating to CHERI may be found in the *CHERI Formal Methods Report*.

## 1.5 CHERI and CHERI2 Reference Prototypes

As a central part of this research, we are developing reference prototypes of the CHERI ISA via several CHERI processor designs. These prototypes allow us to explore, validate, evaluate, and demonstrate the CHERI approach through realistic hardware properties and real-world software stacks. A detailed description of the current prototypes, both from architectural and practical use perspectives, may be found in the companion *BERI Hardware Reference*, *BERI Software Reference*, and *CHERI User’s Guide* documents.

Our first prototype, known simply as CHERI1, is based on Cambridge’s MAMBA research processor, and is a single-threaded, single-core implementation intended to allow us to explore ISA tradeoffs. This prototype is implemented in the Bluespec HDL, a high-level functional programming language for hardware design. CHERI1 is a pipelined baseline processor implementing the 64-bit MIPS ISA, and incorporates an initial prototype of the CHERI capability coprocessor that includes capability registers and a basic capability instruction set.

We have ported the commodity open-source FreeBSD operating system, with support for a wide variety of peripherals on the Terasic tPad and DE4 FPGA development boards; we use these boards in both mobile tablet-style and network configurations. FreeBSD is able to manage the capability coprocessor and maintain additional thread state for capability-aware user applications, although capability features are not yet used within the kernel for its own internal protection. FreeBSD also implements exception-handler support for object-capability invocation, signal delivery when protection faults occur (allowing language runtimes to catch and handle protection violations), and error recovery for in-process sandboxes. We have adapted the Clang and LLVM compiler suite to allow language-level annotations in C to direct capability use. In addition, we have developed a number of capability-enhanced applications demonstrating fine-grained memory protection and in-process compartmentalization to explore security, performance, and programmability tradeoffs. We also have a work-in-progress multi-core version of the CHERI prototype.

Using Bluespec, we are able to run the CPU in simulation, and synthesize the CHERI design to execute in field-programmable gate arrays (FPGAs). In our development work, we are targeting an Altera FPGAs on Terasic development boards. However, in our companion MRC project we have also targeted CHERI at the second-generation NetFPGA 10G research and teaching board, which we hope to use in ongoing research into datacenter network fabrics. That work includes the development of Blueswitch, a Bluespec implementation of an OpenFlow switch that can operate as a tightly coupled CHERI coprocessor. In the future, should it become desirable, we will be able to construct an ASIC design from the same Bluespec source code. We have released the CHERI soft core as *open-source hardware*, making it available for more widespread use in research. This should allow others, especially in the research community, to reproduce and extend our results.

We have also developed a second prototype, known as CHERI2, which is compatible with CHERI but has additional CPU features including fine-grained multi-threading. This will be used as a platform for exploring the synergy between compartmentalization and parallelism: each thread exists in its own security domain so efficient transitions can be achieved by communicating between threads. CHERI2 also employs a more stylized form of Bluespec that is intended to considerably enhance our formal analysis of the hardware architecture.

## 1.6 Historical Context

As with many aspects of contemporary computer and operating system design, the origins of operating system security may be found at the world’s leading research universities, but especially the Massachusetts Institute of Technology (MIT), the University of Cambridge, and Carnegie Mellon University. MIT’s Project MAC, which began with MIT’s Compatible Time Sharing System (CTSS) [13], and continued over the next decade with MIT’s Multics project, described many central tenets of computer security [14, 26]. Dennis and Van Horn’s 1965 *Pro-*

*programming Semantics for Multiprogrammed Computations* [17] laid out principled hardware and software approaches to concurrency, object naming, and security for multi-programmed computer systems – or, as they are known today, multi-tasking and multi-user computer systems. Multics implemented a coherent, unified architecture for processes, virtual memory, and protection, integrating new ideas such as *capabilities*, unforgeable tokens of authority, and *principals*, the end users with whom authentication takes place and to whom resources are accounted [67].

In 1975, Saltzer and Schroeder surveyed the rapidly expanding vocabulary of computer security in *The Protection of Information in Computer Systems* [68]. They enumerated design principles such as the *principle of least privilege* (which demands that computations run with only the privileges they require) and the core security goals of protecting *confidentiality*, *integrity*, and *availability*. The tension between fault tolerance and security (a recurring debate in systems literature) saw its initial analysis in Lampson’s 1974 *Redundancy and Robustness in Memory Protection* [35], which considered ways in which hardware memory protection addressed accidental and intentional types of failure: if it is not reliable, it will not be secure, and if it is not secure, it will not be reliable! Intriguingly, recent work by Nancy Leveson and William Young has unified security and human safety as overarching emergent system properties [39], and allows the threat model to fall out of the top-down analysis, rather than driving it. This work in some sense unifies a long thread of work that considers trustworthiness as a property encompassing security, integrity, reliability, survivability, human safety, and so on (e.g., [54, 55], among others).

The Security Research community also blossomed outside of MIT: Wulf’s Hydra operating system at Carnegie Mellon University (CMU) [88, 12], Needham and Wilkes’ CAP Computer at Cambridge [86], SRI’s Provably Secure Operating System (PSOS) [21, 55] hardware-software co-design that included strongly typed object capabilities, Rushby’s security kernels supported by formal methods at Newcastle [66], and Lampson’s work on formal models of security protection at the Berkeley Computer Corporation all explored the structure of operating system access control, and especially the application of capabilities to the protection problem [36, 37]. Another critical offshoot from the Multics project was Ritchie and Thompson’s UNIX operating system at Bell Labs, which simplified concepts from Multics, and became the basis for countless directly and indirectly derived products such as today’s Solaris, FreeBSD, Mac OS X, and Linux operating systems [64].

The creation of secure software went hand in hand with analysis of security flaws: Anderson’s 1972 US Air Force *Computer Security Technology Planning Study* not only defined new security structures, such as the *reference monitor*, but also analyzed potential attack methodologies such as Trojan horses and inference attacks [3]. Karger and Schell’s 1974 report on a security analysis of the Multics system similarly demonstrated a variety of attacks that bypass hardware and OS protection [31]. In 1978, Bisbey and Hollingworth’s *Protection Analysis: Project final report* at ISI identified common patterns of security vulnerability in operating system design, such as race conditions and incorrectly validated arguments at security boundaries [7]. Adversarial analysis of system security remains as critical to the success of security research as principled engineering and formal methods.

Almost fifty years of research have explored these and other concepts in great detail, bringing new contributions in hardware, software, language design, and formal methods, as well as networking and cryptography technologies that transform the context of operating system security. However, the themes identified in those early years remain topical and highly influential,

structuring current thinking about systems design.

Over the next few sections, we consider three closely related ideas that directly influence our thinking for CTSRD: capability security, microkernel OS design, and language-based constraints. These apparently disparate areas of research are linked by a duality, observed by Morris in 1973, between the enforcement of data types and safety goals in programming languages on one hand, and the hardware and software protection techniques explored in operating systems [49] on the other hand. Each of these approaches blends a combination of limits defined by static analysis (perhaps at compile-time), limits on expression on the execution substrate (such as what programming constructs can even be represented), and dynamically enforced policy that generates runtime exceptions (often driven by the need for configurable policy and labeling not known until the moment of access). Different systems make different uses of these techniques, affecting expressibility, performance, and assurance.

### 1.6.1 Capability Systems

Throughout the 1970s and 1980s, high-assurance systems were expected to employ a capability-oriented design that would map program structure and security policy into hardware enforcement; for example, Lampson’s BCC design exploited this linkage to approximate least privilege [36, 37].

Systems such as the CAP Computer at Cambridge [86] and Ackerman’s DEC PDP-1 architecture at MIT [2] attempted to realize this vision through embedding notions of capabilities in the memory management unit of the CPU, an approach described by Dennis and Van Horn [17]. Levy provides a detailed exploration of segment- and capability-oriented computer system design through the mid-1980s in *Capability-Based Computer Systems* [40].

### 1.6.2 Microkernels

Denning has argued that the failures of capability hardware projects were classic failures of large systems projects, an underestimation of the complexity and cost of reworking an entire system design, rather than fundamental failures of the capability model [16]. However, the benefit of hindsight suggests that the earlier demise of hardware capability systems was a result of three related developments in systems research: microkernel OS design, a related interest from the security research community in security kernel design, and Patterson and Sequin’s Reduced Instruction-Set Computers (RISC) [59].

However, with a transition from complex instruction set computers (CISC) to reduced instruction set computers (RISC), and a shift away from microcode toward operating system implementation of complex CPU functionality, the attention of security researchers turned to microkernels.

Carnegie Mellon’s Hydra [12, 89] embodied this approach, in which microkernel message passing between separate tasks stood in for hardware-assisted security domain crossings at capability invocation. Hydra developed a number of ideas, including the relationship between capabilities and object references, refined the *object-capability* paradigm, and further pursued the separation of policy and mechanism.<sup>3</sup> Jones and Wulf argue through the Hydra design that the capability model allows the representation of a broad range of system policies as a result

---

<sup>3</sup>Miller has expanded on the object-capability philosophy in considerable depth in his 2006 PhD dissertation, *Robust composition: towards a unified approach to access control and concurrency control* [47]

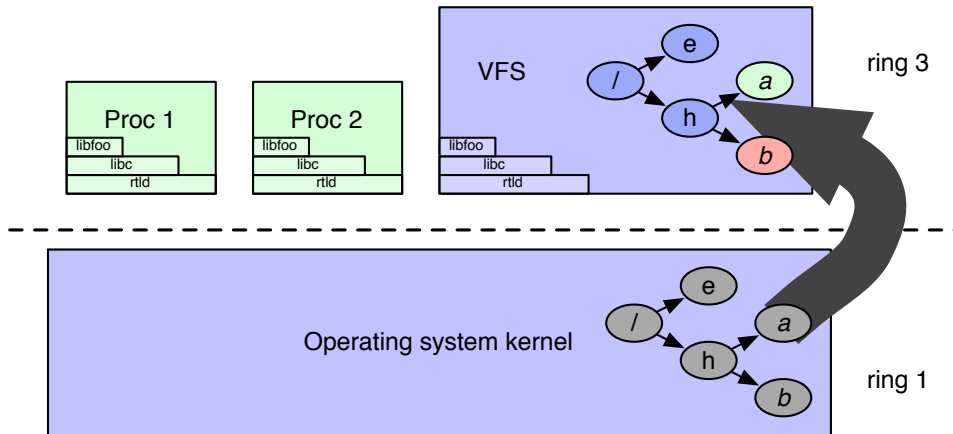


Figure 1.1: The microkernel movement shifted complex OS components, such as file systems, from the kernel to userspace tasks linked by IPC. Microkernels provide a smaller, easier-to-analyze, easier-to-debug, and more robust foundation in the face of dramatic increases in OS complexity.

of integration with the OS object model, which in turn facilitates interposition as a means of imposing policies on object access [29].

Successors to Hydra at CMU include Accent and Mach [61, 1], both microkernel systems intended to explore the decomposition of a large and decidedly un-robust operating system kernel. Figure 1.1 illustrates the principle of microkernel design: traditional OS services, such as the file system, are migrated out of ring 0 and into user processes, improving debuggability and independence of failure modes. They are also based on mapping of capabilities as object references into IPC pipes (*ports*), in which messages on ports represent methods on objects. This shift in operating system design went hand in hand with a related analysis in the security community: Lampson’s model for capability security was, in fact, based on pure message passing between isolated processes [37]. This further aligned with proposals by Andrews [4] and Rushby [66] for a *security kernel*, whose responsibility lies solely in maintaining isolation, rather than the provision of higher-level services such as file systems. Unfortunately, the shift to message passing also invalidated Fabry’s semantic argument for capability systems, namely, that by offering a single namespace shared by all protection domains, the distributed system programming problem could be avoided [20].

A panel at the 1974 National Computer Conference and Exposition (AFIPS) chaired by Lipner brought the design goals and choices for microkernels and security kernels clearly into focus: microkernel developers sought to provide flexible platforms for OS research with an eye towards protection, while security kernel developers aimed for a high assurance platform for separation, supported by hardware, software, and formal methods [42].

The notion that the microkernel, rather than the hardware, is responsible for implementing the protection semantics of capabilities also aligned well with the simultaneous research (and successful technology transfer) of RISC designs, which eschewed microcode by shifting complexity to the compiler and operating system. Without microcode, the complex C-list peregrinations of CAP’s capability unit, and protection domain transitions found in other capability-based systems, become less feasible in hardware. Virtual memory designs based on fixed-size pages and simple semantics have since been standardized throughout the industry.

Security kernel designs, which combine a minimal kernel focused entirely on correctly implementing protection, and rigorous application of formal methods, formed the foundation for several secure OS projects during the 1970s. Schiller's security kernel for the PDP-11/45 [69] and Neumann's Provably Secure Operating System [22] design study were ground-up operating system designs based soundly in formal methodology.<sup>4</sup> In contrast, Schroeder's MLS kernel design for Multics [70], the DoD Kernelized Secure Operating System (KSOS) [44], and Bruce Walker's UCLA UNIX Security Kernel [75] attempted to slide MLS kernels underneath existing Multics and UNIX system designs. Steve Walker's 1980 survey of the state of the art in trusted operating systems provides a summary of the goals and designs of these high-assurance security kernel designs [76].

The advent of CMU's Mach microkernel triggered a wave of new research into security kernels. TIS's Trusted Mach (TMach) project extended Mach to include mandatory access control, relying on enforcement in the microkernel and a small number of security-related servers to implement the TCB to accomplish sufficient assurance for a TCSEC B3 evaluation [10]. Secure Computing Corporation (SCC) and the National Security Agency (NSA) adapted PSOS's type enforcement from LoCK (LOgical Coprocessor Kernel) for use in a new Distributed Trusted Mach (DTMach) prototype, which built on the TMach approach while adding new flexibility [71]. DTMach, adopting ideas from Hydra, separates mechanism (in the microkernel) from policy (implemented in a userspace security server) via a new reference monitor framework, FLASK [73]. A significant focus of the FLASK work was performance: an access vector cache is responsible for caching access control decisions throughout the OS to avoid costly up-calls and message passing (with associated context switches) to the security server. NSA and SCC eventually migrated FLASK to the FLUX microkernel developed by the University of Utah in the search for improved performance. Invigorated by the rise of microkernels and their congruence with security kernels, this flurry of operating system security research also faced the limitations (and eventual rejection) of the microkernel approach by the computer industry – which perceived the performance overheads as too great.

Microkernels and mandatory access control have seen another experimental composition in the form of Decentralized Information Flow Control (DIFC). This model, proposed by Myers, allows applications to assign information flow labels to OS-provided objects, such as communication channels, which are propagated and enforced by a blend of static analysis and runtime OS enforcement, implementing policies such as taint tracking [50] – effectively, a composition of mandatory access control and capabilities in service to application security. This approach is embodied by Efstathopoulos et al.'s Asbestos [19] and Zeldovich et al.'s Histar [92] research operating systems.

Despite the decline of both hardware-oriented and microkernel capability system design, capability models continue to interest both research and industry. Inspired by the proprietary KEYKOS system [27], Shapiro's EROS [72] (now CapROS) continues the investigation of higher-assurance software capability designs, seL4 [33], a formally verified, capability-oriented microkernel, has also continued along this avenue. General-purpose systems also have adopted elements of the microkernel capability design philosophy, such as Apple's Mac OS X [5] (which uses Mach interprocess communication (IPC) objects as capabilities) and Cambridge's Capsicum [81] research project (which attempts to blend capability-oriented design with UNIX).

---

<sup>4</sup>PSOS's ground-up design included ground-up hardware, whereas Schiller's design revised only the software stack.

More influentially, Morris’s suggestion of capabilities at the programming language level has seen widespread deployment. Gosling and Gong’s Java security model blends language-level type safety with a capability-based virtual machine [25, 24]. Java maps language-level constructs (such as object member and method protections) into execution constraints enforced by a combination of a pre-execution bytecode verification and expression constraints in the bytecode itself. Java has seen extensive deployment in containing potentially (and actually) malicious code in the web browser environment. Miller’s development of a capability-oriented E language [47], Wagner’s Joe-E capability-safe subset of Java [46], and Miller’s Caja capability-safe subset of JavaScript continue a language-level exploration of capability security [48].

### 1.6.3 Language and Runtime Approaches

Direct reliance on hardware for enforcement (which is central to both historic and current systems) is not the only approach to isolation enforcement. The notion that limits on expressibility in a programming language can be used to enforce security properties is frequently deployed in contemporary systems to supplement coarse and high-overhead operating-system process models. Two techniques are widely used: virtual-machine instruction sets (or perhaps physical machine instruction subsets) with limited expressibility, and more expressive languages or instruction sets combined with type systems and formal verification techniques.

The Berkeley Packet Filter (BPF) is one of the most frequently cited examples of the virtual machine approach: user processes upload pattern matching programs to the kernel to avoid data copying and context switching when sniffing network packet data [43]. These programs are expressed in a limited packet-filtering virtual-machine instruction set capable of expressing common constructs, such as accumulators, conditional forward jumps, and comparisons, but are incapable of expressing arbitrary pointer arithmetic that could allow escape from confinement, or control structures such as loops that might lead to unbounded execution time. Similar approaches have been used via the type-safe Modula 3 programming language in SPIN [6], and the DTrace instrumentation tool that, like BPF, uses a narrow virtual instruction set to implement the D language [11].

Google’s Native Client (NaCl) model edges towards a verification-oriented approach, in which programs must be implemented using a ‘safe’ (and easily verified) subset of the x86 or ARM instruction sets, which would allow confinement properties to be validated [91]. NaCl is closely related to Software Fault Isolation (SFI) [74], in which safety properties of machine code are enforced through instrumentation to ensure no unsafe access, and Proof-Carrying Code (PCC) in which the safe properties of code are demonstrated through attached and easily verifiable proofs [52]. As mentioned in the previous section, the Java Virtual Machine (JVM) model is similar; it combines runtime execution constraints of a restricted, capability-oriented bytecode with a static verifier run over Java classes before they can be loaded into the execution environment; this ensures that only safe accesses have been expressed. C subsets, such as Cyclone [28], and type-safe languages such as Ruby [65], offer similar safety guarantees, which can be leveraged to provide security confinement of potentially malicious code without hardware support.

These techniques offer a variety of trade-offs relative to CPU enforcement of the process model. For example, some (BPF, D) limit expressibility that may prevent potentially useful constructs from being used, such as loops bounded by invariants rather than instruction limits; in doing so, this can typically impose potentially significant performance overhead. Systems

such as FreeBSD often support just-in-time compilers (JITs) that convert less efficient virtual-machine bytecode into native code subject to similar constraints, addressing performance but not expressibility concerns [45].

Systems like PCC that rely on proof techniques have had limited impact in industry, and often align poorly with widely deployed programming languages (such as C) that make formal reasoning difficult. Type-safe languages have gained significant ground over the last decade, with widespread use of JavaScript and increasing use of functional languages such as OCaml [63]; they offer many of the performance benefits with improved expressibility, yet have had little impact on operating system implementations. However, an interesting twist on this view is described by Wong in *Gazelle*, in which the observation is made that a web browser is effectively an operating system by virtue of hosting significant applications and enforcing confinement between different applications [77]. Web browsers frequently incorporate many of these techniques including Java Virtual Machines and a JavaScript interpreter.

### 1.6.4 Bounds Checking and Fat Pointers

In contrast to prior capability systems, a key design goal for CHERI was to support mapping C-language pointers into capabilities. In earlier prototypes, we did this solely through base and bounds fields within capabilities, which worked well but required substantial changes to existing C software that often contained programming idioms that violated monotonic rights decrease for pointers. In later versions of the ISA, we adopt ideas from the C fat-pointer literature, which differentiate the idea of a delegated region from a current pointer: while the base and bounds are subject to guarded manipulation rules, we allow the offset to float within and beyond the delegated region. Only on dereference are protections enforced, allowing a variety of imaginative pointer operations to be supported. Many of these ideas originate with the type-safe C dialect Cyclone [28], and see increasing adaptation to off-the-shelf C programs with work such as Softbound [51], Hardbound [18], and CCured [53]. This flexibility permits a much broader range of common C idiom to be mapped into the capability-based memory-protection model.

### 1.6.5 Influences of Our Own Past Projects

Our CHERI capability hardware design responds to all these design trends – and their problems. Reliance on traditional paged virtual memory for strong address-space separation, as used in Mach, EROS, and UNIX, comes at significant cost: attempts to compartmentalize system software and applications sacrifice the programmability benefits of a language-based capability design (a point made convincingly by Fabry [20]), and introduce significant performance overhead to cross-domain security boundaries. However, running these existing software designs is critical to improve the odds of technology transfer, and to allow us to incrementally apply ideas in CHERI to large-scale contemporary applications such as office suites. CHERI’s hybrid approach allows a gradual transition from virtual address separation to capability-based separation within a single address space, thus restoring programmability and performance so as to facilitate fine-grained compartmentalization throughout the system and its applications.

We consider some of our own past system designs in greater detail, especially as they relate to CTSRD:



**Multics** The Multics system incorporated many new concepts in hardware, software, and programming [58, 15]. The Multics hardware provided independent virtual memory segments, paging, interprocess and intra-process separation, and cleanly separated address spaces. The Multics software provided symbolically named files that were dynamically linked for efficient execution, rings of protection providing layers of security and system integrity, hierarchical directories, and access-control lists. Input-output was also symbolically named and dynamically linked, with separation of policy and mechanism, and separation of device independence and device dependence. A subsequent redevelopment of the two inner-most rings enabled Multics to support multilevel security in the commercial product. Multics was implemented in a stark subset of PL/I that considerably diminished the likelihood of many common programming errors. In addition, the stack discipline inherently avoided buffer overflows.

**PSOS** SRI's Provably Secure Operating System hardware-software design was formally specified in a single language, with encapsulated modular abstraction, interlayer state mappings, and abstract programs relating each layer to those on which it depended [55, 56]. The hardware design provided tagged, typed, unforgeable capabilities required for every operation, with identifiers that were unique for the lifetime of the system. In addition to a few primitive types, application-specific object types could be defined and their properties enforced with the hardware assistance provided by the capability-based access controls. The design allowed application layers to efficiently execute instructions, with object-oriented capability-based addressing directly to the hardware – despite appearing at a much higher layer of abstraction in the design specifications.

**MAC Framework** The MAC Framework is an OS reference-monitor framework used in FreeBSD, also adopted in Mac OS X and iOS, as well as other FreeBSD-descended operating systems such as Juniper Junos and McAfee Sidewinder [80]. Developed in the DARPA CHATS program, the MAC Framework allows static and dynamic extension of the kernel's access-control model, supporting implementation of *security localization* – that is, the adaptation of the OS security to product and deployment-specific requirements. The MAC Framework (although originally targeted at classical mandatory access control models) found significant use in application sandboxing, especially in Junos, Mac OS X, and iOS. One key lesson from this work is the importance of longer-term thinking about security-interface design, including interface stability and support for multiple policy models; these are especially important in instruction-set design. Another important lesson is the increasing criticality of extensibility of not just the access-control model, but also the means by which remote principals are identified and execute within local systems: not only is consideration of classical UNIX users inadequate, but also there is a need to allow widely varying policies and notions of remote users executing local code across systems. These lessons are taken to heart in capability systems, which carefully separate policy and enforcement, but also support extensible policy through executable code.

**Capsicum** Capsicum is a lightweight OS capability and sandbox framework included in FreeBSD 9.x and later [81, 79]. Capsicum extends (rather than replaces) UNIX APIs, and provides new kernel primitives (sandboxed capability mode and capabilities) and a userspace sandbox API. These tools support compartmentalization of monolithic UNIX applications into logical applications, an increasingly common goal supported poorly by discretionary and manda-

tory access controls. This approach was demonstrated by adapting core FreeBSD utilities and Google’s Chromium web browser to use Capsicum primitives; it showed significant simplicity and robustness benefits to Capsicum over other confinement techniques. Capsicum provides both inspiration and motivation for CHERI: its hybrid capability-system model is transposed into the ISA to provide compatibility with current software designs, and its demand for finer-grained compartmentalization motivates CHERI’s exploration of more scalable approaches.

### 1.6.6 A Fresh Opportunity for Capabilities

Despite an extensive research literature exploring the potential of capability-system approaches, and limited transition to date, we believe that now is the time to revisit these ideas, albeit through the lens of contemporary problems and with insight gained through decades of research into security and systems design. As described earlier in the chapter, a transformed threat environment deriving from ubiquitous computing and networking, and the practical reality of widespread exploitation of software vulnerabilities, provides a strong motivation to investigate improved processor foundations for software security. This change in environment has coincided with improved hardware prototyping techniques and higher-level hardware definition languages that facilitate academic hardware-software system research at larger scales without which we would have been unable to explore the CHERI approach in such detail. Simultaneously, our understanding of operating-system and programming-language security has been vastly enhanced by several decades of research, and recent development of the hybrid capability-system Capsicum model suggests a strong alignment between capability-based techniques and successful mitigation approaches that can be translated into processor design choices.

## 1.7 Publications

As our approach has evolved, and project developed, we have published a number of papers and reports describing aspects of the work. The revisiting of capability-based approaches is described in *Capabilities Revisited: A Holistic Approach to Bottom-to-Top Assurance of Trustworthy Systems*, published at the Layered Assurance Workshop (LAW 2010) [57], shortly after the inception of the project. Mid-way through creation of both the BERI prototyping platform, and CHERI ISA model, we published *CHERI: A Research Platform Deconflating Hardware Virtualization and Protection* at the Workshop on Runtime Environments, Systems, Layering and Virtualized Environments (RESolve 2012) [78].

Most recently, we have published *The CHERI Capability Model: Revisiting RISC in an Age of Risk* at the International Symposium on Computer Architecture (ISCA 2014) [87]. A forthcoming paper in the International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS 2015), *Beyond the PDP-11: Architectural support for a memory-safe C abstract machine*, further explores C-language considerations for the CHERI ISA.

We have additionally prepared a number of technical reports, including this document, describing our approach and prototypes. To date, this includes the second version of the CHERI ISA as described in the *CHERI Instruction-Set Architecture* [84], the *CHERI User’s Guide* [83], the *BERI Hardware Reference* [85], and the *BERI Software Reference* [82]. The current report

extends the CHERI ISA for improved C-language and software object-capability support. Further research publications and technical reports on the topics of the CHERI hardware-software security model, compiler approaches, and applications of formal methods will be forthcoming.

## 1.8 Version History

This report was previously made available as the *CHERI Architecture Document*, but is now the *CHERI Instruction-Set Architecture*.

- 1.0 This first version of the CHERI architecture document was prepared for a six-month deliverable to DARPA. It included a high-level architectural description of CHERI, motivations for our design choices, and an early version of the capability instruction set.
- 1.1 The second version was prepared in preparation for a meeting of the CTSRD External Oversight Group (EOG) in Cambridge during May 2011. The update followed a week-long meeting in Cambridge, UK, in which many aspects of the CHERI architecture were formalized, including details of the capability instruction set.
- 1.2 The third version of the architecture document came as the first annual reports from the CTSRD project were in preparation, including a decision to break out formal-methods appendices into their own *CHERI Formal Methods Report* for the first time. With an in-progress prototype of the CHERI capability unit, we significantly refined the CHERI ISA with respect to object capabilities, and matured notions such as a trusted stack and the role of an operating system supervisor. The formal methods portions of the document was dramatically expanded, with proofs of correctness for many basic security properties. Satisfyingly, many ‘future work’ items in earlier versions of the report were becoming completed work in this version!
- 1.3 The fourth version of the architecture document was released while the first functional CHERI prototype was in testing. It reflects on initial experiences adapting a microkernel to exploit CHERI capability features. This led to minor architectural refinements, such as improvements to instruction opcode layout, some additional instructions (such as allowing CGetPerms retrieve the unsealed bit), and automated generation of opcode descriptions based on our work in creating a CHERI-enhanced MIPS assembler.
- 1.4 This version updated and clarified a number of aspects of CHERI following a prototype implementation used to demonstrate CHERI in November 2011. Changes include updates to the CHERI architecture diagram; replacement of the CDecLen instruction with CSetLen, addition of a CMove instruction; improved descriptions of exception generation; clarification of the in-memory representation of capabilities and byte order of permissions; modified instruction encodings for CGetLen, CMove, and CSetLen; specification of reset state for capability registers; and clarification of the CIncBase instruction.
- 1.5 This version of the document was produced almost two years into the CTSRD project. It documented a significant revision (version 2) to the CHERI ISA, which was motivated by our efforts to introduce C-language extensions and compiler support for CHERI, with improvements resulting from operating system-level work and restructuring the Bluespec

hardware specification to be more amenable to formal analysis. The ISA, programming language, and operating system sections were significantly updated.

- 1.6** This version made incremental refinements to version 2 of the CHERI ISA, and also introduced early discussion of the CHERI2 prototype.
- 1.7** Roughly two and a half years into the project, this version clarified and extended documentation of CHERI ISA features such as CCall/CReturn and its software emulation, Permit\_Set\_Type, the CMove pseudo-op, new load-linked and instructions for store-conditional relative to capabilities, and several bug fixes such as corrections to sign extension for several instructions. A new capability-coprocessor `cause` register, retrieved using a new CGetCause, was added to allow querying information on the most recent CP2 exception (e.g., bounds-check vs type-check violations); priorities were provided, and also clarified with respect to coprocessor exceptions vs. other MIPS ISA exceptions (e.g., unaligned access). This was the first version of the *CHERI Architecture Document* released to early adopters.
- 1.8** Less three and a half years into the project, this version refined the CHERI ISA based on experience with compiler, OS, and userspace development using the CHERI model. To improve C-language compatibility, new instructions CToPtr and CFromPtr were defined. The capability permissions mask was extended to add user-defined permissions. Clarifications were made to the behavior of jump/branch instructions relating to branch-delay slots and the program counter. CClearTag simply cleared a register's tag, not its value. A software-defined capability-cause register range was made available, with a new CSetCause instruction letting software set the cause for testing or control-flow reasons. New CCheckPerm and CCheckType instructions were added, letting software object methods explicitly test for permissions and the types of arguments. TLB permission bits were added to authorize use of loading and storing tagged values from pages. New CGetDefault and CSetDefault pseudo-ops have become the preferred way to control MIPS ISA memory access. CCall/CReturn calling conventions were clarified; CCall now pushes the incremented version of the program counter, as well as stack pointer, to the trusted stack.
- 1.9 - UCAM-CL-TR-850** The document was renamed from the *CHERI Architecture Document* to the *CHERI Instruction-Set Architecture*. This version of the document was made available as a University of Cambridge Technical Report. The high-level ISA description and ISA reference were broken out into separate chapters. A new rationale chapter was added, along with more detailed explanations throughout about design choices. Notes were added in a number of places regarding non-MIPS adaptations of CHERI and 128-bit variants. Potential future directions, such as capability cursors, are discussed in more detail. Further descriptions of the memory-protection model and its use by operating systems and compilers was added. Throughout, content has been updated to reflect more recent work on compiler and operating-system support for CHERI. Bugs have been fixed in the specification of the CJR and CJALR instructions. Definitions and behavior for user-defined permission bits and OS exception handling have been clarified.
- 1.10** This version of the Instruction-Set Architecture is timed for delivery at the end of the fourth year of the CTSRD Project. It reflects a significant further revision to the ISA

(version 3) focused on C-language compatibility, better exception-handling semantics, and reworking of the object-capability mechanism.

The definition of the NULL capability has been revised such that the memory representation is now all zeroes, and with a zeroed tag. This allows zeroed memory (e.g., ELF BSS segments) to be interpreted as being filled with NULL capabilities. To this end, the tag is now defined as unset, and the Unsealed bit has now been inverted to be a Sealed bit; the CGetUnsealed instruction has been renamed to CGetSealed.

A new **offset** field has been added to the capability, which converts CHERI from a simple base/length capability to blending capabilities and fat pointers that associate a base and bounds with an offset. This approach learns from the extensive fat-pointer research literature to improve C-language compatibility. The offset can take on any 64-bit value, and is added to the base on dereference; if the resulting pointer does not fall within the base and length, then an exception will be thrown. New instructions are added to read (CGetOffset) and write (CSetOffset) the field, and the semantics of memory access and other CHERI instructions (e.g., CIncBase) are updated for this new behavior.

A new CPtrCmp instruction has been added, which provides C-friendly comparison of capabilities; the instruction encoding supports various types of comparisons including ‘equal to’, ‘not equal to’, and both signed and unsigned ‘less than’ and ‘less than or equal to’ operators.

GetPCC now returns **PC** as the **offset** field of the returned **PCC** rather than storing it to a general-purpose register. CJR and CJALR now accept target **PC** values via the offsets of their jump-target capability arguments rather than via explicit general-purpose registers. CJALR now allows specification of the return-program-counter capability register in a manner similar to return-address arguments to the MIPS JALR instruction.

CCall and CReturn are updated to save and restore the saved **PC** in the **offset** field of the saved **EPCC** rather than separately. **EPCC** now incorporates the saved exception **PC** in its **offset** field. The behavior of **EPCC** and expectations about software-supervisor behavior are described in greater detail. The security implications of exception cause-code precedence as relates to alignment and the emulation of unaligned loads and stores are clarified. The behavior of CSetCause has been clarified to indicate that the instruction should not raise an exception unless the check for *Access EPCC* fails. When an exception is raised due to the state of an argument register for an instruction, it is now defined which register will be named as the source of the exception in the capability cause register.

The object-capability type field is now 24-bit; while a relationship to addresses is maintained in order to allow delegation of type allocation, that relationship is deemphasized. It is assumed that the software type manager will impose any required semantics on the field, including any necessary uniqueness for the software security model. The CSet-Type instruction has been removed, and a single CSeal instruction replaces the previous separate CSealCode and CSealData instructions.

The validity of capability fields accessed via the ISA is now defined for untagged capabilities; the undefinedness of the in-memory representation of capabilities is now explicit in order to permit ‘non-portable’ microarchitectural optimizations.

There is now a structured description of the pseudocode language used in defining instructions. Format numbers have now been removed from instruction descriptions.

Ephemeral capabilities are renamed to ‘local capabilities,’ and non-ephemeral capabilities are renamed to ‘global capabilities’; the semantics are unchanged.

**1.11 - UCAM-CL-TR-864** This version of the CHERI ISA has been prepared for publication as a University of Cambridge technical report. It includes a number of refinements to CHERI ISA version 3 based on further practical implementation experience with both C-language memory protection and software compartmentalization.

There are a number of updates to the specification reflecting introduction of the **offset** field, including discussion of its semantics. A new **CIncOffset** instruction has been added, which avoids the need to read the offset into a general-purpose register for frequent arithmetic operations on pointers.

Interactions between **EPC** and **EPCC** are now better specified, including that use of untagged capabilities has undefined behaviour. **CBTS** and **CBTU** are now defined to use branch-delay slots, matching other MIPS-ISA branch instructions. **CJALR** is defined as suitably incrementing the returned program counter, along with branch-delay slot semantics. Additional software-path pseudocode is present for **CCall** and **CReturn**.

**CAndPerm** and **CGetPerm** use of argument-register or return-register permission bits has been clarified. Exception priorities and cause-code register values have been defined, clarified, or corrected for **CClearTag**, **CGetPCC**, **CSC**, and **CSeal**. Sign or zero extension for immediates and offsets are now defined **CL**, **CS**, and other instructions.

Exceptions caused due to TLB bits controlling loading and storing of capabilities are now **CP2** rather than TLB exceptions, reducing code-path changes for MIPS exception handlers. These TLB bits now have modified semantics: **LC** now discards tag bits on the underlying line rather than throwing an exception; **SC** will throw an exception only if a tagged store would result, rather than whenever a write occurs from a capability register. These affect **CLC** and **CSC**.

Pseudocode definitions now appear earlier in the chapter, and have now been extended to describe **EPCC** behavior. The ISA reference has been sorted alphabetically by instruction name.

## 1.9 Document Structure

This document is an introduction to, and reference manual for, the CHERI instruction-set architecture.

Chapter 1 introduces the CHERI instruction-set architecture, its motivations, goals, context, philosophy, and design.

Chapter 2 provides a detailed description of the CHERI architecture, including its register and memory capability models, new instructions, procedure capabilities, and use of message-passing primitives.

Chapter 3 describes the CHERI capability coprocessor, its register file, tagged memory, and other ISA-related semantics.

Chapter 4 provides a detailed description of each new CHERI instruction, its pseudo-operations, and how compilers should handle floating-point loads and stores via capabilities.

Chapter 5 discusses some of the design choices that we have made thus far.

Chapter 6 discusses the programming language and operating system implications of CHERI, including its impact on operating-system kernels, language runtimes, and compilers.

Chapter 7 discusses our short- and long-term plans for the CHERI architecture, considering both our specific plans and open research questions that must be answered as we proceed. It also addresses some CHERI-related work that is emerging from our companion project for the DARPA Mission-oriented Resilient Cloud program.

Future versions of this document will continue to expand our consideration of the CHERI instruction-set architecture and its impact on software, as well as evaluation strategies and results. Additional information on our CHERI hardware and software implementations, as well as formal methods work, can be found in accompanying reports.

# Chapter 2

## CHERI Capability-System Model

In this chapter we describe the high-level design for the CHERI Instruction-Set Architecture (ISA) and consider both the semantics and mechanism of CHERI’s memory and object capabilities. We discuss CHERI in relative isolation from the general-purpose ISA, as our approach might reasonably apply to a number of RISC ISAs (e.g., MIPS, ARM, and RISC-V), but potentially also to CISC ISAs (such as Intel and AMD 32-bit and 64-bit ISAs).<sup>1</sup> In Chapters 3 and 4, we consider in detail an instantiation of the CHERI model as an extension to the 64-bit MIPS ISA.

### 2.1 Design Goals

As described in Chapter 1, the key observation motivating the CHERI design is that page-oriented virtual memory, nearly universal in commodity CPUs, is neither an efficient nor a programmer-friendly primitive for fine-grained memory protection or scalable hardware-supported compartmentalization. Virtual addressing, implemented by a memory management unit (MMU) and translation look-aside buffer (TLB), clearly plays an important role by dissociating physical memory allocation and address-space management, facilitating software features such as strong separation, OS virtualization, and virtual-memory concepts such as swapping and paging. However, with a pressing need for scalable and fine-grained separation, the overheads and programmability difficulties imposed by virtual addressing as the sole primitive for hardware isolation actively deter employment of the principle of least privilege. These concerns translate into three high-level security design goals for CHERI:

1. Management of security context must be a “fast path” that avoids expensive operations such as TLB entry invalidation, frequent ring transitions, and cache-busting OS supervisor paths. This is a natural consequence of the integral role security functions (such as creation, refinement, and delegation of memory and object rights) play in fine-grained compartmentalized code.

---

<sup>1</sup>One idea we have considered is that CHERI-like semantics might be accomplished through extensions to Google’s Native Client ISA [90], which uses a strict and statically analyzable subset of Intel and ARM ISAs to ensure memory safety. An exciting possibility is that we might extend the LLVM intermediate representation [38] to capture notions of segmentation and capability protection, in which case either NaCl or CHERI back ends might be targeted as underlying execution substrates. This naturally raises the question, “why new hardware” – one that we have constantly in mind, and believe will be constructively answerable in terms of functionality, performance, and formal assurances.



2. Security domain switches must be inexpensive and efficient, with cost scaling linearly with the number of switches and actual code/data footprint (and hence general-purpose cache performance), rather than scaling as a product of the number of security domains and controlled objects regardless of code and data cache footprints. CHERI is intended to support at least two orders of magnitude more active security domains per CPU than current MMU-based systems (going from tens or hundreds to at least tens of thousands of domains).
3. Security domain switches must allow shared object namespaces that provide a unified view that connotes both efficient and programmer-comprehensible delegation. Compartmentalized applications should be able to be programmed and debugged without unnecessary recourse to distributed-system methodology.

These security goals, combined with observations about TLB performance and a desire to compartmentalize existing single-address-space applications, led us to the conclusion that new instruction set primitives for memory and object control *within an address space* would usefully complement existing address-space-based separation. In this view, security state associated with a thread should be captured as a set of registers that can be explicitly managed by code, and be preserved and restored cheaply on either side of security domain transitions – in effect, part of a thread’s register file. In the parlance of contemporary CPU and OS design, this establishes a link between hardware threads (OS threads) and security domains, rather than address spaces (OS processes) and security domains.

Because we wish to consider delegation of memory and object references within an address space as a first-class operation, we choose to expose these registers to the programmer (or, more desirably, the compiler) so that they can be directly manipulated and passed as arguments. Previous systems built along these principles have been referred to as *capability systems*, a term that also usefully describes CHERI.

CHERI’s capability model represents an explicit capability system, in which common capability manipulation operations are unprivileged instructions and transfer of control to a supervisor during regular operations is avoided. In historic capability systems, microcode (or even the operating system) was used to implement complex capability operations, some of which were privileged. In contemporary RISC CPU designs, the intuitive functional equivalent has an exception that triggers the supervisor. However, entrance to a supervisor usually remains an expensive operation, and hence one to avoid in high-performance paths. In keeping with the RISC design philosophy, we are willing to delegate significant responsibility for safety to the compiler and run-time linker to minimize hardware knowledge of higher-level language constructs.

CHERI capabilities may refer to *regions of memory*, with bounded memory access (as in *segments*). Memory capabilities will frequently refer to programmer-described data structures such as strings of bytes, structures consisting of multiple fields, and entries in arrays, although they might also refer to larger extents of memory (e.g., the entire address space). While compatibility features in the CHERI ISA allow programmers to continue to use pointers in legacy code, we anticipate that capabilities will displace pointer use as code is migrated to CHERI code generation, providing stronger integrity for data references, bounds checking, permission checking, and so on. In our prototype extensions to the C language, programmers can explicitly request that capabilities be used instead of pointers, providing stronger protection, or in some cases rely on the compiler to automatically generate capability-aware code – for example, when

code accessing the stack is compiled with a suitable application binary interface (ABI). We are exploring further static analysis and compilation techniques that will allow us to automate deployment of capability-aware code to a greater extent, minimizing disruption of current source code while allowing programs to experience protection improvements.

Alternatively, capabilities may refer to *objects* that can be *invoked*, which allows the implementation of *protected subsystems* – i.e., services that execute in a security domain other than the caller’s. At the moment of object invocation, caller capabilities are *sealed* to protect them from inappropriate use by the callee, and the invoked object is *unsealed* to allow the object callee to access private resources it requires to implement its services. The caller and callee experience a controlled delegation of resources across object invocation and return. For example, the caller might delegate access to a memory buffer, and the callee might then write a Unicode string to the buffer describing the contents of the protected object, implementing call-by-reference.<sup>2</sup> A key goal has been to allow capabilities passed across protection-domain boundaries to refer to ordinary C data on the stack or heap, allowing easier adaptation of existing programs and libraries to use CHERI’s features. The semantics of capabilities are discussed in greater detail later in this and the following chapter.

## 2.2 A Hybrid Capability-System Architecture

Despite our complaints about the implications of virtual addressing for compartmentalization, we feel that virtual memory is a valuable hardware facility: it provides a strong separation model; it makes implementing facilities such as swapping and paging easier; and by virtue of its virtual layout, it can significantly improve software maintenance and system performance. CHERI therefore adopts a *hybrid capability-system model*: we retain support for a commodity virtual-memory model, implemented using an MMU with a TLB, while also introducing new primitives to permit multiple security domains within address spaces (Figure 2.1). Each address space becomes its own decomposition domain, within which protected subsystems can interact using both hierarchical and non-hierarchical security models. In effect, each address space is its own *virtual capability machine*.

To summarize our approach, CHERI draws on two distinct, and previously uncombined, designs for processor architecture:

- *Page-oriented virtual memory systems* allow an executive (often the operating system kernel) to create a *process abstraction* via the MMU. In this model, the kernel is responsible for maintaining separation using this relatively coarse tool, and then providing system calls that allow spanning process isolation, subject to access control. Systems such as this make only weak distinctions between code and data, and in the mapping from programming language to machine code discard most typing and security information.
- *Capability systems*, often based on a single global address space, map programming-language type information and protection constraints into instruction selection. Code at any given moment in execution exists in a protection domain consisting of a dynamic set

---

<sup>2</sup>CHERI does not implement implicit rights *amplification*, a property of some past systems including HYDRA. Callers across protected subsystem boundaries may choose to pass all rights they hold, but it is our expectation that they will generally not do so – otherwise, they would use regular function calls within a single protected subsystem.

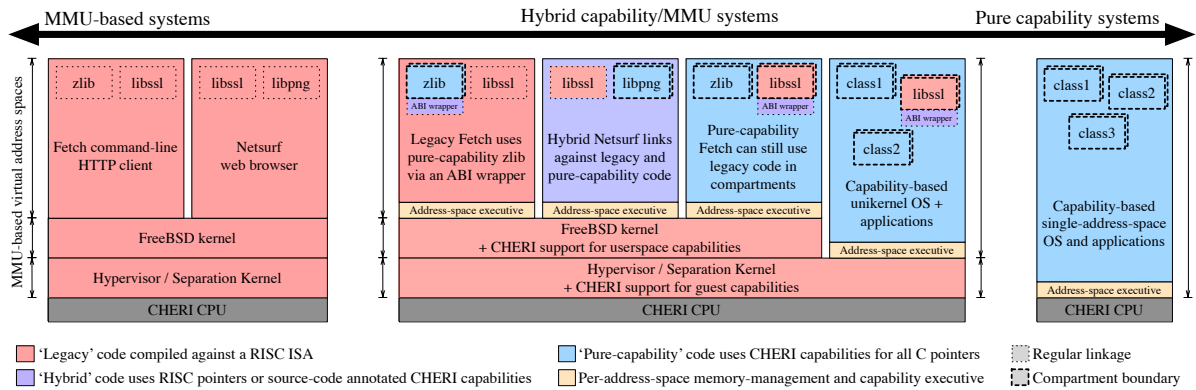


Figure 2.1: CHERI’s hybrid capability architecture: initially, legacy software components execute without capability awareness, but security-sensitive TCB elements or particularly risky code bases are converted. In the long term, all packages are converted, implementing least privilege throughout the system.

of rights whose delegation is controlled by the flow of code. (These *instantaneous rights* are sometimes referred to as *spheres of protection* in the operating system and security literature.) Such a design generally offers greater assurance, because the *principle of least privilege* can be applied at a finer granularity.

Figure 2.1 illustrates the following alternative ways in which the CHERI architecture might be used. In CHERI, even within an address space, existing and capability-aware code can be hybridized, as reads and writes via general-purpose MIPS registers are automatically indirected through a reserved capability register before being processed by the MMU. This allows a number of interesting compositions, including the execution of capability-aware, and hence significantly more robust, libraries within a legacy application. Another possibility is a capability-aware application running one or more instances of capability-unaware code within sandboxes, such as legacy application components or libraries – effectively allowing the trivial implementation of the Google Native Client model.

Finally, applications can be compiled to be fully capability-aware, i.e., able to utilize the capability features for robustness and security throughout their structure. The notion of a capability-aware *executive* also becomes valuable – likely as some blend of the run-time linker and low-level system libraries (such as `libc`): the executive will set up safe linkage between mutually untrusting components (potentially with differing degrees of capability support, and hence differing ABIs), and ensure that memory is safely managed to prevent memory-reuse bugs from escalating to security vulnerabilities.<sup>3</sup> Useful comparison might also be made between our notion of an in-address-space executive and a microkernel, as the executive will similarly take responsibility for configuring protection and facilitating controlled sharing of data. As microkernels are frequently capability-based, we might find that not only are ideas from the microkernel space reusable, but also portions of their implementations. This is an exciting prospect, especially considering that significant effort has been made to apply formal verification techniques to microkernels.

<sup>3</sup>Similar observations about the criticality of the run-time linker for both security and performance in capability systems have been made by Karger [32].

## 2.3 The CHERI Software Stack

The notion of hybrid design is key to our adoption argument: CHERI systems are able to execute today's commodity operating systems and applications with few modifications. Use of capability features can then be selectively introduced in order to raise confidence in the robustness and security of individual system components, which are able to fluidly interact with other un-enhanced components. This notion of hybrid design first arose in Cambridge's Capsicum [81] (which blends the POSIX Application Programming Interface (API), as implemented in the FreeBSD operating system) with a capability design by allowing processes to execute in hybrid mode or in *capability mode*. Traditional POSIX code can run along side capability-mode processes, allowing the construction of sandboxes; using a capability model, rights can be delegated to these sandboxes by applications that embody complex security policies. One such example from our USENIX Security 2010 Capsicum paper [81] is the Chromium web browser, which must map the distributed World Wide Web security model into local OS containment primitives.

CHERI's software stack will employ hybrid design principles from the bottom up: capability-enhanced separation kernels will be able to implement both conventional virtual-machine interfaces to guest operating systems, or directly host capability-aware operating systems or applications, ensuring robustness. This would provide an execution substrate on which both commodity systems built on traditional RISC instruction models (such as FreeBSD) can run side by side with a pure capability-oriented software stack, such as capability-adapted language runtimes. Further, CheriBSD, a CHERI-enhanced version of the FreeBSD operating system, and its applications, will be able to employ CHERI features in their own implementations. For example, key data-processing libraries, such as image compression or video decoding, might use CHERI features to limit the impact of programming errors through fine-grained memory protection, but also apply compartmentalization to mitigate logical errors through the principle of least privilege. We have extended the existing Clang/LLVM compiler suite to support C-language extensions for capabilities, allowing current code to be recompiled to use capability protections based on language-level annotations, but also to link against unmodified code.

To this end, the CHERI ISA design allows software context to address memory either via legacy MIPS ISA load and store instructions, which implicitly indirect through a reserved capability register configurable by software, or via new capability load and store instructions that allow the compiler to explicitly name the object to be used. In either case, access is permitted to memory only if it is authorized by a capability that is held in the register file (or, by transitivity, any further capability that can be retrieved using those registers and the memory or objects that it can reach). New ABIs and calling conventions are defined to allow transition between (and across) CHERI-ISA and MIPS-ISA code to allow legacy code to invoke capability-aware code, and vice versa. For example, in this model CheriBSD might employ capability-oriented instructions in the implementation of risky data manipulations (such as network-packet processing), while still relying on traditionally written and compiled code for the remainder of the kernel. Similarly, within the Chromium web browser, the JavaScript interpreter might be implemented in terms of capability-oriented instructions to offer greater robustness, while the remainder of Chromium would use traditional instructions.

One particularly interesting property of our hardware design is that capabilities can take on different semantics within different address spaces, with each address-space's executive integrating memory management and capability generation. In the CheriBSD kernel, for example,

virtual addressing and capability use can be blended; the compiler and kernel memory allocator can use capabilities for certain object types, but not for others. In various userspace processes, a hybrid UNIX / C runtime might implement limited pools of capabilities for specially compiled components, but another process might use just-in-time (JIT) compilation techniques to map Java bytecode into CHERI instructions, offering improved performance and a significantly smaller and stronger Java TCB.

Capabilities supplement the purely hierarchical ring model with a non-hierarchical mechanism – as rings support traps, capabilities support protected subsystems. One corollary is that the capability model could be used to implement rings within address spaces. This offers some interesting opportunities, not least the ability to implement purely hierarchical models where desired; for example, a separation kernel might use the TLB to support traditional OS instances, but only capability protections to constrain an entirely capability-based OS. A further extreme is to use the TLB only for paging support, and to implement a single-address-space operating system as envisioned by the designers of many historic capability systems.

This hybrid view offers a vision for a gradual transition to stronger protections, in which individual libraries, applications, and even whole operating systems can incrementally adopt stronger hardware memory protections without sacrificing the existing software stack. Discussion of these approaches also makes clear the close tie between memory-oriented protection schemes and the role of the memory allocator, an issue discussed in greater depth later in this chapter.

## 2.4 Capability Model

Chapter 3 provides detailed documentation of the registers, capabilities, and new instructions currently defined in CHERI. These concepts are briefly introduced here.

### 2.4.1 Capabilities are for Compilers

Throughout, we stress the distinction between the notion of the hardware security model and the programming model; unlike in historic CISC designs, and more in keeping with historic RISC designs, CHERI instructions are intended to support the activities of the compiler, rather than be directly programmed by application authors. While there is a necessary alignment between programming language models for computation (and in the case of CHERI, security) and the hardware execution substrate, the purpose of CHERI instructions is to make it possible for the compiler to *cleanly* and *efficiently* implement higher-level models, and not implement them directly. As such, we differentiate the idea of a *hardware capability type* from a *programming language type* – the compiler writer may choose to conflate the two, but this is an option rather than a requirement.

### 2.4.2 Capabilities

Capabilities are unforgeable tokens of authority through which programs access all memory and services within an address space. Capabilities may be held in capability registers, where they can be manipulated or dereferenced using capability coprocessor instructions, or in memory. Capabilities themselves may refer to memory (unsealed capabilities) or objects (sealed

capabilities). Memory capabilities are used as arguments to load and store instructions, to access either data or further capabilities. Object capabilities may be invoked to transition between protection domains using call and return instructions.

Unforgeability is implemented by two means: tag bits and controlled manipulation. Each capability register, and each capability-aligned physical memory location, is associated with a tag bit indicating that a capability is valid. Attempts to directly overwrite a capability in memory using data instructions automatically clear the tag bit. When data is loaded into a capability register, its tag bit is also loaded; while data without a valid tag can be loaded into a capability register, attempts to dereference or invoke such a register will trigger an exception.

Controlled manipulation is enforced by virtue of the ISA: instructions that manipulate capability register fields (e.g., base, length, permissions, type) are not able to increase the rights associated with a capability. Similarly, sealed capabilities can be unsealed only via the invocation mechanism, or via the unseal instruction subject to similar monotonicity rules. This enforces encapsulation, and prevents unauthorized access to the internal state of objects.

We anticipate that many languages will expose capabilities to the programmer via pointers or references – e.g., as qualified pointers in C, or mapped from object references in Java. In general, we expect that languages will not expose capability registers to management by programmers, instead using them for instruction operands and as a cache of active values, as is the case for general-purpose registers today. On the other hand, we expect that there will be some programmers using the equivalent of assembly-language operations, so that system security cannot rely solely on the goodness of compilers.

### 2.4.3 Capability Registers

CHERI supplements the 32 general-purpose, per-hardware thread registers provided by the MIPS ISA with 32 additional capability registers. Where general-purpose registers describe the computation state of a software thread, capability registers describe its instantaneous rights within an address space. A thread's capabilities potentially imply a larger set of rights, which may be loaded via held capabilities, which may notionally be considered as the protection domain of a thread.

There are also several implicit capability registers associated with each hardware thread, including a memory capability that corresponds to the instruction pointer, and capabilities used during exception handling. This is structurally congruent to implied registers and system control coprocessor (CPO) registers found in the base MIPS ISA.

Each capability register has 256 bits; unlike general-purpose registers, capability registers are structured, and contain a number of fields with defined semantics and constrained values:

**Sealed bit** If the sealed bit is unset, the capability describes a *memory segment* that is accessible via load and store instructions. If it is set, the capability describes an *object capability*, which can be accessed only via *object invocation*.

**Permissions** The permissions mask controls operations that may be performed using the capability; some permissions are hardware-defined, controlling instructions that may be used with the capability; others are software-defined and intended for use with the object-capability mechanism.

**Object type** Notionally, a software-managed *object type* used to ensure that corresponding code and data capabilities for an object are used together correctly.

**Base** This is the base address of a memory region.

**Length** This defines the length of a memory region.

**Offset** A free-floating pointer that will be added to the base when dereferencing a capability; the value can float outside of the range described by the capability, but an exception will be thrown if a requested access is out of range.

**Reserved fields** These bits are reserved for future experimentation.

**Tag bit** The tag bit is not part of the base 256 bits. It indicates whether or not the capability register holds a valid capability; this allows non-capability values to be moved via capability registers, making it possible to implement software functions that, for example, copy memory oblivious to capabilities being present.

We have discussed a number of schemes to reduce overhead implied by the quite sizable capability register file. 32 registers is nicely symmetric with the MIPS ISA, but in practice leads to a substantial overhead; we have considered reducing the number to 16 or even 8 to reduce hardware resource, cache footprint, and context-switch time. We have also pondered schemes to reduce the size of capability registers themselves; simple reduction of the addresses to smaller numbers to 40-bit from 64-bit, reflecting the largest virtual addresses supportable in the MIPS TLB, might allow reduction to 128-bit. Another approach might be to differentiate larger object capabilities, which require types, from pure memory capabilities, which could be represented more compactly, but would require compilers to handle multiple capability sizes. Finally, more complex techniques, such as Low-Fat Pointers [34] might also prove useful. Once again, hardware specifications written in Bluespec allow considerable flexibility and ease of understanding among these options, which will be particularly valuable as we get further into detailed experimentation, simulation, and modeling.

*Object invocation* is a central operation in the CHERI ISA, as it implements protected subsystem domain transitions that atomically update the set of rights (capabilities) held by a hardware thread, and that provide a trustworthy return path for later use. When an object capability is invoked, its data and code capabilities are *unsealed* to allow access to per-object instance data and code execution. Rights may be both acquired and dropped on invocation, allowing non-hierarchical security models to be implemented. Strong typing and type checking of objects, a notion first introduced in PSOS's *type enforcement*, serves functions both at the ISA level – providing object atomicity despite the use of multiple independent capabilities to describe an object – and to support language-level type features. For example, types can be used to check whether additional object arguments passed to a method are as they should be. As indicated earlier, the hardware capability type may be used to support language-level types, but should not be confused with language-level types.

## 2.4.4 Memory Model

In the abstract, capabilities are unforgeable tokens of authority. In the most reductionist sense, the CHERI capability namespace is the virtual address space, as all capabilities name (and authorize) actions on addresses. CHERI capabilities are unforgeable by virtue of capability register semantics and tagged memory, and act as tokens of authority by virtue of memory segments and object capability invocation.

However, enforcement of uniqueness over time is a property of the software memory allocation policy. More accurately, it is a property of virtual address-space allocation and reuse, which rests in a memory model composed from the capability mechanism, virtual address space configuration, and software language-runtime memory allocation.

This issue has presented a significant challenge in the design of CHERI: how can we provide sufficient mechanism to allow memory management, fundamentally a security operation in capability systems, while not overly constraining software runtimes regarding the semantics they can implement? Should we provide hardware-assisted garbage collection along the lines of the Java Virtual Machine’s garbage collection model? Should we implement explicit revocation functionality, along the lines of Redell’s capability revocation scheme (effectively, a level of indirection for all capabilities, or selectively when the need for revocation is anticipated)?

We have instead opted for dual semantics grounded in the requirements of real-world low-level system software: CHERI lacks a general revocation scheme; however, in coordination with the software stack, it can provide for both strict limitations on the extent of hardware-supported delegation periods, and software-supported generalized revocation using interposition. The former is intended to support the brief delegation of arguments from callers to callees across object-capability invocation; the latter allows arbitrary object reference revocation at a greater price.

### 2.4.5 Local Capabilities and Revocation

To this end, capabilities may be further tagged as *local*, which allows them to be processed in registers, stored in constrained memory regions, and passed on via invocation of other objects. The goal of local capabilities is to introduce a limited form of *revocation* that is appropriate for temporary delegation across protected subsystem invocations, which are not permitted to persist beyond that invocation. Among other beneficial properties, local capabilities allow the brief delegation of access to arguments passed by reference, such as regions of the caller’s stack (a common paradigm in C language programming).

In effect, *local capabilities* inspire a single-bit information-flow model, bounding the potential spread of capabilities for ephemeral objects to capability registers and limited portions of memory. The desired protection property can be enforced through appropriate memory management by the address-space executive: that is, local capabilities can be limited to a particular thread, with bounded delegation time down the (logical) stack.<sup>4</sup>

Generalized revocation is not supported directly by the CHERI ISA; instead, we rely on the language runtime to implement either a policy of *virtual address non-reuse* or *garbage collection*. A useful observation is that address space non-reuse is not the same as memory non-reuse: the meta-data required to support sparse use of a 64-bit address space scales with actual allocation, rather than the span of consumed address space. For many practical purposes, a 64-bit address space is *virtually* infinite<sup>5</sup>, so causing the C runtime to not reuse address space is now

---

<sup>4</sup>It has been recommended that we substitute a generalized generation count-based model for an information flow model. This would be functionally identical in the local capability case, used to protect per-stack data. However, it would also allow us to implement protection of thread-local state, as well as garbage collection, if desired. The current ISA does not yet reflect this planned change.

<sup>5</sup>As is 640K of memory. It has also not escaped our notice that there is a real OS cost to maintaining the abstraction of virtual memory; one merit to our approach is that it will deemphasize the virtual memory as a protection system, potentially reducing that overhead.



a realistic option. Software can, however, make use of interposition to implement revocation or other more semantically rich notions of privilege narrowing, as proposed in HYDRA.

### 2.4.6 Notions of Privilege

In operating-system design, *privileges* are a special set of rights exempting a component from the normal protection and access-control models – perhaps for the purposes of system bootstrapping, system management, or low-level functionality such as direct hardware access. In CHERI, three notions of privilege are defined – two in hardware, and a new notion of privilege in software relating to the interactions of capability security models between rings.

*Ring-based privilege* is derived from the commodity hardware notion that a series of successively higher-level rings provides progressively fewer rights to manage hardware protection features, such as TLB entries – and consequently potentially greater integrity, reliability, and resilience overall (as in Multics). Attempts to perform privileged instructions will trap to a lower ring level, which may then proceed with the operation, or reject it. CHERI extends this notion of privilege into the new capability coprocessor, authorizing certain operations based on the ring in which a processor executes, and potentially trapping to the next lower ring if an operation is not permitted. The trap mechanism itself is modified in CHERI, in order to save and restore the capability register state required within the execution of each ring – to authorize appropriate access for the trap handler.

*Hardware capability context privilege* is a new notion of privilege that operates within rings, and is managed by the capability coprocessor. When a new address space is instantiated, code executing in the address space is provided with adequate initial capabilities to fully manage the address space, and derive any required capabilities for memory allocation, code linking, and object-capability type management. In CHERI, all capability-related privileges are captured by capabilities, and capability operations never refer to the current processor ring to authorize operations, although violation of a security property (i.e., an attempt to broaden a memory capability) will lead to a trap, and allow a software supervisor in a lower ring to provide alternative semantics. This approach follows the spirit of Paul Karger’s paper on limiting the damage potential of discretionary Trojan horses [30], and extends it further.

*Supervisor-enforced capability context privilege* is a similar notion of privilege that may also be implemented in software trap handlers. For example, an operating system kernel may choose to accept system-call traps only from appropriately privileged userspace code (e.g., by virtue of holding a capability with full access to the userspace address space, rather than just narrow access, or that has a reserved user-defined permission bit set), and therefore can check the capability registers of the saved context to determine whether the trap was from an appropriate execution context. This might be used to limit system call invocation to a specific protected subsystem that imposes its own authorization policy on application components by safely wrapping system calls from userspace.

### 2.4.7 Traps, Interrupts, and Exception Handling

As in MIPS, traps and interrupts remain the means by which ring transitions are triggered in CHERI. They are affected in a number of ways by the introduction of capability features:

**New exceptions** New exception opportunities are introduced for both existing and new instructions, which may trap if insufficient rights are held, or an invalid operation is requested. For example, attempts to read a capability from memory using a capability without the read capability permission will trigger a trap.

**Reserved capability registers for exception handling** New exception-handling functionality is required to ensure that exception handlers themselves can execute properly. We reserve several capability registers for use both by the exception-handling mechanism itself (describing the rights that the exception handler will run with) and for use by software exception handlers (a pair of reserved registers that can be used safely during context switching). This approach is not dissimilar from the current notion of exception-handling registers in the MIPS ABI, which reserves two general-purpose registers for this purpose. However, whereas the MIPS ABI simply dictates that user code cannot rely on the two reserved exception registers being preserved, CHERI requires that access is blocked, as capability registers delegate rights and also hold data. We currently grant access to exception-related capability registers by virtue of special permission bits on the capability that describe the currently executing code; attempting to access reserved registers without suitable permission will trigger an exception.

**Saved program-counter capability** Exception handlers must also be able to inspect exception state; for example, as **PC**, the program counter, is preserved today in a control register, **EPC**, the program counter capability must be preserved as **EPCC** so that it can be queried.

**Implications for pipelining** Another area of concern in the implementation is the interaction between capability registers and pipelining. Normally, writing to TLB control registers in CP0 occurs only in privileged rings, and the MIPS ISA specifies that a number of no-op instructions follow TLB register writes in order to flush the pipeline of any inconsistent or intermediate results. Capability registers, on the other hand, may be modified from unprivileged code, which cannot be relied upon to issue the required no-ops. This case can be handled through the squashing of in-flight instructions, which may add complexity to pipeline processing because incorrect handling could otherwise lead to serious vulnerabilities.

## 2.4.8 Tagged Memory

As with general-purpose registers, storage capability register values in memory is desirable – for example, to push capabilities onto the stack, or manipulate arrays of capabilities. To this end, each capability-aligned and capability-sized word in memory has an additional *tag bit*. The bit is set whenever a capability is atomically written from a register to an authorized memory location, and cleared if a write occurs to any byte in the word using a general-purpose store instruction. Capabilities may be read only from capability-aligned words, and only if the tag bit is set at the moment of load; otherwise, a capability load exception is thrown. Tags are associated with physical memory locations, rather than virtual ones, such that the same memory mapped at multiple points in the address space, or in different address spaces, will have the same tags.

Tags require strong coherency with the data they protect, and it is expected that tags will be cached with the memory they describe within the cache hierarchy. Strong atomicity properties are required such that it is not possible to partially overwrite a capability value in memory while

retaining the tag. This proves a set of properties that falls out naturally from current coherent memory-subsystem designs.

Additional bits are present in TLB entries to indicate whether a given memory page is configured to have capabilities loaded or stored for the pertinent address space identifier (ASID). For example, this allows the kernel to set up data sharing between two address spaces without permitting capability sharing (which, as capability interpretation is scoped to address spaces, might lead to undesirable security or programmability properties). Special instructions allow the supervisor to efficiently extract and set tag bits for ranges of words within a page for the purpose of more easily implemented paging of capability memory pages. Use of these instructions is conditioned on notions of ring and capability context privilege.

## 2.4.9 Capability Instructions

Various newly added instructions are documented in detail in Chapter 3. Briefly, these instructions are used to load and store via capabilities, load and store capabilities themselves, manage capability fields, invoke object capabilities, and create capabilities. Where possible, the structure and semantics of capability instructions have been aligned with similar core MIPS instructions, similar calling conventions, and so on. The number of instructions has also been minimized to the extent possible.

## 2.4.10 Object Capabilities

As noted above, the CHERI design calls for two forms of capabilities: capabilities that describe regions of memory and offer bounded-buffer “segment” semantics, and object capabilities that permit the implementation of protected subsystems. In our model, object capabilities are represented by a pair of sealed code and data capabilities, which provide the necessary information to implement a protected subsystem domain transition. Object capabilities are “invoked” using the CCall instruction (which is responsible for unsealing the capabilities, performing a safe security-domain transition, and argument passing), followed by CReturn (which reverses this process and handles return values).

In traditional capability designs, invocation of an object capability triggered microcode responsible for state management. Initially, we have implemented CCall and CReturn as software exception handlers in the kernel, but are now exploring optimizations in which CCall and CReturn perform a number of checks and transformations to minimize software overhead. In the longer term, we hope to investigate the congruence of object-capability invocation with message-passing primitives between hardware threads: if each register context represents a security domain, and one domain invokes a service offered by another domain, passing a small number of general-purpose and capability registers, then message passing may offer a way to provide significantly enhanced performance.<sup>6</sup> In this view, hardware thread contexts, or register files, are simply caches of thread state to be managed by the processor.

---

<sup>6</sup>This appears to be another instance of the isomorphism between explicit message passing and shared memory design. If we introduce hardware message passing, then it will in fact blend aspects of both models and use the explicit message-passing primitive to cleanly isolate the two contexts, while still allowing shared arguments using pointers to common storage, or delegation using explicit capabilities. This approach would allow application developers additional flexibility for optimization.

Significant questions then arise regarding rendezvous: how can messages be constrained so that they are delivered only as required, and what are the interactions regarding scheduling? While this structure might appear more efficient than a TLB (by virtue of not requiring objects with multiple names to appear multiple times), it still requires an efficient lookup structure (such as a TCAM).

In either instantiation, a number of design challenges arise. How can we ensure safe invocation and return behavior? How can callers safely delegate arguments by reference for the duration of the call to bound the period of retention of a capability by a callee (which is particularly important if arguments from the call stack are passed by reference)?

How should stacks themselves be handled in this light, since a single logical stack will arguably be reused by many different security domains, and it is undesirable that one domain in execution might ‘pop’ rights from another domain off of the stack, or reuse a capability to access memory previously used as a call-by-reference argument.

These concerns argue for at least three features: a logical stack spanning many stack fragments bound to individual security domains, a fresh source of ephemeral stacks ready for reuse, and some notion of a do-not-transfer facility in order to prevent the further propagation of a capability (perhaps implemented via a revocation mechanism, but other options are readily apparent). PSOS explored similar notions of propagation-limited capabilities with similar motivations.

Our current software CCall/CReturn maintains a ‘trusted stack’ in the kernel address space and provides for reliable return, but it is clear that further exploration is required. Our goal is to support many different semantics as required by different programming languages, from an enhanced C language to Java. By adopting a RISC-like approach, in which traps to a lower ring occur when hardware-supported semantics is exceeded, we will be able to supplement the hardware model through modifications to the supervisor.

### **2.4.11 Peripheral Devices**

As described in this chapter, our capability model is a property of the instruction set architecture of a CHERI CPU, and imposed on code executing on the CPU. However, in most computer systems, Direct Memory Access (DMA) is used by peripheral devices to transfer data into and out of system memory without explicit instruction execution for each byte transferred: device drivers configure and start DMA using control registers, and then await completion notification through an interrupt or by polling. Used in isolation, nothing about the CHERI ISA implies that device memory access would be constrained by capabilities.

This raises a number of interesting questions. Should DMA be forced to pass through the capability equivalent of an I/O MMU in order to be appropriately constrained? How might this change the interface to peripheral devices, which currently assume that physical addresses are passed to them? Certainly, reuse of current peripheral networking and video devices with CHERI CPUs while maintaining desired security properties is desirable.

For the time being, device drivers continue to hold the privilege for DMA to use arbitrary physical memory addresses, although hybrid models – such as allowing DMA access only to specific portions of physical memory – may prove appropriate. Similar problems have plagued virtualization in commodity CPUs, where guest operating systems require DMA memory performance but cannot be allowed arbitrary access to physical memory. Exploring I/O MMU-like models and their integration with capabilities is high on our todo list; one thing is certain, how-

ever: a combination of hardware- and software-provided cache and memory management must ensure that tags are suitably cleared when capability-oblivious devices write to memory, in order to avoid violation of capability integrity properties.

In the longer term, one quite interesting idea is embedding CHERI support in peripheral devices themselves, to require the device to implement a CHERI-aware TCB that would synchronize protection information with the host OS. This type of model appeals to ideas from heterogeneous computing, and is one we hope to explore in greater detail in the future. Another alternative would be to pursue the notion of smart buses used by peripherals – for example, making them capability aware.

# Chapter 3

## CHERI Instruction-Set Architecture

This chapter describes an application of the CHERI approach to the 64-bit MIPS ISA. New instructions are implemented as a MIPS coprocessor – coprocessor 2 – an encoding space reserved for ISA extensions. In addition to adding new instructions, some behaviors have been modified in CHERI – notably, those of some standard MIPS instructions, TLBs, and exception handling. For example, existing memory load and store instructions are now implicitly indirectioned through a capability in order to enforce permissions, rebasing, and bounds checking on legacy code.

### 3.1 Capability Registers

Table 3.1 illustrates capability registers defined by the capability coprocessor. CHERI defines 28 general-purpose capability registers, which may be named using most capability register instructions. These registers are intended to hold the working set of rights required by in-execution code, intermediate values used in constructing new capabilities, and copies of capabilities retrieved from **EPCC** and **PCC** as part of the normal flow of code execution, which is congruent with current MIPS-ISA exception handling via coprocessor 0. Four capability registers have special functions and are accessible only if allowed by the permissions field **C0**. Note that **C0** and **C27 (IDC)** also have hardware-specific functions, but are otherwise general-purpose capability registers.

Each capability register also has an associated tag indicating whether it currently contains a valid capability. Any load and store operations via an invalid capability will trap.

#### Conventions for Capability Register Use

We are developing a set of ABI conventions regarding use of the other software-managed capability registers similar to those for general-purpose registers: caller-save, callee-save, a stack capability register, etc.

The current convention used by LLVM makes the following reservations for calls within a protection domain:

- **C1-C2** are caller-save. During a cross-domain call, these are used to pass the **PCC** and **IDC** values, respectively. In the invoked context, they are always available as temporaries, irrespective of whether the function was invoked as the result of a cross-domain call.

Register(s)	Description
<b>PCC</b>	Program counter capability ( <b>PCC</b> ); the capability through which <b>PC</b> is indirected by the processor when fetching instructions.
<b>C0</b>	Capability register through which all non-capability load and store instructions are indirected. This allows legacy MIPS code to be controlled using the capability coprocessor.
<b>C1...C23</b>	General-purpose capability registers referenced explicitly by capability-aware instructions.
<b>RCC (C24)</b>	Return code capability; after a <b>CJALR</b> instruction, the previous value of <b>PCC</b> is saved in <b>RCC</b> .
<b>C25</b>	General-purpose capability register reserved for use in exception handling.
<b>IDC (C26)</b>	Invoked data capability; the capability that was unsealed at the last protected procedure call. This capability holds the unlimited capability at boot time.
<b>KR1C (C27)</b>	A capability reserved for use during kernel exception handling.
<b>KR2C (C28)</b>	A capability reserved for use during kernel exception handling.
<b>KCC (C29)</b>	Kernel code capability; the code capability moved to <b>PCC</b> when entering the kernel for exception handling.
<b>KDC (C30)</b>	Kernel data capability; the data capability containing the security domain for the kernel exception handler.
<b>EPCC (C31)</b>	Capability register associated with the exception program counter ( <b>EPC</b> ) required by exception handlers to save, interpret, and store the value of <b>PCC</b> at the time the exception fired.

Table 3.1: Capability registers defined by the capability coprocessor.

- **C3-C10** are used to pass arguments and are not preserved across calls. Capability returns are placed in **C3**.
- **C11-C16** are caller-save registers.
- **C17-C24** are callee-save registers.

When calling less-trusted code, there is no guarantee that a non-malicious callee will abide by these conventions. Thus, all registers should be regarded as caller-save. Additionally, all capability registers that are not part of the explicit argument set should be invalidated using the **CClearTag** instruction. This will prevent leakage of rights to untrustworthy callees. Where rights are explicitly passed to a callee, it may be desirable to clear the global bit that will (in a suitably configured runtime) prevent further propagation of the capability. Similar concerns apply to general-purpose registers, which should be preserved by the caller if their correct preservation is important, and cleared if they might leak sensitive data.

The CHERI version of Clang supports a second ABI for use in sandboxes. This ABI is intended to explore full use of capabilities, giving maximum protection at the expense of capabilities. In the Sandbox ABI, **C11** is used for the stack. Optionally, capability instructions are used for call and returns, with **CC12** used for calls and **Cc17** as the link register.

### Protected Procedure Calls

A protected procedure call, instruction **CCall**, escapes to a handler which takes a sealed executable (“code”) and sealed non-executable (“data”) capability with matching types. If the types match, the unsealed code capability is placed in **PCC** and the unsealed data capability is placed in **IDC**. The handler will also push the previous **PCC**, **IDC**, **PC + 4**, and **SP** to a stack pointed to by **TSC**. The stack pointer **TSC** may be implemented either as a hardware register or as a variable internal to a software implementation of **CCall**. The caller should invalidate all registers that are not intended to be passed to the callee before the call.

A protected procedure return, instruction **CReturn**, also escapes to a handler which pops the code and data capabilities from the stack at **TSC** and places them in **PCC** and **IDC** respectively; it likewise pops **PC** and **SP**. The callee should invalidate all registers that are not intended to be passed to the caller before the return.

The caller is responsible for ensuring that its protection domain is entirely embodied in the capability in **IDC** so that it can restore its state upon return.

These semantics are software defined, and we anticipate that different operating-system and programming-language security models might handle these, and other behaviors, in different ways. For example, in our prototype CheriBSD implementation, the operating-system kernel maintains a “trusted stack” onto which which values are pushed during invocation, and from which values are popped on return. Over time, we anticipate providing multiple sets of semantics, perhaps corresponding to less synchronous domain-transition models, and allowing different userspace runtimes to select (or implement) the specific semantics their programming model requires. This is particularly important in order to provide flexible error handling: if a sandbox suffers a fault, or exceeds its execution-time budget, it is the OS and programming language that will define how recovery takes place, rather than the ISA definition. Basic hardware acceleration of capability invocation and return is easy to envision regardless of the specific semantic: many of the checks performed against capability permissions and types will be shared by all of these systems.



## Capabilities and Exception Handling

**KCC** and **KDC** hold the code capability and data capability that describe the protection domain of the system exception handler. When an exception occurs, the victim **PCC** is copied to **EPCC** so that the exception may return to the correct address, and **KCC**, excepting its **offset** field, which will be set to the appropriate MIPS exception-vector address, is moved to **PCC** to grant execution rights for kernel code. When an exception handler returns with `eret`, **EPCC**, possibly after having been updated by the software exception handler, is moved into **PCC**. **KDC** may be manually installed by the exception handler if needed, and will typically be moved into **C0** in order to allow otherwise unmodified MIPS exception handlers to be used. This may also need to be restored before returning from the exception.

## 3.2 Capabilities

Each capability register contains the following fields:

- Tag bit (“**tag**”, 1 bit)
- Sealed bit (“**sealed**”, 1 bit)
- Permissions mask (“**perms**”, 31 bits)
- Object type (“**otype**”, 24 bits)
- Offset (“**offset**”, 64 bits)
- Base virtual address (“**base**”, 64 bits)
- Length in bytes (“**length**”, 64 bits)

A capability register can be in one of four possible states:

- It contains a valid capability. The **tag** bit is set, and all of the above fields contain defined values.
- It contains all of the fields of a capability (as defined above), but the capability is not valid for use (for example, because the program that set the values has not yet demonstrated that it is authorized to create a valid capability with those values). The **tag** bit is not set.
- It contains a 64-bit integer stored in the **offset** field. The **tag** bit is not set, and the **offset** field contains the integer.
- It contains 256 bits of non-capability data that have been loaded from memory. The **tag** bit is not set. Programs should not rely on there being any particular relationship between the bytes that have been loaded from memory and the **offset** field, although see below for the current in-memory representation of capabilities.

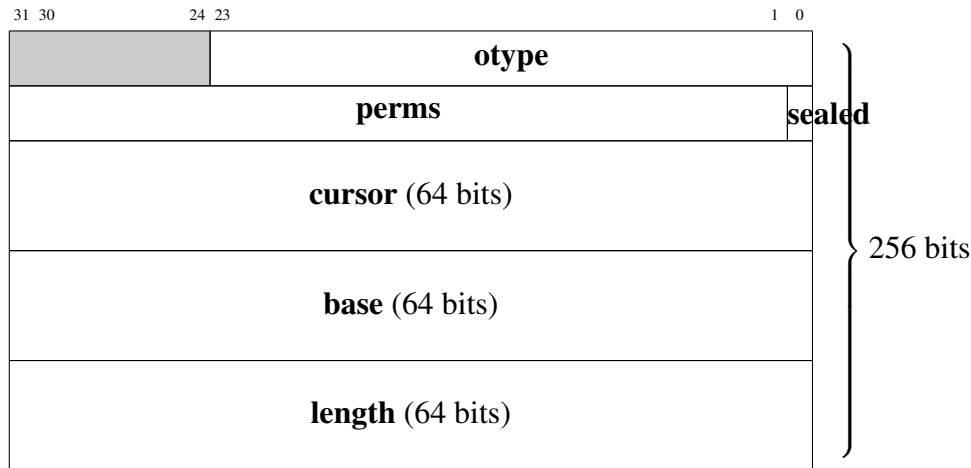


Figure 3.1: Memory representation of a capability

When capabilities are stored in memory, they are saved in the format shown in Figure 3.1. Programs should not rely on this memory representation, as it may change in future, and should instead access the fields through the instructions provided in the ISA. Note that there is a significant difference between the architecturally defined fields and the in-memory representation: **offset** is currently implemented as **cursor** – **base**, where the **cursor** field is internal to the implementation. These fields are stored in memory in a big-endian format. The CHERI processor is currently always defined to be big-endian, in contrast to the traditional MIPS ISA, which allows endianness to be selected by the supervisor. This is not fundamental to our approach; rather, it is expedient for early prototyping purposes.

### 3.2.1 tag

The **tag** bit indicates whether a capability register contains a capability or normal data. If **tag** is set, the register contains a capability. If **tag** is cleared, the rest of the register contains 256 bits of normal data.

### 3.2.2 sealed

The **sealed** flag indicates whether a capability is usable for general-purpose capability operations. If this flag is set, the capability is sealed and it may be used only by a `CCall` instruction. If the `CCall` instruction receives a sealed executable capability and a sealed non-executable capability with matching **otype** fields, both capabilities will have their **sealed** flag cleared and control flow branches to the code indicated by the code capability, thus entering a new security domain.

### 3.2.3 perms

The 31-bit **perms** bit vector governs the permissions of the capability including read, write and execute permissions. The contents of this field are listed in table 3.2. Bits 15–30 may

be used by application programs for user-defined permissions; they can be checked using the **CCheckPerm** instruction.

### 3.2.4 otype

This 24-bit field holds the “type” of a sealed capability; this field allows an unforgeable link to be created between associated data and object capabilities. While earlier versions of the CHERI ISA interpreted this field as an address, recent versions treat this as an opaque software-managed value.

### 3.2.5 offset

This 64-bit field holds an offset between the base of the capability and a byte of memory that is currently of interest to the program that created the capability. Effectively, it is a convenient place for a program to take an index into an array and store it inside a capability that grants access to the array. No bounds checks are performed on **offset** when its value is set by the **CSetOffset** instruction: programs are free to set its value beyond the end of the capability as defined by the **length** field. (Bounds checks are, however, performed when a program attempts to use the capability to access memory at the address given by **base** + **offset** + register offset.)

### 3.2.6 base

This 64-bit field is the base virtual address of the segment described by a capability.

### 3.2.7 length

This 64-bit field is the length of the segment described by a capability.

## 3.3 Capability Permissions

Table 3.2 shows constants currently defined for memory permissions; remaining bits are software-defined.

**Global** Allow this capability to persist beyond a protected procedure return.

**Permit\_Execute** Allow this capability to be used in the **PCC** register as a capability for the program counter.

**Permit\_Store\_Capability** Allow this capability to be used as a pointer for storing other capabilities.

**Permit\_Load\_Capability** Allow this capability to be used as a pointer for loading other capabilities.

**Permit\_Store** Allow this capability to be used as a pointer for storing data from general-purpose registers.

Value	Name
0	Global
1	Permit_Execute
2	Permit_Load
3	Permit_Store
4	Permit_Load_Capability
5	Permit_Store_Capability
6	Permit_Store_Local_Capability
7	Permit_Seal
8	Reserved
9	Reserved
10	Access_EPCC
11	Access_KDC
12	Access_KCC
13	Access_KR1C
14	Access_KR2C

Table 3.2: Memory permission bits for the **perms** capability field

**Permit\_Load** Allow this capability to be used as a pointer for loading data into general-purpose registers.

**Permit\_Store\_Local\_Capability** Allow this capability to be used as a pointer for storing local capabilities.

**Permit\_Seal** Allow this capability to be used to seal or unseal capabilities that have the same **otype**.

**Access\_EPCC** Allow access to **EPCC** when this capability is in **PCC**.

**Access\_KR1C** Allow access to **KR1C** when this capability is in **PCC**.

**Access\_KR2C** Allow access to **KR2C** when this capability is in **PCC**.

**Access\_KCC** Allow access to **KCC** when this capability is in **PCC**.

**Access\_KDC** Allow access to **KDC** when this capability is in **PCC**.

Local capabilities can be stored only via capabilities that have the **Permit\_Store\_Local\_Capability** permission bit set; normally, this permission will be set only on capabilities that, themselves, have the **Global** bit cleared.

### 3.4 Capability Exceptions

Many of the capability instructions can cause an exception (e.g., if the program attempts a load or a store that is not permitted by the capability system). When the cause of an exception is that

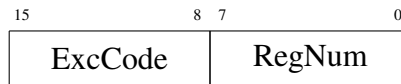


Figure 3.2: Capability Cause Register

the attempted operation is prohibited by the capability system, the *ExcCode* field within the **cause** register of coprocessor 0 are set to 18 (*C2E*, coprocessor 2 exception), **PCC** and **EPCC** are set as described in section 3.5 and **capcause** is set as described below.

### Capability Cause Register

The capability coprocessor has a **capcause** register that gives additional information on the reason for the exception. It is formatted as shown in figure 3.2. The possible values for the *ExcCode* of **capcause** are shown in table 3.4. If the last instruction to throw an exception did not throw a capability exception, then the *ExcCode* field of **capcause** will be *None*. *ExcCode* values from 128 to 255 are reserved for use by application programs. (A program can use *CSetCause* to set *ExcCode* to a user-defined value).

The *RegNum* field of **capcause** will hold the number of the capability register whose permission was violated in the last exception, if this register was not the unnumbered register **PCC**. If the capability exception was raised because **PCC** did not grant access to a numbered reserved register, then **capcause** will contain the number of the reserved register to which access was denied. If the exception was raised because **PCC** did not grant some other permission (e.g., permission to read **capcause** was required, but not granted) then *RegNum* will hold 0xff.

The **CGetCause** instruction can be used by an exception handler to read the **capcause** register. **CGetCause** will raise an exception if **PCC.perms.Access\_EPCC** is not set, so the operating system can prevent user space programs from reading **capcause** directly by not granting them *Access\_EPCC* permission.

### Exception Priority

If an instruction throws more than one capability exception, **capcause** is set to the highest priority exception (numerically lowest priority number) as shown in table 3.4. The *RegNum* field of **capcause** is set to the register which caused the highest priority exception.

All capability exceptions (*C2E*) have higher priority than address error exceptions (*AdEL*, *AdES*).

If an instruction throws more than one capability exception with the same priority (e.g., both the source and destination register are reserved registers), then the register which is furthest to the left in the assembly language opcode has priority for setting the *RegNum* field.

Some of these priority rules are security critical. In particular, an exception caused by a register being reserved must have priority over other capability exceptions (e.g., *AdEL* and *AdES*) to prevent a process from discovering information about the contents of a register that it is not allowed to access.

Other priority rules are not security critical, but are defined by this specification so that exception processing is deterministic.

An operating system might implement unaligned loads and stores by catching the *AdEL* and *AdES* exceptions and emulating the load or store. As capability exceptions have higher priority

Value	Description
0x00	None
0x01	Length Violation
0x02	Tag Violation
0x03	Seal Violation
0x04	Type Violation
0x05	Call Trap
0x06	Return Trap
0x07	Underflow of trusted system stack
0x08	User-defined Permission Violation
0x09	TLB prohibits store capability
0x0a	<i>reserved</i>
0x0b	<i>reserved</i>
0x0c	<i>reserved</i>
0x0d	<i>reserved</i>
0x0e	<i>reserved</i>
0x0f	<i>reserved</i>
0x10	Global Violation
0x11	Permit_Execute Violation
0x12	Permit_Load Violation
0x13	Permit_Store Violation
0x14	Permit_Load_Capability Violation
0x15	Permit_Store_Capability Violation
0x16	Permit_Store_Local_Capability Violation
0x17	Permit_Seal Violation
0x18	<i>reserved</i>
0x19	<i>reserved</i>
0x1a	Access_EPCC Violation
0x1b	Access_KDC Violation
0x1c	Access_KCC Violation
0x1d	Access_KR1C Violation
0x1e	Access_KR2C Violation
0x1f	<i>reserved</i>

Table 3.3: Capability Exception Codes

Priority	Description
1	Access_EPCC Violation Access_KDC Violation Access_KCC Violation Access_KR1C Violation Access_KR2C Violation
2	Tag Violation
3	Seal Violation
4	Type Violation
5	Permit_Seal Violation
6	Permit_Execute Violation
7	Permit_Load Violation Permit_Store Violation
8	Permit_Load_Capability Violation Permit_Store_Capability Violation
9	Permit_Store_Local_Capability Violation
10	Global Violation
11	Length Violation
12	User-defined Permission Violation
13	TLB prohibits store capability
14	Call Trap Return Trap

Table 3.4: Exception Priority

than alignment exceptions, this exception handler would not need to check the permissions (and base, length, etc.) of the capability before emulating the load/store, because it would be guaranteed that all capability checks had already been done by the hardware, and had passed.

### Exceptions and indirect addressing

If an exception is caused by the combination of the values of a capability register and a general purpose register (e.g., if an expression such as `c1b t1, t0(c0)` raises an exception because the offset `t0` is trying to read beyond `c0`'s length), the number of the capability register (not of the general-purpose register) will be stored in `capcause.RegNum`.

### Software Emulation of CCall and CReturn

In the current hardware implementation of CHERI, the CCall and CReturn instructions always raise an exception, so that the details of the call or return operation can be implemented in software by a trap handler. This exception uses a different trap handler vector, at 0x100 above the general purpose exception handler. The exception cause will be *C2E* and `capcause` will be *Call Trap* for CCall and *Return Trap* for CReturn.

## 3.5 CPU Reset

When the CPU is hard reset, all capability registers will be initialized to the following values:

- The **tag** bit is set.
- The **sealed** bit is unset.
- **offset** = 0 (except for **PCC.offset**, which will be initialized to the boot vector address)
- **base** = 0
- **length** =  $2^{64} - 1$
- **otype** = 0
- All permissions bits are set.
- All unused bits are cleared.

The initial values of **PCC** and **KCC** will allow the system to initially execute code relative to virtual address 0. The initial value of **C0** will allow general-purpose loads and stores to all of virtual memory for the bootstrapping process.

## 3.6 Changes to Standard MIPS ISA Processing

The following changes are made to the behavior of instructions from the standard MIPS ISA when a capability coprocessor is present:

**Instruction fetch** The MIPS-ISA program counter (**PC**) is extended to a full program-counter capability (**PCC**), which incorporates the historic **PC** as **PCC.offset**. Instruction fetch is controlled by the *Permit\_Execute* permission, as well as bounds checks, tag checks, and a requirement that the capability not be sealed. Failures will cause a coprocessor 2 exception (*C2E*) to be thrown. If an exception occurs during instruction fetch (e.g., AdEL, or a TLB miss) then *BadVAddr* is set equal to **PCC.base** + **PCC.offset**, providing the absolute virtual address rather than a **PCC**-relative virtual address to the supervisor, avoiding the need for capability awareness in TLB fault handling.

**Load and Store instructions** Standard MIPS load and store instructions are interposed on by the *default data capability*, **C0**. Addresses provided for load and store will be transformed and bounds checked by **C0.base**, **C0.offset**, and **C0.length**. **C0** must have the appropriate permission (*Permit\_Store* or *Permit\_Load*) set, the full range of addresses covered by the load or store must be in range, **C0.tag** must be set, and **C0.sealed** must not be set. Failures will cause a coprocessor 2 exception (*C2E*) to be thrown. As with instruction fetch, *BadVAddr* values provided to the supervisor will be absolute virtual addresses, avoiding the need for capability awareness in TLB fault handling.



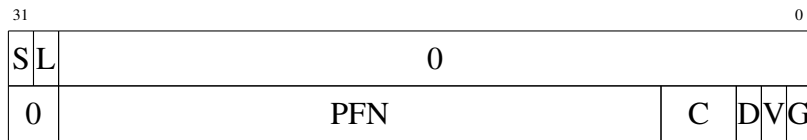


Table 3.5: EntryLo Register

**Floating-point Load and Store instructions** If the CPU is configured with a floating-point unit, all loads and stores between the floating-point unit and memory are also relative to **C0.base** and **C0.offset**, and are checked against the permissions, bounds, tag, and sealed state of **C0**.

**Jump and link register** After a **jair** instruction, the return address is relative to **PCC.base**.

**Exceptions** The MIPS exception program counter (**EPC**) is extended to a full exception program-counter capability (**EPCC**), which incorporates the historic **EPC** as **EPCC.offset**. Once **PCC** has been preserved, the contents of the *kernel code capability* (**KCC**), excluding **KCC.offset**, are moved into **PCC**. **EPCC.offset (EPC)** will be set to the exception vector address normally used by MIPS. This allows the exception handler to run with the permissions granted by **KCC**, which may be greater than the permissions granted by **PCC** before the exception occurred.

On return from an exception (`eret`), **PCC** is restored from **EPCC**, which will include **EPCC.offset** (also visible as **EPC**)<sup>1</sup> This allows exception handlers that are not aware of capabilities to continue to work on a CPU with the **CHERI** extensions. The result of `eret` is **UNDEFINED** if **EPCC.tag** is not set or **EPCC.sealed** is set. Similarly, the result of an exception or interrupt is **UNDEFINED** if **KCC.tag** is not set or **KCC.sealed** is set.

### 3.7 Changes to the TLB

**CHERI** adds two new fields to the **EntryLo** register, shown as **L** and **S** in Table 3.5.

If **L** is set, capability loads are disabled for the page. If a `CLC` instruction is used on a page with the **L** bit set, and the load succeeds, the value loaded into the destination register will have its tag bit cleared, even if the tag bit was set in memory.

If **S** is set, capability stores are disabled for the page. If a `CSC` instruction is used on a page with the **S** bit set, and the capability register to be stored has the tag bit set, then a **CP2** exception will be raised, with the **CP2** cause register set to `0x9` (**TLB** prohibits store capability). If the capability register to be stored does not have the tag bit set (i.e. it contains non-capability data), then this exception will not be raised and the store will proceed.

<sup>1</sup>In our current **CHERI1** prototype, for reasons of critical path length, **EPCC.offset** will be updated to be the value of the MIPS **EPC** on exception entry, but writes to **EPCC.offset** will not be propagated to **EPCC.offset** on exception return. As described later in this chapter, we have proposed shifting **EPCC** out of the ordering capability register file and instead using special registers in order to eliminate this problem. Our **L3** formal model of the **CHERI** ISA implements the specified behavior. **CheriBSD** is able to operate on both as it is careful to write to update both **EPC** and **EPCC** before returning.

At with other TLB-related exceptions, *BadVAddr* will be set to the absolute virtual address that has triggered the fault.

## 3.8 Possible Extensions to the CHERI ISA

The following changes have been discussed and are targeted for short-term implementation in the CHERI ISA:

- Move MIPS memory-access interposition from **C0** to a control register accessed via special instructions, as is the case for **PCC**; similarly migrate exception-handling registers such as **EPCC** to control registers.
- At present, capability load and store instructions can handle both untagged data and capabilities, supporting capability-oblivious operations. A related idea would be to have tag-stripping loads that would avoid triggering permission limitations on loading tagged data. This might be done via load instructions that explicitly request the clearing of tags to be used in efficient implementations of `memcpy`.
- Allow **CCall/CReturn** instructions to partially implement capability call/return semantics while still throwing an exception to allow software processing. For example, they might implement all necessary checks so that software does not need to do this.

The following changes have been discussed for longer-term consideration:

- Allow **CReturn** to accept code/data capability arguments, which might be ignored for the time being.
- 32-byte capabilities impose measurable overhead; implementing a 16-byte “compressed capability” representation, usable for pure data capabilities, might reduce this overhead. In order for this to be useful, the compiler must be able to statically determine and control use of compressed capabilities. Initial simulations of 128-bit capabilities suggest a substantial reduction in cache footprint.
- Consider further the effects of combining general-purpose and capability register files, which would avoid adding a new register file, but make some forms of ABI compatibility more challenging.

# Chapter 4

## CHERI Instruction Reference

CHERI's instructions express a variety of operations affecting capability and general-purpose registers as well as memory access and control flow. A key design concern is *guarded manipulation*, a set of constraints across all instructions that ensure monotonic non-increase in rights through capability manipulations. These instructions also assume, and specify, the presence of *tagged memory*, described in the previous chapter, which protects in-memory representations of capability values. Many further behaviors, such as reset state and exception handling, taken for granted in these instruction descriptions, are also described in the previous chapter.

The instructions fall into a number of categories: instructions to copy fields from capability registers into general-purpose registers so that they can be computed on, instructions for refining fields within capabilities, instructions for memory access via capabilities, instructions for jumps via capabilities, instructions for sealing capabilities, and instructions for capability invocation. In this chapter, we specify each instruction via both informal descriptions and pseudocode. To allow for more succinct pseudocode descriptions we rely on a number of common definitions also described in this chapter.

### 4.1 Common Definitions

The following functions are used in the pseudocode for more than one instruction, and are collected here for convenience. The & notation means use the number of a register, rather than its contents.

**function** REGISTER\_INACCESSIBLE(cb)

**return**

(&cb = EPCC **and not** PCC.access\_epcc) **or**

(&cb = KDC **and not** PCC.access\_kdc) **or**

(&cb = KCC **and not** PCC.access\_kcc) **or**

(&cb = KR1C **and not** PCC.access\_kr1c) **or**

(&cb = KR2C **and not** PCC.access\_kr2c)

**end function**

to\_signed64 converts an unsigned 64 bit integer into a signed 64 bit integer:

**function** TO\_SIGNED64(x)

**if**  $x < 2^{32}$  **then**

```

    return x
else
    return  $x - 2^{64}$ 
end if
end function

```

`raise_c2_exception` is used when an instruction raises an exception. The following pseudocode omits the details of MIPS exception handling (branching to the address of the exception handler, etc.)

```

procedure RAISE_C2_EXCEPTION(cause, reg)
    cp0.cause ← 18
    capcause.reg ← &reg
    capcause.cause ← cause
    ...
end procedure

procedure RAISE_C2_EXCEPTION_V(reg)
    if &reg = EPCC then
        raise_c2_exception(exceptionAccessEPCC, reg)
    else if &reg = KDC then
        raise_c2_exception(exceptionAccessKDC, reg)
    else if &reg = KCC then
        raise_c2_exception(exceptionAccessKCC, reg)
    else if &reg = KR1C then
        raise_c2_exception(exceptionAccessKR1C, reg)
    else if &reg = KR2C then
        raise_c2_exception(exceptionAccessKR2C, reg)
    end if
end procedure

procedure RAISE_C2_EXCEPTION_NOREG(cause)
    cp0.cause ← 18
    capcause.reg ← 0xff
    capcause.cause ← cause
    ...
end procedure

```

`execute_branch` is used when an instruction branches to an address. The MIPS ISA includes branch delay slots, so the instruction in the branch delay slot will be executed before the branch is taken; this is omitted in the following pseudocode:

```

procedure EXECUTE_BRANCH(pc)
    ...
end procedure

```

`execute_branch_pcc` is used when an instruction branches to an address and changes PCC. The change to PCC does not take effect until after the instruction in the branch delay slot has been executed.

```

procedure EXECUTE_BRANCH_PCC(pc)
    ...

```

**end procedure**

## **4.2 Details of Individual Instructions**

The following sections provide a detailed description of each CHERI ISA instructions. Each instruction description includes the following information:

- Instruction opcode format number
- Assembly language syntax
- Bitwise figure of the instruction layout
- Text description of the instruction
- Pseudo-code description of the instruction
- Enumeration of any exceptions that the instruction can trigger

Table 4.1 lists available capability coprocessor instructions.

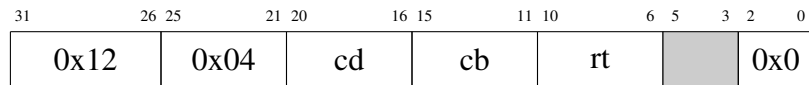
Mnemonic	Description
CGetBase	Move base to a general-purpose register
CGetOffset	Move offset to a general-purpose register
CGetLen	Move length to a general-purpose register
CGetTag	Move tag bit to a general-purpose register
CGetSealed	Move sealed bit to a general-purpose register
CGetPerm	Move permissions field to a general-purpose register
CGetType	Move object type field to a general-purpose register
CToPtr	Capability to pointer
CPtrCmp	Compare capability pointers
CIncBase	Increase Base
CIncOffset	Increase Offset
CMove	Pseudo-instruction for CIncBase with no change to the base
CSetLen	Set Length
CClearTag	Clear the tag bit
CAndPerm	Restrict Permissions
CSetOffset	Set cursor to an offset from base
CGetPCC	Move PCC to capability register
CFromPtr	Create capability from pointer
CSC	Store Capability Register
CLC	Load Capability Register
CL[BHWD][U]	Load Byte, Half-Word, Word or Double Via Capability Register (Unsigned)
CS[BHWD]	Store Byte, Half-Word, Word or Double Via Capability Register
CLLD	Load linked doubleword via capability register
CSCD	Store conditional doubleword via capability register
CBTU	Branch if capability tag is unset
CBTS	Branch if capability tag is set
CJR	Jump Capability Register
CJALR	Jump and link Capability Register
CCheckPerm	Check <b>perms</b> field
CCheckType	Check <b>otype</b> field
CSeal	Seal an unsealed capability
CUnseal	Unseal a sealed capability
CCall	Protected procedure call into a new security domain
CReturn	Return to the previous security domain
CGetCause	Move capability exception cause register to a general-purpose register
CSetCause	Set the capability exception cause register

Figure 4.1: Capability coprocessor instruction summary

## CAndPerm: Restrict Permissions

### Format

CAndPerm *cd*, *cb*, *rt*



### Description

Capability register *cd* is replaced with the contents of capability register *cb* with the **perms** field set to the bitwise AND of its previous value and the contents of general-purpose register *rt*.

### Pseudocode

```
if register_inaccessible(cd) then
    raise_c2_exception_v(cd)
else if register_inaccessible(cb) then
    raise_c2_exception_v(cb)
else if not cb.tag then
    raise_c2_exception(exceptionTag, cb)
else if cb.sealed then
    raise_c2_exception(exceptionSealed, cb)
else
    cd ← cb
    cd.perms ← cb.perms ∩ rt[0 .. 30]
end if
```

### Exceptions

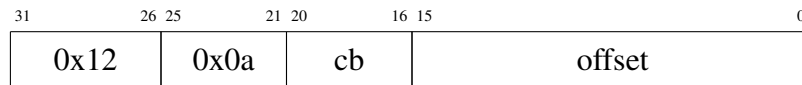
A coprocessor 2 exception is raised if:

- *cb* or *cd* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.
- *cb.tag* is not set.
- *cb.sealed* is set.

## CBTS: Branch if tag is set

### Format

CBTS *cb*, *offset*



### Description

Sets the **PC** to  $\text{PC} + 4 * \text{offset} + 4$ , where *offset* is sign extended, if *cb.tag* is set.

The instruction following the branch, in the delay slot, is executed before branching.

### Pseudocode

```
if register_inaccessible(cb) then
    raise_c2_exception_v(cb)
else if cb.tag then
    if  $\text{PC} + \text{offset} + 4 > \text{PCC.length}$  then
        raise_c2_exception_noreg(exceptionLength)
    else
        execute_branch( $\text{PC} + 4 * \text{sign\_extend}(\text{offset}) + 4$ )
    end if
end if
```

### Exceptions

A coprocessor 2 exception is raised if:

- *cb* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.
- $\text{PC} + 4 * \text{offset} + 4$  is greater than **PCC.length** and *cb.tag* is set.

### Notes

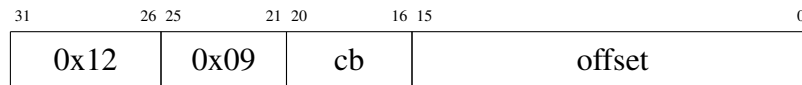
1. Like all MIPS branch instructions, CBTS has a branch delay slot. The instruction after it will always be executed, regardless of whether the branch is taken or not.



## CBTU: Branch if tag is unset

### Format

CBTU *cb*, *offset*



### Description

Sets the **PC** to  $\text{PC} + 4 * \text{offset} + 4$ , where *offset* is sign extended, if *cb.tag* is not set.  
The instruction following the branch, in the delay slot, is executed before branching.

### Pseudocode

```
if register_inaccessible(cb) then
    raise_c2_exception_v(cb)
else if not cb.tag then
    if  $\text{PC} + 4 * \text{sign\_extend}(\text{offset}) + 4 > \text{PCC.length}$  then
        raise_c2_exception_noreg(exceptionLength)
    else
        execute_branch( $\text{PC} + 4 * \text{sign\_extend}(\text{offset}) + 4$ )
    end if
end if
```

### Exceptions

A coprocessor 2 exception is raised if:

- *cb* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.
- $\text{PC} + 4 * \text{sign\_extend}(\text{offset}) + 4$  is greater than **PCC.length** and *cb.tag* is not set.

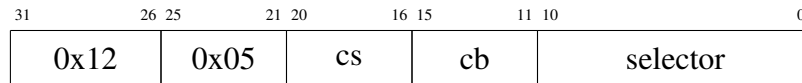
### Notes

1. Like all MIPS branch instructions, CBTU has a branch delay slot. The instruction after it will always be executed, regardless of whether the branch is taken or not.
2. This instruction is intended to resemble the conditional branch instructions from the MIPS ISA. In particular, the shift left of the offset by 2 bits and adding 4 is the same as MIPS conditional branches.

## CCall: Call into a new security domain

### Format

CCall *cs*, *cb*[, *selector*]



### Description

CCall is used to make a call into a protected subsystem (which may have access to a different set of capabilities than its caller). *cs* contains a code capability for the subsystem to be called, and *cb* contains a sealed data capability which will be unsealed for use by the called subsystem. In terms of object-oriented programming, *cb* is a capability for an *object* and *cs* is a capability for the methods of the object's class.

This instruction can be implemented in a number of ways split over hardware and software; we have experimented with several. A simple implementation might have CCall throw a software exception, with all other behavior implemented via a software trap handler. A hybrid implementation could perform various checks in hardware, deferring only trusted stack manipulation (or other behaviors, such as asynchronous calling conventions) to the software trap handler. Further defensive coding conventions (beyond instruction semantics) may also sensibly be shifted to the exception handler in order to avoid redundancy – e.g., the clearing of the same registers to prevent leaks in either direction. A significant tension exists in the hardware optimization of this instruction between using a flexible calling convention and semantics versus exploiting hardware optimization opportunities. Authors of compilers or assembly language programs should not rely on CCall being implemented in any particular blend of hardware and software.

1. **PCC** (with its **offset** field set to the program counter (**PC**) + 4) is pushed onto the trusted system stack.
2. **IDC** is pushed onto the trusted system stack.
3. *cs* is unsealed and the result placed in **PCC**.
4. *cb* is unsealed and the result placed in **IDC**.
5. The program counter is set to *cs.offset*. (i.e. control branches to virtual address *cs.base* + *cs.offset*, because the program counter is relative to **PCC.base**).

The *selector* field can be omitted in assembly and, in the current version of the ISA, must be 0 if explicitly provided. Issuing a CCall instruction with any other value is undefined behavior.

### Pseudocode (hardware)

```
raise_c2_exception(exceptionCall, cs)
```

## Pseudocode (software)

```
if register_inaccessible(cs) then
    raise_c2_exception_v(cs)
else if register_inaccessible(cb) then
    raise_c2_exception_v(cb)
else if not cs.tag then
    raise_c2_exception(exceptionTag, cs)
else if not cb.tag then
    raise_c2_exception(exceptionTag, cb)
else if not cs.sealed then
    raise_c2_exception(exceptionSealed, cs)
else if not cb.sealed then
    raise_c2_exception(exceptionSealed, cb)
else if cs.otype  $\neq$  cb.otype then
    raise_c2_exception(exceptionType, cs)
else if not cs.perms.Permit_Execute then
    raise_c2_exception(exceptionPermitExecute, cs)
else if not cb.perms.Permit_Execute then
    raise_c2_exception(exceptionPermitExecute, cb)
else if cs.offset  $\geq$  cs.length then
    raise_c2_exception(exceptionLength, cs)
else
    PCC.offset  $\leftarrow$  PC + 4
    TSS  $\leftarrow$  TSS - 32
    mem[TSS .. TSS + 31]  $\leftarrow$  PCC
    tags[toTag(TSS)]  $\leftarrow$  PCC.tag
    TSS  $\leftarrow$  TSS - 32
    mem[TSS .. TSS + 31]  $\leftarrow$  IDC
    tags[toTag(TSS)]  $\leftarrow$  IDC.tag
    PCC  $\leftarrow$  cs
    PCC.sealed  $\leftarrow$  false
    IDC  $\leftarrow$  cb
    IDC.sealed  $\leftarrow$  false
    PC  $\leftarrow$  cs.offset
end if
```

## Exceptions

A coprocessor 2 exception will be raised so that the desired semantics can be implemented in a trap handler.

The capability exception code will be 0x05 and the handler vector will be 0x100 above the general purpose exception handler.

A further coprocessor 2 exception raised if:

- *cs* or *cb* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.

- *cs.sealed* is not set.
- *cb.sealed* is not set.
- *cs.otype*  $\neq$  *cb.otype*
- *cs.perms.Permit\_Execute* is not set.
- *cb.perms.Permit\_Execute* is set.
- *cs.offset*  $\geq$  *cs.length*.
- The trusted system stack would overflow (i.e., if **PCC** and **IDC** were pushed onto the system stack, it would overflow the bounds of **TSC**).

## Notes

- From the point of view of security, `CCall` needs to be an atomic operation (i.e. the caller cannot decide to just do some of it, because partial execution could put the system into an insecure state). From the point of view of hardware design, `CCall` needs to write two capabilities to memory, which might take more than one clock cycle. One possible way to satisfy both of these constraints is to make `CCall` cause a software trap, and the trap handler uses its access to **KCC** and **KDC** to implement `CCall`.
- Implementations may choose to restrict the register numbers that may be passed as *cs* and *cb* in order to avoid the need to decode the instruction and identify the register arguments. The software implementation in CheriBSD at the time of writing requires the *cs* be **C1**, and that *cb* be **C2**, consistent with the CHERI ABI.
- In the current implementation, where `CCall` is implemented by the hardware raising an exception, `CCall` does not have a branch delay slot.
- The 10 bit selector in the `CCall` instruction is reserved for future iterations to provide different domain transition semantics, for example an asynchronous message send rather than a synchronous call.

## CCheckPerm: Raise exception if don't have permission

### Format

CCheckPerm cs, rt



### Description

A exception is raised (and the capability cause set to “user defined permission violation”) if there is a bit set in *rt* which is not set in *cs.perms* (i.e. *rt* describes a set of permissions, and an exception is raised if *cs* does not grant all of those permissions).

### Pseudocode

```
if register_inaccessible(cs) then
    raise_c2_exception_v(cs)
else if not cs.tag then
    raise_c2_exception(exceptionTag, cs)
else if cs.perms ∩ rt ≠ rt then
    raise_c2_exception(exceptionUserDefined, cs)
end if
```

### Exceptions

A coprocessor 2 exception is raised if:

- *cs* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.
- *cs.tag* is not set.
- There is a bit which is set in *rt* and is not set in *cs.perms*.

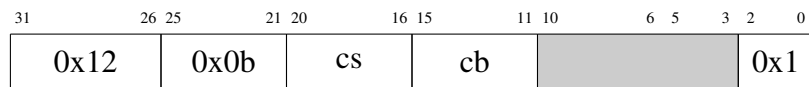
### Notes

- If *cs.tag* is not set, then *cs* does not contain a capability, *cs.perms* might not be meaningful as a permissions field, and so a *tagViolation* exception is raised.
- This instruction can be used to check the permissions field of a sealed capability, so the instruction does not check *cs.sealed*.

## CCheckType: Raise exception if otypes don't match

### Format

CCheckType cs, cb



### Description

An exception is raised if *cs.otype* is not equal to *cb.otype*.

### Pseudocode

```
if register_inaccessible(cs) then
    raise_c2_exception_v(cs)
else if register_inaccessible(cb) then
    raise_c2_exception_v(cb)
else if not cs.tag then
    raise_c2_exception(exceptionTag, cs)
else if not cb.tag then
    raise_c2_exception(exceptionTag, cb)
else if not cs.sealed then
    raise_c2_exception(exceptionSealed, cs)
else if not cb.sealed then
    raise_c2_exception(exceptionSealed, cb)
else if cs.otype ≠ cb.otype then
    raise_c2_exception(exceptionType, cs)
end if
```

### Exceptions

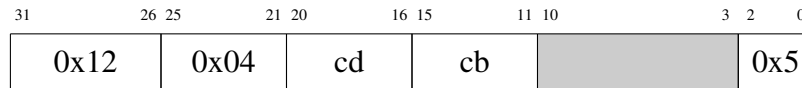
A coprocessor 2 exception is raised if:

- *cs* or *cb* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.
- *cs.tag* is not set.
- *cb.tag* is not set.
- *cs.sealed* is not set.
- *cb.sealed* is not set.
- *cs.otype* ≠ *cb.otype*.

## CClearTag: Clear the tag bit

### Format

CClearTag *cd*, *cb*



### Description

Capability register *cd* is replaced with the contents of *cb*, with the tag bit cleared.

### Pseudocode

```
if register_inaccessible(cd) then  
    raise_c2_exception_v(cd)  
else if register_inaccessible(cb) then  
    raise_c2_exception_v(cb)  
else  
    cd ← cb  
    cd.tag ← false  
end if
```

### Exceptions

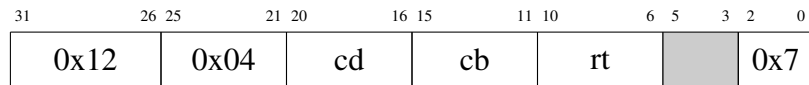
A coprocessor 2 exception is raised if:

- *cb* or *cd* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.

## CFromPtr: Create capability from pointer

### Format

CFromPtr *cd*, *cb*, *rt*



### Description

*rt* is a pointer using the C-language convention that a zero value represents the NULL pointer. If *rt* is zero, then *cd* will be the NULL capability (tag bit not set, all other bits also not set). If *rt* is non-zero, then *cd* will be a capability whose base is *cb.base* + *rt*.

### Pseudocode

```
if register_inaccessible(cd) then
    raise_c2_exception_v(cd)
else if register_inaccessible(cb) then
    raise_c2_exception_v(cb)
else if rt = 0 then
    cd.tag ← false
    cd.sealed ← false
    cd.perms ← ∅
    cd.base ← 0
    cd.length ← 0
    cd.offset ← 0
    cd.otype ← 0
    cd.reserved ← 0
else if not cb.tag then
    raise_c2_exception(exceptionTag, cb)
else if cb.sealed then
    raise_c2_exception(exceptionSealed, cb)
else if rt > cb.length then
    raise_c2_exception(exceptionLength, cb)
else
    cd ← cb
    cd.base ← cb.base + rt
    cd.length ← cb.length - rt
end if
```

### Exceptions

A coprocessor 2 exception is raised if:

- *cb* or *cd* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.

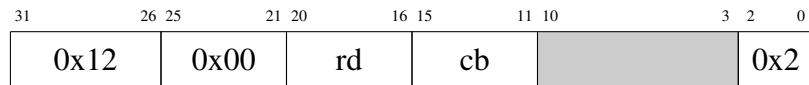


- *cb.tag* is not set.
- *cb.sealed* is set.
- *rt > cb.length*.

## CGetBase: Move Base to a General-Purpose Register

### Format

CGetBase *rd*, *cb*



### Description

General-purpose register *rd* is set equal to the **base** field of capability register *cb*.

### Pseudocode

```
if register_inaccessible(cb) then
    raise_c2_exception_v(cb)
else
    rd ← cb.base
end if
```

### Exceptions

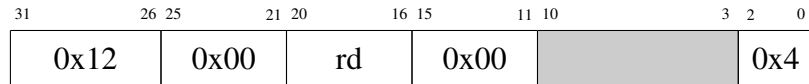
A coprocessor 2 exception is raised if:

- *cb* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.

## CGetCause: Move the Capability Exception Cause Register to a General-Purpose Register

### Format

CGetCause rd



### Description

General-purpose register *rd* is set equal to the capability cause register.

### Pseudocode

```
if not PCC.perms.Access_EPCC then  
    raise_c2_exception_noreg(exceptionAccessEPCC)  
else  
    rd[0 .. 7] ← capcause.reg  
    rd[8 .. 15] ← capcause.cause  
    rd[16 .. 63] ← 0  
end if
```

### Exceptions

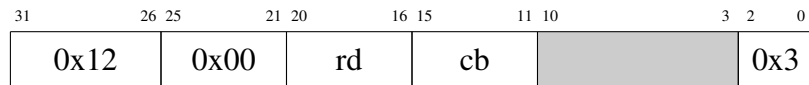
A coprocessor 2 exception is raised if:

- **PCC.perms**.Access\_EPCC is not set.

## CGetLen: Move Length to a General-Purpose Register

### Format

CGetLen rd, cb



### Description

General-purpose register *rd* is set equal to the **length** field of capability register *cb*.

### Pseudocode

```
if register_inaccessible(cb) then
    raise_c2_exception_v(cb)
else
    rd ← cb.length
end if
```

### Exceptions

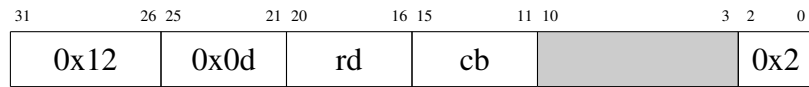
A coprocessor 2 exception is raised if:

- *cb* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.

## CGetOffset: Move Offset to a General-Purpose Register

### Format

CGetOffset rd, cb



### Description

General-purpose register *rd* is set equal to the **offset** fields of capability register *cb*.

### Pseudocode

```
if register_inaccessible(cb) then
    raise_c2_exception_v(cb)
else
    rd ← cb.offset
end if
```

### Exceptions

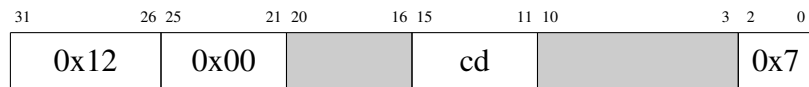
A coprocessor 2 exception is raised if:

- *cb* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.

## CGetPCC: Move PCC to capability register

### Format

CGetPCC *cd*



### Description

Capability register *cd* is set equal to the **PCC**, with *cd.offset* set equal to **PC**.

### Pseudocode

```
if register_inaccessible(cd) then
    raise_c2_exception_v(cd)
else
    cd ← PCC
    cd.offset ← PC
end if
```

### Exceptions

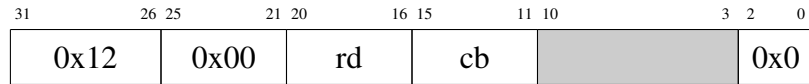
A coprocessor 2 exception is raised if:

- *cd* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.

## CGetPerm: Move Memory Permissions Field to a General-Purpose Register

### Format

CGetPerm rd, cb



### Description

The least significant 31 bits (bits 0 to 30) of general-purpose register *rd* are set equal to the **perms** field of capability register *cb*. The other bits of *rd* are set to zero.

### Pseudocode

```
if register_inaccessible(cb) then  
    raise_c2_exception_(cb)  
else  
    rd[0:30] ← cb.perms  
    rd[31:63] ← 0  
end if
```

### Exceptions

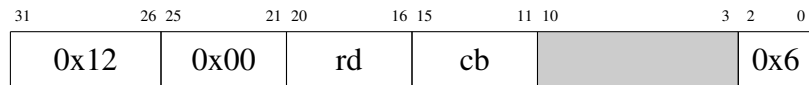
A coprocessor 2 exception is raised if:

- *cb* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.

## CGetSealed: Move sealed bit to a General-Purpose Register

### Format

CGetSealed rd, cb



### Description

The low-order bit of *rd* is set to *cb.sealed*. All other bits of *rd* are cleared.

### Pseudocode

```
if register_inaccessible(cb) then
    raise_c2_exception_v(cb)
else
    rd[0] ← cb.sealed
    rd[1:63] ← 0
end if
```

### Exceptions

A coprocessor 2 exception is raised if:

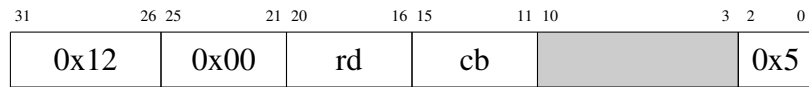
- *cb* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.



## CGetTag: Move Tag to a General-Purpose Register

### Format

CGetTag rd, cb



### Description

The low bit of *rd* is set to the tag value of *cb*. All other bits are cleared.

### Pseudocode

```
if register_inaccessible(cb) then  
    raise_c2_exception_v(cb)  
else  
    rd[0] ← cb.tag  
    rd[1:63] ← 0  
end if
```

### Exceptions

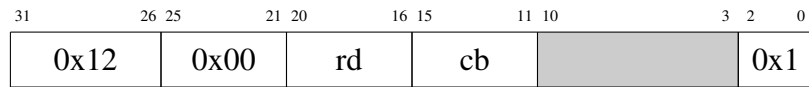
A coprocessor 2 exception is raised if:

- *cb* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.

## CGetType: Move Object Type Field to a General-Purpose Register

### Format

CGetType *rd*, *cb*



### Description

General-purpose register *rd* is set equal to the **otype** field of capability register *cb*.

### Pseudocode

```
if register_inaccessible(cb) then  
    raise_c2_exception_v(cb)  
else  
    rd[0:23] ← cb.otype  
    rd[24:63] ← 0  
end if
```

### Exceptions

A coprocessor 2 exception is raised if:

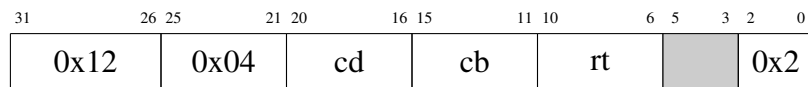
- *cb* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.

## CIncBase: Increase Base

### Format

CIncBase *cd*, *cb*, *rt*

CMove *cd*, *cb*



### Description

Capability register *cd* is replaced with the contents of capability register *cb* with the **base** field set to the sum of its previous value and the contents of general-purpose register *rt*. The **length** field of capability register *cd* is replaced with *cb.length* minus the contents of general-purpose register *rt*, ensuring that capability register *cd* points to a subset of the original memory region.

### Pseudocode

```
if register_inaccessible(cd) then
    raise_c2_exception_v(cd)
else if register_inaccessible(cb) then
    raise_c2_exception_v(cb)
else if not cb.tag and rt ≠ 0 then
    raise_c2_exception(exceptionTag, cb)
else if cb.sealed and rt ≠ 0 then
    raise_c2_exception(exceptionSealed, cb)
else if rt > cb.length then
    raise_c2_exception(exceptionLength, cb)
else
    cd ← cb
    cd.base ← cb.base + rt
    cd.length ← cb.length - rt
end if
```

### Exceptions

A coprocessor 2 exception is raised if:

- *cb* or *cd* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.
- *cb.tag* is not set and *rt* ≠ 0.
- *cb.sealed* is set and *rt* ≠ 0.
- *rt* > *cb.length*

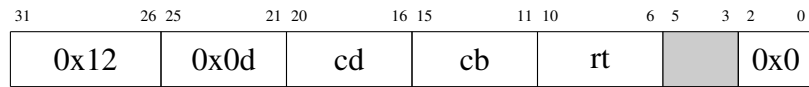
## Notes

- `CIncBase` can be used to copy one register to another by setting `rt` equal to zero. If `rt` is zero, the operation will succeed even if `cb.sealed` is set, allowing it to be used to copy sealed capabilities. `CIncBase` also succeeds if `rt` is zero and `cb.tag` is unset, allowing it to be used to copy non-capability data items between capability registers.
- In assembly language, `CMove cd, cb` is a pseudo-instruction which the assembler converts to `CIncBase cd, cb, $zero`.
- Assembly programmers must exercise caution when using `CIncBase` or `CSetLen` on capabilities such as C pointers that have non-zero **offset** values. Incorrect operation orders may product empty capabilities.

## CIncOffset: Increase Offset

### Format

CIncOffset *cd*, *cb*, *rt*



### Description

Capability register *cd* is set equal to capability register *cb* with its **offset** field replaced with *cb.offset* + *rt*.

### Pseudocode

```
if register_inaccessible(cd) then  
    raise_c2_exception_v(cd)  
else if register_inaccessible(cb) then  
    raise_c2_exception_v(cb)  
else if cb.tag and cb.sealed then  
    raise_c2_exception(cb, exceptionSealed)  
else  
    cd ← cb  
    cd.offset ← (cb.offset + rt) mod 264  
end if
```

### Exceptions

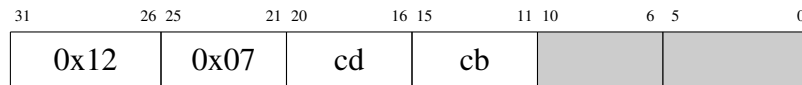
A coprocessor 2 exception is raised if:

- *cb* or *cd* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.
- *cb.tag* and *cb.sealed* are both set.

## CJALR: Jump and Link Capability Register

### Format

CJALR *cd*, *cb*



### Description

The current **PCC** (with an offset of the current **PC** + 8) is saved in *cd*. **PCC** is then loaded from capability register *cb* and **PC** is set from its offset.

### Pseudocode

```
if register_inaccessible(cd) then
    raise_c2_exception_v(cd)
else if register_inaccessible(cb) then
    raise_c2_exception_v(cb)
else if not cb.tag then
    raise_c2_exception(exceptionTag, cb)
else if cb.sealed then
    raise_c2_exception(exceptionSealed, cb)
else if not cb.perms.Permit_Execute then
    raise_c2_exception(exceptionPermitExecute, cb)
else if not cb.perms.Global then
    raise_c2_exception(exceptionGlobal, cb)
else if cb.offset + 4 > cb.length then
    raise_c2_exception(exceptionLength, cb)
else if align_of(cb.base + cb.offset) < 4 then
    raise_exception(exceptionAdEL)
else
    cd ← PCC
    cd.offset ← PC + 8
    execute_branch_pcc(cb.offset, cb)
end if
```

### Exceptions

A coprocessor 2 exception will be raised if:

- *cd* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.
- *cb* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.
- *cb*.perms.Permit\_Execute is not set.

- ***cb.offset*** + 4 is greater than ***cb.length***.
- ***cb.sealed*** is set.
- ***cb.tag*** is not set.
- ***cb.perms.Global*** is not set.

An address error exception will be raised if

- ***cb.base*** + ***cb.offset*** is not 4-byte word aligned.

## Notes

- `cjalr` has a branch delay slot.
- The change to **PCC** does not take effect until the instruction in the branch delay slot has been executed.

## CJR: Jump Capability Register

### Format

CJR *cb*



### Description

**PCC** is loaded from *cb*, and **PC** is loaded from *cb.offset*.

### Pseudocode

```
if register_inaccessible(cb) then
    raise_c2_exception_v(cb)
else if not cb.tag then
    raise_c2_exception(exceptionTag, cb)
else if cb.sealed then
    raise_c2_exception(exceptionSealed, cb)
else if not cb.perms.Permit_Execute then
    raise_c2_exception(exceptionPermitExecute, cb)
else if not cb.perms.Global then
    raise_c2_exception(exceptionGlobal, cb)
else if cb.offset + 4 > cb.length then
    raise_c2_exception(exceptionLength, cb)
else if align_of(cb.base + cb.offset) < 4 then
    raise_exception(exceptionAdEL)
else
    execute_branch_pcc(cb.offset, cb)
end if
```

### Exceptions

A coprocessor 2 exception is raised if:

- *cb* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.
- *cb.tag* is not set.
- *cb.sealed* is set.
- *cb.perms.Permit\_Execute* is not set.
- *cb.perms.Global* is not set.
- Register *cb.offset* + 4 is greater than *cb.length*.



An address error exception is raised if:

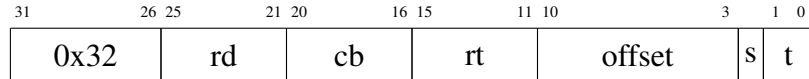
- *cb.base* + *cb.offset* is not 4-byte word aligned.

*cb.base* and *cb.length* are treated as unsigned integers, and the result of the addition does not wrap around (i.e., an exception is raised if *cb.base+cb.offset* is greater than `maxaddr`).

## Load Via Capability Register

### Format

CLB *rd*, *rt*, *offset(cb)*  
CLH *rd*, *rt*, *offset(cb)*  
CLW *rd*, *rt*, *offset(cb)*  
CLD *rd*, *rt*, *offset(cb)*  
CLBU *rd*, *rt*, *offset(cb)*  
CLHU *rd*, *rt*, *offset(cb)*  
CLWU *rd*, *rt*, *offset(cb)*  
CLBR *rd*, *rt(cb)*  
CLHR *rd*, *rt(cb)*  
CLWR *rd*, *rt(cb)*  
CLDR *rd*, *rt(cb)*  
CLBUR *rd*, *rt(cb)*  
CLHUR *rd*, *rt(cb)*  
CLWUR *rd*, *rt(cb)*  
CLBI *rd*, *offset(cb)*  
CLHI *rd*, *offset(cb)*  
CLWI *rd*, *offset(cb)*  
CLDI *rd*, *offset(cb)*  
CLBUI *rd*, *offset(cb)*  
CLHUI *rd*, *offset(cb)*  
CLWUI *rd*, *offset(cb)*



### Purpose

Loads a data value via a capability register, and extends the value to fit the target register.

### Description

The lower part of general-purpose register *rd* is loaded from the memory location specified by *cb.base* + *cb.offset* + *rt* + *offset*. Capability register *cb* must contain a valid capability that grants permission to load data.

The size of the value loaded depends on the value of the *t* field:

- 0** byte (8 bits)
- 1** halfword (16 bits)
- 2** word (32 bits)
- 3** doubleword (64 bits)

The extension behavior depends on the value of the *s* field: 1 indicates sign extend, 0 indicates zero extend. For example, CLWU is encoded by setting *s* to 0 and *t* to 2, CLB is encoded by setting *s* to 1 and *t* to 0.

## Pseudocode

```
if register_inaccessible(cb) then
    raise_c2_exception_v(cb)
else if not cb.tag then
    raise_c2_exception(exceptionTag, cb)
else if cb.sealed then
    raise_c2_exception(exceptionSealed, cb)
else if not cb.perms.Permit_Load then
    raise_c2_exception(exceptionPermitLoad, cb)
else
    if t = 0 then
        size ← 1
    else if t = 1 then
        size ← 2
    else if t = 2 then
        size ← 4
    else if t = 3 then
        size ← 8
    end if
    cursor ← (cb.base + cb.offset) mod 264
    addr ← cursor + rt + offset
    if sign_extend(offset) + rt + size > cb.length then
        raise_c2_exception(exceptionLength, cb)
    else if offset + rt < 0 then
        raise_c2_exception(exceptionLength, cb)
    else if align_of(addr) < size then
        raise_exception(exceptionAdEL)
    else if s = 0 then
        rd ← zero_extend(mem[addr:addr + size - 1])
    else
        rd ← sign_extend(mem[addr:addr + size - 1])
    end if
end if
```

## Exceptions

A coprocessor 2 exception is raised if:

- *cb.tag* is not set.
- *cb* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.
- Immediate *cb.offset* + *offset* + *rt* + *size* is greater than *cb.length*. **Check depends on the size of the data loaded.**
- *cb.perms.Permit\_Load* is not set.

- *cb.sealed* is set.

## Notes

- This instruction reuses the opcode from the Load Word to Coprocessor 2 (LWC2) instruction in the MIPS Specification.
- *rt* and *offset* are treated as signed integers.
- The result of the addition does not wrap around (i.e., an exception is raised if  $cb.base + cb.offset + rt + offset$  is less than zero, or greater than *maxaddr*).

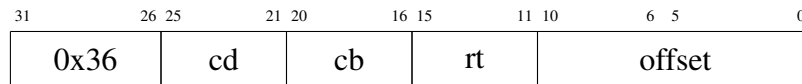
## CLC: Load Capability Register

### Format

CLC *cd*, *rt*, *offset(cb)*

CLCR *cd*, *rt(cb)*

CLCI *cd*, *offset(cb)*



### Description

Capability register *cd* is loaded from the memory location specified by *cb.base* + *cb.offset* + general-purpose register *rt*. Capability register *cb* must contain a capability that grants permission to load capabilities. The virtual address *cb.base* + *cb.offset* + *rt* must be 32-byte word aligned.

The bit in the tag memory corresponding to *cb.base* + *cb.offset* + *rt* is loaded into the tag bit associated with *cd*.

### Pseudocode

```
if register_inaccessible(cd) then
    raise_c2_exception_v(cd)
else if register_inaccessible(cb) then
    raise_c2_exception_v(cb)
else if not cb.tag then
    raise_c2_exception(exceptionTag, cb)
else if cb.sealed then
    raise_c2_exception(exceptionSealed, cb)
else if not cb.perms.Permit_Load_Capability then
    raise_c2_exception(exceptionPermitLoadCapability, cb)
end if
cursor ← (cb.base + cb.offset) mod 264
addr ← cursor + rt + offset
if rt + offset + 32 > cb.length then
    raise_c2_exception(exceptionLength, cb)
else if rt + offset < 0 then
    raise_c2_exception(exceptionLength, cb)
else if align_of(addr) < 32 then
    raise_exception(exceptionAdEL)
else if TLB(addr).L then
    cd ← mem[addr]
    cd.tag ← false
else
    cd ← mem[addr]
    cd.tag ← tags[toTag(addr)]
end if
```

## Exceptions

A coprocessor 2 exception is raised if:

- *cb* or *cd* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.
- *cb.tag* is not set.
- *cb.perms.Permit\_Load\_Capability* is not set.
- *cb.sealed* is set.
- $rt + offset + 32$  is greater than *cb.length*.
- $rt + offset < 0$ .

An address error during load (AdEL) exception is raised if:

1. The virtual address  $cb.base + cb.offset + rt$  is not 32-byte word aligned.

## Notes

- This instruction reuses the opcode from the Load Doubleword to Coprocessor 2 (LDC2) instruction in the MIPS Specification.
- *offset* is interpreted as a signed integer.
- The CLCI mnemonic is equivalent to CLC with *cb* being the zero register (\$zero). The CLCR mnemonic is equivalent to CLC with *offset* set to zero.

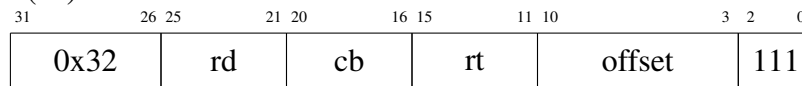
## CLLD: Load Linked Doubleword via Capability

### Format

CLLD rd, rt, offset(cb)

CLLDR rd, rt(cb)

CLLDI rd, offset(cb)



### Description

CLLD and CSCD are used to implement safe access to data shared between different threads. The typical usage is that CLLD is followed (an arbitrary number of instructions later) by CSCD to the same address; the CSCD will only succeed if the memory loaded from in the CLLD has not been modified.

The exact conditions under which CSCD fails are implementation dependent, particularly in multicore or multiprocessor implementations). The following pseudocode is intended to represent the security semantics of the instruction correctly, but should not be taken as a definition of the CPU's memory coherence model.

### Pseudocode

```
cursor ← (cb.base + cb.offset) mod 264
addr ← cursor + rt + offset
if register_inaccessible(cb) then
    raise_c2_exception_v(cb)
else if not cb.tag then
    raise_c2_exception(exceptionTag, cb)
else if cb.sealed then
    raise_c2_exception(exceptionSealed, cb)
else if not cb.perms.Permit_Load then
    raise_c2_exception(exceptionPermitLoad, cb)
else if offset + rt + 8 > cb.length then
    raise_c2_exception(exceptionLength, cb)
else if offset + rt < 0 then
    raise_c2_exception(exceptionLength, cb)
else if align_of(addr) < 8 then
    raise_c2_exception(exceptionAdEL)
else
    rd ← mem[addr:addr+7]
    linkedFlag ← true
end if
```

### Notes

- The current version of the CHERI ISA only provides a doubleword load linked instruction. Load linked instructions for other sizes (byte, half-word, word), may eventually be

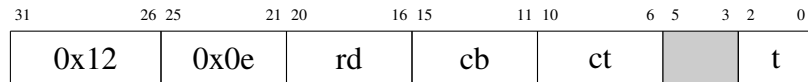
needed, as a program might need to do load linked via a capability that grants access to fewer than 16 bytes of memory.



## CPtrCmp: Capability Pointer Compare

### Format

CEQ rd, cb, ct  
 CNE rd, cb, ct  
 CLT rd, cb, ct  
 CLE rd, cb, ct  
 CLTU rd, cb, ct  
 CLEU rd, cb, ct



### Description

Capability registers *cb* and *ct* are compared, and the result of the comparison is placed in general purpose register *rd*. The rules for comparison are as follows:

- A capability with the **tag** bit unset is less than any capability with the **tag** bit set.
- Otherwise, the result of comparison is the result of comparing **base + offset** mod  $2^{64}$  for the two capabilities. Numerical comparison is signed for CLT and CLE, and unsigned for CLTU and CLEU.

This instruction can be used to compare capabilities so that capabilities can replace pointers in C executables.

Mnemonic	<i>t</i>	Comparison
CEQ	0	=
CNE	1	≠
CLT	2	< (signed)
CLE	3	≤ (signed)
CLTU	4	< (unsigned)
CLEU	5	≤ (unsigned)

### Pseudocode

```

if register_inaccessible(cb) then
  raise_c2_exception_v(cb)
else if register_inaccessible(ct) then
  raise_c2_exception_v(ct)
else
  if cb.tag ≠ ct.tag then
    equal ← false
    if cb.tag then
      less ← false
  
```

```

        signed_less ← false
    else
        less ← true
        signed_less ← true
    end if
else
    cursor1 ← (cb.base + cb.offset) mod 264
    cursor2 ← (ct.base + ct.offset) mod 264
    equal ← cursor1 = cursor2
    less ← cursor1 < cursor2
    signed_less ← to_signed64(cursor1) < to_signed64(cursor2)
end if
if t = 0 then
    rd ← equal
else if t = 1 then
    rd ← not equal
else if t = 2 then
    rd ← signed_less
else if t = 3 then
    rd ← signed_less or equal
else if t = 4 then
    rd ← less
else if t = 5 then
    rd ← less or equal
end if
end if

```

## Exceptions

A coprocessor 2 exception will be raised if:

- *cb* or *ct* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.

## Notes

- `cltu` can be used by a C compiler to compile code that compares two non-null pointers (e.g., to detect whether a pointer to a character within a buffer has reached the end of the buffer). When two pointers to addresses within the same object (e.g., to different offsets within an array) are compared, the pointer to the earlier part of the object will be compared as less. (Signed comparison would also work as long as the object did not span address 2<sup>63</sup>; the MIPS address space layout makes it unlikely that objects spanning 2<sup>63</sup> will exist in user-space C code).
- Although the ANSI C standard does not specify whether a NULL pointer is less than or greater than a non-NULL pointer (clearly, they must not be equal), the comparison instructions have been designed so that when C pointers are represented by capabilities, NULL will be less than any non-NULL pointer.

- A C compiler may also use these instructions to compare two values of type `uintptr_t` that have been obtained by casting from an integer value. If the cast is compiled as a `CFromPtr` of zero followed by `CSetOffset` to the integer value, the result of `CPtrCmp` will be the same as comparing the original integer values, because `CFromPtr` will have set **base** to zero. Signed and unsigned capability comparison operations are provided so that both signed and unsigned integer comparisons can be performed on capability registers.
- A program could use pointer comparison to determine the value of **base**, by setting **offset** to different values and testing which values cause **base** + **offset** to wrap around and be less than **base** + a zero offset. This is not an attack against a security property of the ISA, because **base** is not a secret.
- One possible way in which garbage collection could be implemented is for the garbage collector to move an object and fix up all capabilities that refer to it. If there are appropriate restrictions on which capabilities the program has to start with, the garbage collector can be sure that the program does not have any references to the object stored as integers, and so can know that it is safe to move the object. With this type of garbage collection, comparing pointers by extracting their base and offset with `cgetbase` and `cgetoffset` and comparing the integer values is not guaranteed to work, because the garbage collector might have moved the object part-way through. `CPtrCmp` is atomic, and so will work in this scenario.

## CReturn: Return to the previous security domain

### Format

CReturn



### Description

CReturn is used by a protected subsystem to return to its caller.

1. **IDC** is popped off the trusted system stack.
2. **PCC** is popped off the trusted system stack.

In the current implementation of CHERI, CReturn is implemented by the hardware raising an exception, while the rest of the behavior is implemented in software by the exception handler. Later versions of CHERI may implement more of this instruction in hardware, for improved performance. Authors of compilers or assembly language programs should not rely on CReturn being implemented in either hardware or software.

### Pseudocode (hardware)

```
raise_c2_exception_noreg(exceptionReturn)
```

### Pseudocode (software)

```
IDC ← mem[TSS .. TSS + 31]
IDC.tag ← tags[toTag(TSS)]
TSS ← TSS + 32
PCC ← mem[TSS .. TSS + 31]
PCC.tag ← tags[toTag(TSS)]
TSS ← TSS + 32
PC ← PCC.offset
```

### Exceptions

The exception raised when CReturn is implemented in software is a coprocessor 2 exception (C2E) with the capability cause code set to 0x6 (exceptionReturn) and RegNum set to *cs*. The handler vector for this exception is 0x100 above the general purpose exception handler.

An additional coprocessor 2 exception is raised if:

- The trusted system stack would underflow.
- The tag bits are not set on the memory location that are popped from the stack into **IDC** and **PCC**.

## Notes

The `CReturn` instruction will be removed in a future version of the ISA specification (though it will continue to exist as a pseudo-instruction in the assembler) and will be replaced by a specific selector in the `CCall` instruction.

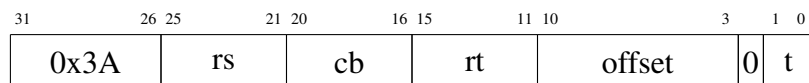
As with `CCall`, `CReturn` can have multiple implementation choices that span hardware and software. Pure software implementations can address more broad types of calling convention; pulling register clearing into the software handler provides guarantees that require less effort that might otherwise be performed redundantly by the caller and callee – although it would also require greater knowledge (and hard coding) of the ABI.

In the future, an explicit `CReturn` instruction may be replaced with invocation of `CCall` with a specific selector value.

## Store Via Capability Register

### Format

CSB *rs*, *rt*, *offset(cb)*  
CSH *rs*, *rt*, *offset(cb)*  
CSW *rs*, *rt*, *offset(cb)*  
CSD *rs*, *rt*, *offset(cb)*  
CSBR *rs*, *rt(cb)*  
CSHR *rs*, *rt(cb)*  
CSWR *rs*, *rt(cb)*  
CSDR *rs*, *rt(cb)*  
CSBI *rs*, *offset(cb)*  
CSHI *rs*, *offset(cb)*  
CSWI *rs*, *offset(cb)*  
CSDI *rs*, *offset(cb)*



### Purpose

Stores some or all of a register into a memory location.

### Description

Part of general-purpose register *rs* is stored to the memory location specified by *cb.base* + *cb.offset* + *rt* + *offset*. Capability register *cb* must contain a capability that grants permission to store data.

The *t* field determines how many bits of the register are stored to memory:

- 0** byte (8 bits)
- 1** halfword (16 bits)
- 2** word (32 bits)
- 3** doubleword (64 bits)

If less than 64 bits are stored, they are taken from the least-significant end of the register.

### Pseudocode

```
if register_inaccessible(cb) then  
    raise_c2_exception_v(cb)  
else if not cb.tag then  
    raise_c2_exception(exceptionTag, cb)  
else if cb.sealed then
```

```

        raise_c2_exception(exceptionSealed, cb)
else if not cb.Permit_Store then
        raise_c2_exception(exceptionPermitStore, cb)
end if
if t = 0 then
        size ← 1
else if t = 1 then
        size ← 2
else if t = 2 then
        size ← 4
else if t = 3 then
        size ← 8
end if
cursor ← (cb.base + cb.offset) mod 264
addr ← cursor + rt + offset
if rt + offset + size > cb.length then
        raise_c2_exception(exceptionLength, cb)
else if rt + offset < 0 then
        raise_c2_exception(exceptionLength, cb)
else if align_of(addr) < size then
        raise_exception(exceptionAdES)
else
        mem[addr:addr + size - 1] ← rs[0:size - 1]
        tags[toTag(addr)] ← false
end if

```

## Exceptions

A coprocessor 2 exception is raised if:

- *cb.tag* is not set.
- *cb* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.
- Immediate *offset* + *cb.offset* + *size* is greater than *cb.length*.
- *cb.perms.Permit\_Store* is not set.
- *cb.sealed* is set.

## Notes

- This instruction reuses the opcode from the Store Word from Coprocessor 2 (SWC2) instruction in the MIPS Specification.
- If *t* is 3 and *e* is 1, then the instruction is CSCD (Store Conditional Doubleword via Capability).

- *rt* and *offset* are treated as signed integers.
- The result of the addition does not wrap around (i.e., an exception is raised if `cb.base+cb.offset+rt+offset` is less than zero, or greater than *maxaddr*).



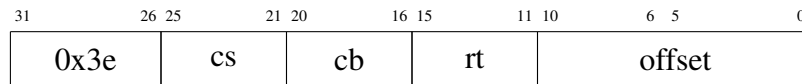
## CSC: Store Capability Register

### Format

CSC *cs*, *rt*, *offset(cb)*

CSCR *cs*, *rt(cb)*

CSCI *cs*, *offset(cb)*



### Description

Capability register *cs* is stored at the memory location specified by *cb.base* + *cb.offset* + general-purpose register *rt*, and the bit in the tag memory associated with *cb.base* + *cb.offset* + *rt* is set to the value of *cs.tag*. Capability register *cb* must contain a capability that grants permission to store capabilities. The virtual address *cb.base* + *cb.offset* + *rt* must be 32-byte word aligned.

The capability is stored in memory in the format described in Figure 3.1. **base**, **length** and **otype** are stored in memory with the same endianness that the CPU uses for double-word stores, i.e., big-endian. The bits of **perms** are stored with bit zero being the least significant bit, so that the least significant bit of the eighth byte stored is the **sealed** bit, the next significant bit is the *Global* bit, the next is *Permit\_Execute* and so on.

### Pseudocode

```
if register_inaccessible(cs) then
    raise_c2_exception_v(cs)
else if register_inaccessible(cb) then
    raise_c2_exception_v(cb)
else if not cb.tag then
    raise_c2_exception(exceptionTag, cb)
else if cb.sealed then
    raise_c2_exception(exceptionSealed, cb)
else if not cb.perms.Permit_Store_Capability then
    raise_c2_exception(exceptionPermitStoreCapability, cb)
else if not cb.perms.Permit_Store_Local_Capability and cs.tag and not cs.perms.Global then
    raise_c2_exception(exceptionPermitStoreLocalCapability, cb)
end if
cursor ← (cb.base + cb.offset) mod 264
addr ← cursor + rt + offset
if rt + sign_extend(offset) + 32 > cb.length then
    raise_c2_exception(exceptionLength, cb)
else if rt + sign_extend(offset) < 0 then
    raise_c2_exception(exceptionLength, cb)
else if cs.tag and TLB(addr).S then
    raise_c2_exception(exceptionTLBStore, cs)
else if align_of(addr) < 32 then
```

```

        raise_exception(exceptionAdES)
    else
        mem[addr] ← cs
        tags[toTag(addr)] ← cs.tag
    end if

```

## Exceptions

A coprocessor 2 exception is raised if:

- *cb.tag* is not set.
- *cs* or *cb* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.
- The virtual address  $rt + \text{offset} + 32$  is greater than *cb.length*.
- *cb.perms.Permit\_Store\_Capability* is not set.
- *cb.sealed* is set.
- *cb.perms.Permit\_Store\_Local* is not set and *cs.tag* is set and *cs.perms.Global* is not set.
- *cs.tag* is set and the *S* bit in the TLB entry for the destination address is set.

An address error during store (AdES) exception is raised if:

- The virtual address  $cb.\text{base} + cb.\text{offset} + rt + \text{offset}$  is not 32-byte word aligned.

## Notes

- If the address alignment check fails and one of the security checks fails, a coprocessor 2 exception (and not an address error exception) is raised. The priority of the exceptions is security-critical, because otherwise a malicious program could use the type of the exception that is raised to test the bottom bits of a register that it is not permitted to access.
- It is permitted to store a local capability with the tag bit unset even if the permit store local bit is not set in *cb*. This is because if the tag bit is not set then the permissions have no meaning.
- *offset* is interpreted as a signed integer.
- This instruction reuses the opcode from the Store Doubleword from Coprocessor 2 (SDC2) instruction in the MIPS Specification.
- The CSCI mnemonic is equivalent to CSC with *cb* being the zero register (\$zero). The CSCR mnemonic is equivalent to CSC with *offset* set to zero.

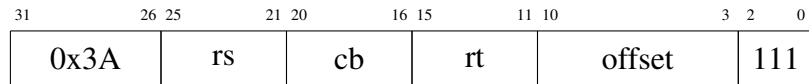
## CSCD: Store Conditional Doubleword via Capability

### Format

CSCD rs, rt, offset(cb)

CSCDR rs, rt(cb)

CSCRI rs, offset(cb)



### Pseudocode

```
cursor ← (cb.base + cb.offset) mod 264
addr ← cursor + rt + offset
if register_inaccessible(cb) then
    raise_c2_exception_v(cb)
else if not cb.tag then
    raise_c2_exception(exceptionTag, cb)
else if cb.sealed then
    raise_c2_exception(exceptionSealed, cb)
else if not cb.perms.Permit_Store then
    raise_c2_exception(exceptionPermitStore, cb)
else if rt + offset + 32 > cb.length then
    raise_c2_exception(exceptionLength, cb)
else if rt + offset < 0 then
    raise_c2_exception(exceptionLength, cb)
else if align_of(addr) < 32 then
    raise_exception(AdES)
else if not linkedFlag then
    rs ← 0
else
    mem[addr:addr+7] ← rs
    tags[toTag(addr)] ← false
    rs ← 1
end if
```

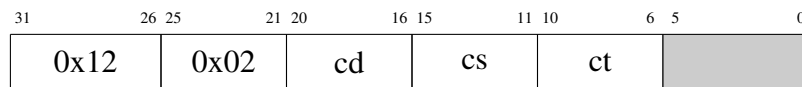
### Notes

- The current version of the CHERI ISA only provides a doubleword store conditional instruction. Store conditional instructions for other sizes (byte, half-word, word) may eventually be needed, as a program might need to do a store conditional via a capability that grants access to fewer than 16 bytes of memory.

## CSeal: Seal a Capability

### Format

CSeal *cd*, *cs*, *ct*



### Description

Capability register *cs* is sealed with an **otype** of  $ct.\text{base} + ct.\text{offset}$  and the result is placed in *cd*:

- *cd.otype* is set to  $ct.\text{base} + ct.\text{offset}$ ;
- *cd.sealed* is set;
- and the other fields of *cd* are copied from *cs*.

*ct* must grant *Permit\_Seal* permission, and the new **otype** of *cd* must be between  $ct.\text{base}$  and  $ct.\text{base} + ct.\text{length} - 1$ .

### Pseudocode

```
if register_inaccessible(cd) then
    raise_c2_exception_v(cd)
else if register_inaccessible(cs) then
    raise_c2_exception_v(cs)
else if register_inaccessible(ct) then
    raise_c2_exception_v(ct)
else if not cs.tag then
    raise_c2_exception(exceptionTag, cs)
else if not ct.tag then
    raise_c2_exception(exceptionTag, ct)
else if cs.sealed then
    raise_c2_exception(exceptionSealed, cs)
else if ct.sealed then
    raise_c2_exception(exceptionSealed, ct)
else if not ct.perms.Permitt_Seal then
    raise_c2_exception(exceptionPermitt_Seal, ct)
else if ct.offset ≥ ct.length then
    raise_c2_exception(exceptionLength, ct)
else if ct.base + ct.offset ≥ 224 then
    raise_c2_exception(exceptionLength, ct)
else
    cd ← cs
    cd.sealed ← true
    cd.otype ← ct.base + ct.offset
end if
```

## Exceptions

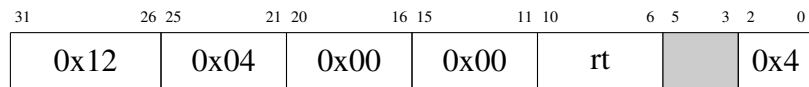
A coprocessor 2 exception is raised if:

- *cd*, *cs*, or *ct* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.
- *ct.tag* is not set.
- *ct.sealed* is set.
- *ct.perms.Permit\_Seal* is not set.
- *cs.tag* is not set.
- *cs.sealed* is set.
- *ct.offset*  $\geq$  *ct.length*
- *ct.base* + *ct.offset*  $\geq 2^{24}$

## CSetCause: Set the Capability Exception Cause Register

### Format

CSetCause *rt*



### Description

The capability cause register value is set to the low 16 bits of general-purpose register *rt*.

### Pseudocode

```
if not PCC.perms.Access_EPCC then  
    raise_c2_exception_noreg(exceptionAccessEPCC)  
else  
    CapCause ← rt[0 .. 15]  
end if
```

### Exceptions

A coprocessor 2 exception is raised if:

- **PCC.perms.Access\_EPCC** is not set.

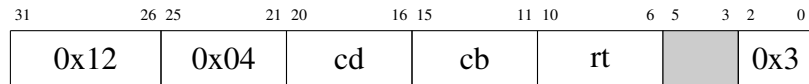
### Notes

- CSetCause does not cause an exception to be raised (unless the permission check for *Access\_EPCC* fails). CSetCause will typically be used in an exception handler, where the exception handler wants to change the cause code set by the hardware before doing further exception handling. (e.g. when the original cause code was *CCall*, the *CCall* handler detects that *CCall* should fail, and it sets *CapCause* to the reason it failed). In cases like this, it is important that EPC (etc.) are not overwritten by CSetCause.

## CSetLen: Set Length

### Format

CSetLen *cd*, *cb*, *rt*



### Description

Capability register *cd* is replaced with the contents of capability register *cb* with the **length** field set to the contents of general-purpose register *rt*.

### Pseudocode

```
if register_inaccessible(cd) then
    raise_c2_exception_v(cd)
else if register_inaccessible(cb) then
    raise_c2_exception_v(cb)
else if not cb.tag then
    raise_c2_exception(exceptionTag, cb)
else if cb.sealed then
    raise_c2_exception(exceptionSealed, cb)
else if rt > cb.length then
    raise_c2_exception(exceptionLength, cb)
else
    cd ← cb
    cd.length ← rt
end if
```

### Exceptions

A coprocessor 2 exception is raised if:

- *cb.tag* is not set.()
- *cb* or *cd* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.
- *cb.sealed* is set.
- *rt* > *cb.length*

### Notes

- Unlike **CIncBase**, this operation will always raise an exception if *cb.tag* is unset or *cb.sealed* is set, even if the length is unchanged.

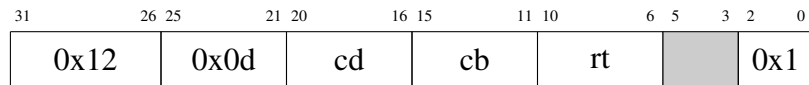
- Assembly programmers must exercise caution when using `CIncBase` or `CSetLen` on capabilities such as C pointers that have non-zero **offset** values. Incorrect operation orders may produce empty capabilities.



## CSetOffset: Set cursor to an offset from base

### Format

CSetOffset *cd*, *cb*, *rt*



### Description

Capability register *cd* is replaced with the contents of capability register *cb* with the **offset** field set to the contents of general purpose register *rt*.

### Pseudocode

```
if register_inaccessible(cd) then
    raise_c2_exception_v(cd)
else if register_inaccessible(cb) then
    raise_c2_exception_v(cb)
else if cb.tag and cb.sealed then
    raise_c2_exception(exceptionSealed, cb)
else
    cd ← cb
    cd.offset ← rt
end if
```

### Exceptions

A coprocessor 2 exception is raised if:

- *cb* or *cd* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.
- *cb.tag* is set and *cb.sealed* is set.

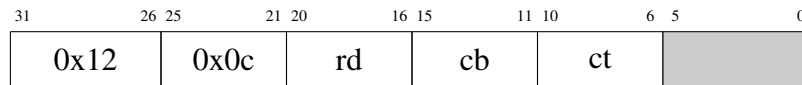
### Notes

- `CSetOffset` can be used on a capability register whose tag bit is not set. This can be used to store an integer value in a capability register, and is useful when implementing a variable that is a union of a capability and an integer (`intcap_t` in C). The in-memory representation that will be used if the capability register is stored to memory might be surprising to some users (**base** + **offset** is stored in the **cursor** field in memory) and may change if the memory representation of capabilities changes, so compilers should not rely on it.

## CToPtr: Capability to Pointer

### Format

CToPtr rd, cb, ct



### Description

If *cb* has its tag bit unset (i.e. it is either the NULL capability, or contains some other non-capability data), then *rd* is set to zero. Otherwise, *rd* is set to  $cb.\mathbf{base} + cb.\mathbf{offset} - ct.\mathbf{base}$

This instruction can be used to convert a capability into a pointer that uses the C language convention that a zero value represents the NULL pointer. Note that *rd* will also be zero if *cb* and *ct* have the same base; this is similar to the C language not being able to distinguish a NULL pointer from a pointer to a structure at address 0.

### Pseudocode

```
if register_inaccessible(cb) then
    raise_c2_exception_v(cb)
else if register_inaccessible(ct) then
    raise_c2_exception_v(ct)
else if not ct.tag then
    raise_c2_exception(exceptionTag, ct)
else if not cb.tag then
    rd ← 0
else
    rd ← (cb.base + cb.offset - ct.base) mod 264
end if
```

### Exceptions

A coprocessor 2 exception will be raised if:

- *cb* or *ct* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.
- *ct.tag* is not set.

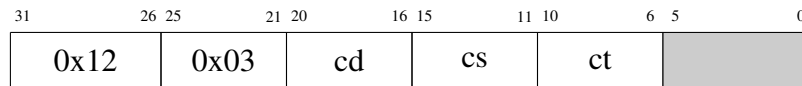
### Notes

- *cb* or *ct* being sealed will not cause an exception to be raised. This is for further study.

## CUnseal: Unseal a sealed capability

### Format

CUnseal *cd*, *cs*, *ct*



### Description

The sealed capability in *cs* is unsealed with *ct* and the result placed in *cd*. The global bit of *cd* is the AND of the global bits of *cs* and *ct*. *ct* must be unsealed, have Permit\_Seal permission, and *ct*.base + *ct*.offset must equal *cs*.otype.

### Pseudocode

```
if register_inaccessible(cd) then
    raise_c2_exception_v(cd)
else if register_inaccessible(cs) then
    raise_c2_exception_v(cs)
else if register_inaccessible(ct) then
    raise_c2_exception_v(ct)
else if not cs.tag then
    raise_c2_exception(exceptionTag, cs)
else if not ct.tag then
    raise_c2_exception(exceptionTag, ct)
else if not cs.sealed then
    raise_c2_exception(exceptionSealed, cs)
else if ct.sealed then
    raise_c2_exception(exceptionSealed, ct)
else if ct.base + ct.offset ≠ cs.otype then
    raise_c2_exception(exceptionType, ct)
else if not ct.perms.Permits_Seal then
    raise_c2_exception(exceptionPermitsSeal, ct)
else if ct.offset ≥ ct.length then
    raise_c2_exception(exceptionLength, ct)
else
    cd ← cs
    cd.sealed ← false
    cd.otype ← 0
    cd.perms.Global ← cs.perms.Global and ct.perms.Global
end if
```

### Exceptions

A coprocessor 2 exception is raised if:

- *cd*, *cs*, or *ct* is one of the reserved registers (**KR1C**, **KR2C**, **KCC**, **KDC** or **EPCC**) and the corresponding bit in **PCC.perms** is not set.
- *ct.tag* is not set.
- *ct.sealed* is set.
- *ct.perms.Permit\_Seal* is not set.
- *ct.offset*  $\geq$  *ct.length*
- *ct.base* + *ct.offset*  $\neq$  *cs.otype*.
- *cs.tag* is not set.
- *cs.sealed* is not set.

## 4.3 Assembler Pseudo-Instructions

For convenience, several pseudo-instructions are accepted by the assembler. These expand to either single instructions or short sequences of instructions.

### 4.3.1 Capability Move

`CMove` is a pseudo operation that moves a capability from one register to another. It expands to a `CIncBase` instruction, with `$zero` as the increment operand.

```
1 # The following are equivalent:
2   CMove $c1, $c2
3   CIncBase $c1, $c2, $zero
```

### 4.3.2 Get/Set Default Capability

`CGetDefault` and `CSetDefault` get and set the capability register that is used by the legacy MIPS load and store instructions. In the current version of the ISA, this register is `C0`. These pseudo-operations are provided for the benefit of the LLVM compiler: the compiler can more easily detect that a write to `C0` affects the meaning of subsequent legacy MIPS instructions if these are separate pseudo-operations.

```
1 # The following are equivalent:
2   CGetDefault $c1
3   CIncBase $c1, $c0, $zero
```

```
1 # The following are equivalent:
2   CSetDefault $c1
3   CIncBase $c0, $c1, $zero
```

### 4.3.3 Capability Loads and Stores of Floating-Point Values

The current revision of the CHERI ISA does not have instructions for loading floating point values directly via capabilities. MIPS does provide instructions for moving values between integer and floating point registers, so a load or store of a floating point value via a capability can be implemented in two instructions.

Four pseudo-instructions are defined to implement these patterns. These are `clwc1` and `cldc1` for loading 32-bit and 64-bit floating point values, and `cswc1` and `csdc1` as the equivalent store operations. The load operations expand as follows:

```
1 | cldc1    $f7, $zero, 0($c2)
2 | # Expands to:
3 | cld     $1, $zero, 0($c2)
4 | dmtc1   $1, $f7
```

Note that integer register `$1` (`$at`) is used; this pseudo-op is unavailable if the `noat` directive is used. The 32-bit variant (`clwc1`) has a similar expansion, using `clwu` and `mtc1`.

The store operations are similar:

```
1 | csdc1    $f7, $zero, 0($c2)
2 | # Expands to:
3 | dmfc1   $1, $f7
4 | csd     $1, $zero, 0($c2)
```

The specified floating point value is moved from the floating point register to `$at` and then stored using the correct-sized capability instruction.

# Chapter 5

## Design Rationale

During the design of CHERI, we considered many different capability architectures and design approaches. This chapter describes the design choices that were made, briefly outlines some possible alternatives, and provides a rationale for the choices that were made.

### High-Level Design Approach: Capabilities as Pointers

Our goals of providing fine-grained memory protection and compartmentalization led to an early design choice to approach capabilities as a form of pointer. This rapidly led to a number of conclusions:

- Capabilities are within virtual address spaces, imposing an ordering in which capability protections are evaluated before virtual-memory protections; this in turn had implications for the hardware composition of the capability coprocessor and conventional MMU interact.
- Capabilities are treated by the compiler in much the same way as pointers, meaning that they will be loaded, manipulated, dereferenced, and stored via registers and to/from general-purpose memory by explicit instructions, which we used modeled on similar conventional RISC instructions.
- Incremental deployment within programs meant that not all pointers would immediately be converted to capabilities, implying that both might coexist in the same memory; also, there was a strong desire to embed capabilities within data structures, rather than store them in separate segments, requiring fine-granularity tagging.
- Incremental deployment and compatibility with the UNIX model implied retaining the general-purpose memory management unit (MMU) pretty much as currently designed, including support for variable page size, TLB layout, etc. The MIPS ISA describes a software-managed TLB rather than hardware page-table walking as is present in most other ISAs; this is not fundamental to our approach, and either model would work.

## Capability-Register File

The decision to separate the capability-register file from the general-purpose register file is somewhat arbitrary from a software-facing perspective: we envision capabilities gradually displacing general-purpose registers as pointers, but that management of the two register files will remain largely the same, stack spilling will behave the same way, and so on. We selected the separate representation for a few pragmatic reasons:

- Coprocessor interfaces frequently make the assumption of additional register files (a la floating-point registers).
- Capability registers are quite large, and by giving the capability coprocessor its own pipeline for manipulations, we could avoid enforcing a 256-wide path through the main pipeline.
- It is more obvious, given a coprocessor-based interface, how to provide compatibility support in which the capability coprocessor is “disabled,” the default configuration in order to support unmodified MIPS compilers and operating systems.

However, it is entirely possible to imagine a variation on the CHERI design in which, more similar to the manner in which the 32-bit x86 ISA was extended to support 64-bit registers, the two files were conflated and able to hold both general-purpose and capability registers. Early in our design cycle, capability registers were able to hold only true capabilities (i.e., with tags); later, we weakened this requirement by adding an explicit tag bit to each register in order to improve support for capability-oblivious code such as memory-copy routines able to copy data structures consisting of both capabilities and ordinary data. This shifts our approach somewhat more towards a conflated approach; our view is that efficiency of implementation and compatibility, rather than negligible effect on the software model, would be the primary reasons to select one approach or another for a particular starting-point ISA.

Another design variation might have more tightly coupled specific capability registers with general-purpose registers – an approach we discussed extensively, especially when comparing with the bounds-checking literature which has explored techniques based on *sidecar registers* or associative look-aside buffers. Many of these approaches did not adopt tags as a means of strong integrity protection, which we require for the compartmentalization model, making associative techniques less suitable. Further, we felt that the working-set properties of the two register files might be quite different, and effectively pinning the two to one another would reduce the efficiency of both.

It is worth considering, however, that our recent interest in cursors within capabilities revisits both of these ideas.

## Representation of Memory Segments

CHERI capabilities represent a region of memory by its base address and length; memory accesses are relative to the base address. An alternative representation would have been for capabilities to contain an upper and lower bound on addresses within the memory region, with memory accesses being given in terms of absolute addresses but checked against the upper and lower bound.

The base and length representation was chosen because it is more convenient for arrays and structures in the C language. Given a capability for an array and an index into the array, the array element can be read with (for example) `CLB` without the need for an addition in software. (In C, all arrays are zero based. This is not the case in other languages, e.g., Ada.) The length of a structure is usually known at compile time, and the length of a capability can be set to the length of a structure with `CSetLen`; setting an upper bound would require an additional addition instruction to compute it.

Although CHERI does not attempt to keep the base address of a capability secret, the use of base-relative (rather than absolute) addresses for memory accesses reduces the need to keep the absolute base address of a capability in a general purpose register, and possibly might facilitate code migration to a stricter version of the architecture in which absolute addresses are secret.

The disadvantages of the base and length representation are that:

- There is no way to grant access to the very last byte of the virtual address space (a base of 0 and a length of  $2^{64} - 1$  grants access to addresses 0 to  $2^{64} - 2$ ).
- Base-relative addressing is cumbersome for code capabilities. If a program wants to call a subroutine, and to grant the subroutine execute access only to its own instructions and not to the entire program text, then the subroutine needs to be linked differently from the calling program, because branches within the subroutine will be relative to a different base.

A key concern with the current representation is its substantial size – simulation suggests that cache footprint is a dominant factor in performance, although optimization techniques such as CCured would reduce this effect. We believe that a reduction to 128-bit capability registers would come at an observable cost to both protection scalability (e.g., limiting the number of bits in a pointer to 40-48 bits rather than the full 64) as well as compartmentalization functionality (e.g., having fewer software-defined permission bits). However, in practice this may prove necessary to support widespread adoption. Some care must be taken to retain current software flexibility, especially regarding very fine-grained regions of memory, which are highly desirable to support critical protection properties for C – e.g., granular stack protection and arbitrary subdivision of character-based strings into separate bounded regions. It could be that pointer compression techniques eliding specific middle bits in the address space, or possibly trading off size and granularity (e.g., bits might be invested either in describing very small objects at arbitrary alignment, or very large objects at more coarse alignment) provide a useful middle ground.

## Signed and Unsigned Offsets

In the CHERI instructions that take both a register offset and an immediate offset, the register offset is treated as unsigned integer but the immediate offset is treated as a signed integer.

Register offsets are treated as unsigned so that given a capability to the entire address space (except for the very last byte, as explained above), a register offset can be used to access any byte within it. Signed register offsets would have the disadvantage that negative offsets would fail the capability bounds check, and memory at offsets within the capability greater than  $2^{63}$  would not be accessible.



Immediate offsets, on the other hand, are signed, because the C compiler often refers to items on the stack using the stack pointer as register offset plus a negative immediate offset. We have already encountered observable difficulty due to a reduced number of bits available for immediate offsets in capability-relative memory operations when dealing with larger stack-frame sizes; it is unclear what real performance cost this might have (if any), but does reemphasize the importance of careful investment of encoding bits for instructions.

## Overwriting Capabilities

In CHERI, if a capability in memory is partly overwritten with non-capability data, then the memory contents afterwards will be the capability converted to a byte representation and then overwritten.

Alternative designs would have been for the capability to be zeroed first before being overwritten; or for the write to raise an exception (with an explicit “clear tag in memory” operation for the case when a program really intends to overwrite a capability with non-capability data).

The chosen approach is simpler to implement in hardware. If store instructions needed to check the tag bit of the memory location that was being written, then they would need to have a read-modify-write cycle to the memory, rather than just a write; in general, the MIPS architecture carefully avoids the need for a read-modify-write cycle within a single instruction. (Although, once the memory system needs to deal with cache coherence, a write is not that much simpler than a read-modify-write.)

The CHERI behavior also has the advantage that programs can write to a memory location (e.g., when spilling a register on to the stack) without needing to worry about whether that location previously contained a capability or non-capability data.

A potential disadvantage is that the contents of capabilities cannot be kept secret from a program that uses them. A program can always discover the contents of a capability by overwriting part of it, then reading the result as non-capability data. In CHERI, there are other, more direct, ways for a program to discover the contents of a capability it owns, so this is not a security vulnerability.

However, there are ABI concerns: we have tried to design the ISA in such a way that software does not need to be aware of the in-memory layout of capabilities, but as it is necessarily exposed, there is a risk that software might become dependent on a specific layout. One case of particular note is in the operating-system paging code, which must save and restore capabilities and their tags separately; this can be accomplished by using instructions such as `CGetBase` on untagged values loaded from disk and then refining an in-hand capability using `CSetBase` – an important reason not to limit capability field retrieval instructions to tagged values.

## Reading Capabilities as Bytes

In CHERI, if a data load instruction such as `CLB` is used on a memory location containing a capability, the internal representation of the capability is read. An alternative architecture would have such loads return zero, or raise an exception.

Because the contents of capabilities are not secret, allowing them to be read as raw data is not a security vulnerability.

## Capability registers are dynamically tagged

In CHERI, capability registers and memory locations have a tag bit that indicates whether they hold a capability or non-capability data. (An alternative architecture would give memory locations a tag bit, where capability registers could contain only capabilities – with an exception raised if an attempt were made to load non-capability data into a capability register with `CLC`.)

Giving capability registers and memory locations a tag bit simplifies the implementation of `memcpy()`. `memcpy()` is a variant of `memcpy()` that copies the tag bit as well as the data, and so can be used to copy structures containing capabilities. As capability registers are dynamically tagged, `memcpy()` can copy a structure by loading it into a capability register and storing it to memory, without needing to know at compile time whether it is copying a capability or non-capability data.

Tag bits on capability registers may also be useful for dynamically typed languages in which a parameter to a function can be (at run time) either a capability or an integer. `memcpy()` can be regarded as a function whose parameter (technically a `void`) is dynamically typed.

## Separate Permissions for Storing Capabilities and Data

CHERI has separate permission bits for storing a capability versus storing non-capability data (and similarly, for loading a capability versus loading non-capability data).

(An alternative design would be just one `Permit_Load` and just one `Permit_Store` permission that were used for both capabilities and non-capability data.)

The advantage of separate permissions bits for capabilities is that there can be two protected subsystems that communicate via a memory buffer to which they have `Permit_Load` and `Permit_Store` permissions, but do not have `Permit_Load_Capability` or `Permit_Store_Capability`. Such communicating subsystems cannot pass capabilities via the shared buffer, even if they collude. (We realized that this was potentially a requirement when trying to formally model the security guarantees provided by CHERI.)

## Capabilities Contain a Cursor

In the C language, pointers can be both incremented and decremented. C pointers are sometimes used as a cursor that points to the current working element of an array, and is moved up and down as the computation progresses.

CHERI capabilities include an offset field, which gives the difference between the base of the capability and the memory address that is currently of interest. The offset can be both incremented and decremented without changing `base`, so that it can be used to implement C pointers.

In the ANSI C standard, the behavior is undefined if a pointer is incremented more than *one* beyond the end of the object to which it points. However, we have found that many existing C programs rely on being able to increment a pointer beyond the end of an array, decrement it back within range, and then dereference it. In particular, network packet processing software often does this. In order to support programs that do this, CHERI offsets are allowed to take on any value. A range check is performed when the capability is dereferenced, so buffer overflows

are prevented, but the offset can take on intermediate out of range values as long as it is not dereferenced.

An alternative architecture would have not included an offset within the capability. This could have been supported by two different capability types in C, one that could not be decremented (but was represented by just a capability) and one that supported decrementing (but was represented by a pair of a capability and a separate integer for the offset). Programming languages that did not have pointer arithmetic could have their pointers compiled as just a capability.

The disadvantage of including offsets within capabilities is that it wastes 64 bits in each capability in cases where offsets are not needed (e.g., when compiling languages that don't have pointer arithmetic, or when compiling C pointers that are statically known to never be decremented). The alternative (no offset) architecture could have used those 64 bits of the capability for other purposes, and stored an extra offset outside the capability when it was known to be needed. The disadvantage of the no-offset architecture is that C pointers become either unable to support decrement or larger: because capabilities need to be aligned, a pair of a capability and an integer will usually end up being padded to the size of two capabilities, doubling the size of a C pointer, and this is a serious performance consideration.

Another disadvantage of the no-offset alternative is that it makes the seal/unseal mechanism considerably more complicated and hard to explain. A program that has a capability for a range of types has to somehow select which type within its permitted range of types it wishes to use when sealing a particular data capability. The CHERI architecture uses the offset for this purpose; not having an offset field leads to more complex encodings when creating sealed capabilities.

By comparison, the CCured language includes both `FSEQ` and `SEQ` pointers. CHERI capabilities are analogous to CCured's `SEQ` pointers. The alternative (no offset) architecture would have capabilities that acted like CCured's `FSEQ`, and used an extra offset when implementing `SEQ` semantics.

## NULL does not have the Tag Bit Set

In some programming languages, pointer variables must always point to a valid object. In C, pointers can either point to an object or be NULL; by convention, NULL is the integer value zero cast to a pointer type.

If hardware capabilities are used to implement a language that has NULL pointers, how is the NULL pointer represented? CHERI capabilities have a **tag** bit; if the **tag** bit is set, a valid capability follows, otherwise the remaining data can be interpreted as (for example) bytes or integers. The representation we have chosen for NULL is that the **tag** bit is not set and the **base** and **length** fields are zero; effectively, NULL is the integer zero stored as a non-capability value in a capability register.

An alternative representation we could have chosen for NULL would have been with the **tag** bit set, and zero in the **base** field and **length** fields. Effectively, NULL would have been a capability for an array of length zero.

Many of the CHERI instructions are agnostic as to which of these two conventions for NULL is employed, but the `CFromPtr`, `CToPtr` and `CPtrCmp` operations are aware of the convention.

The advantages of NULL's **tag** bit being unset are:

- Initializing a region of memory by writing zero bytes to it will initialize all capability variables within the region to the NULL capability. Initializing memory by writing zeros is, for example, done by the C `calloc()` function, and by some operating systems.
- It is possible for code to conditionally branch on a capability being NULL by using the `CBTS` or `CBTU` instruction.

## Permission Bits Determine the Type of a Capability

In CHERI, a capability's permission bits together with the **sealed** bit determine what kind of capability it is. A capability for a region of memory has **sealed** unset and *Permit\_Load* and/or *Permit\_Store* set; a capability for an object has **sealed** set and *Permit\_Execute* unset; a capability to call a protected subsystem (a "call gate") has **sealed** set and *Permit\_Execute* set; a capability that allows the owner to create objects whose type identifier (**otype**) falls within a range has **sealed** set and *Permit\_Seal* set.

An alternative architecture would have included a separate *capability type* field, as well as the **perms** field, within each capability; the meaning of the rest of the bits in the capability would have been dependent on the value of the *capability type* field.

A potential disadvantage of not having a *capability type* field is that different kinds of capability cannot use the remaining bits of the capability in different ways.

A consequence of the architecture we have chosen is that it is possible to create many different kinds of capability (2 to the power of the number of permission bits plus **sealed**). Some of the kinds of capability that it is possible to create do not have a clear use case; they just exist as a consequence of the representation chosen for capabilities.

## Object Types are not Addresses

In CHERI, we make a distinction between the unique identifier for an object type (the **otype** field) and the address of the executable code that implements a method on the type (the **base** + **offset** fields in a sealed executable capability).

An alternative architecture would have been to use the same fields for both, and take the entry address of an object's methods as a convenient unique identifier for the type itself.

The architecture we have chosen is conceptually simpler and easier to explain. It has the disadvantage that the type field is only 24 bits, as there is insufficient space inside the capability for more.

The alternative of treating the set of object type identifiers as being the same as the set of memory addresses enables the saving of some bits within a capability by using the same field for both. It also simplifies assigning type identifiers to protected subsystems: each subsystem can use its start address as the unique identifier for the type it implements. Subsystems that need to implement multiple types, or create new types dynamically can be given a capability with the permission *Permit\_Set\_Type* set for a range of memory addresses, and they are then able to use types within that range. (The current CHERI ISA does not include the *Permit\_Set\_Type* permission; it would only be needed for this alternative approach). This avoids the need for

some sort of privileged type manager that creates new type identifiers; such a type manager is potentially a source of covert channels. (Suppose that the type manager and allocated type identifiers in numerically ascending order. A subsystem that asks the type manager twice for a new type id and gets back  $n$  and  $n + 1$  knows that no other subsystem has asked for a new type id in between the two calls; this could in principle be used for covert communication between two subsystems that were supposed to be kept isolated by the capability mechanism.)

## Unseal is an Explicit Operation

In CHERI, converting a pointer to an opaque object into a pointer that allows the object's contents to be inspected or modified directly is an explicit operation. It can be done directly with the `CUnseal` operation, or by using `CCall` to run the result of unsealing the first argument on the result of unsealing the second argument.

An alternative architecture would have been one with “implicit” unsealing, where a sealed capability (**sealed set**) could be dereferenced without explicitly unsealing it first, provided that the subsystem attempting the dereference had some kind of ambient authority that permitted it to dereference sealed capabilities of that type. This ambient authority could have taken the form of a protection ring or the **otype** field of **PCC**.

The disadvantage of the architecture we have chosen is that protected subsystems need to be careful not to leak capabilities that they have unsealed, for example by leaving them on the stack when they return to their caller. In an architecture with “implicit unseal”, protected subsystems would just need to delete their ambient authority for the type before returning, and would not need to explicitly clean up all the unsealed capabilities that they had created.

## EPCC is a numbered register

The exception program counter (**EPCC**) is a numbered register. An alternative architecture would have been to make **EPCC** like **PCC**, and only accessible via a special instruction. This alternative architecture would have had several advantages:

- If **EPCC** is set to a bad value (e.g. tag bit unset) and then an exception occurs, it is not obvious what the CPU should do. The behavior in this situation is explicitly undefined in the current CHERI ISA. Raising an exception is problematic because **EPCC** is already invalid at that point. If **EPCC** could only be set using a special instruction, then that instruction could check that the proposed new value of **EPCC** was valid, and raise an exception (using the old, valid, **EPCC**) if it wasn't.
- For compatibility with legacy operating systems that are unaware of capabilities, it is desirable that an exception handler that just sets `CP0.EPC` (and not **EPCC.offset**) should work. But that then raises the question of what happens if operating system code changes both `CP0.EPC` and **EPCC.offset**. Ideally, these should behave as if they are the same register. But this can be complex to implement in hardware: **EPCC** and **EPC** are in different register files, and **EPCC** can be modified by many different capability instructions, because it is a numbered register; it can be complex to make all of these update

CP0.**EPC** as well. If **EPCC** was not a numbered register, then only the special instruction for changing **EPCC** should need to be aware that CP0.**EPC** needs to be changed whenever **EPCC** changes.

# Chapter 6

## Programming Languages and Operating Systems

We capture some of our early thoughts on the topic of use of CHERI instructions in programming languages and operating systems. The goal of our software work is to test several fundamental hypotheses underlying the CHERI architecture, such as the following:

- That a hardware-software object-capability model provides superior performance when large numbers of protection domains are required, as well as large numbers of protection-domain transitions.
- That protection domains within address spaces offer improved programmability and debuggability for compartmentalized TCB components.
- That there are (fairly) natural mappings from higher-level language pointer and reference models into memory-capability semantics.
- That capability adaptation of common TCB components offers dramatically improved robustness and security.
- That a hardware capability model and MMU-based virtual addressing can not only co-exist, but can also facilitate an adoption of capability approaches – offering both an incremental adoption path with immediate security benefits and a long-term vision for software security improvement.
- That a fully virtualizable per-address-space capability system is feasible and practical, while allowing use of capability models within individual hierarchical rings and not necessarily hierarchical address spaces and execution domains.

To this end, we are developing a significantly rich software stack that utilizes the CHERI feature set, implementing and exercising various aspects of the hybrid capability model.

### 6.1 Development Plan and Status

At the time of writing, we are a little more into the fifth year of a five-year research project in hardware and software security. We have successfully prototyped a fully pipelined 64-bit CPU

implementing the 64-bit MIPS and CHERI ISAs; detailed information on this prototype can be found in the accompanying *BERI Hardware Reference*. The following sections document our recent accomplishments and continuing strategy for exploring the software implications of the CHERI architecture.

## 6.2 Open-Source Foundations

During the initial bring-up phase of our prototype CHERI CPU, CHERI's support for incremental adoption has proven invaluable: we been able to rely on current open-source boot loaders, operating systems, programming languages, compilers, debuggers, and applications, to selectively deploy capability features in the most critical software foundations and the most vulnerable services.

## 6.3 Current Software Implementation

We have extended existing 64-bit MIPS versions of FreeBSD operating system, GNU assembler, and LLVM/Clang compiler suite to support CHERI ISA features.

## 6.4 CheriBSD

We have extended the existing 64-bit MIPS port of FreeBSD to support a range of Altera/Terasic hardware peripherals, and our CHERI ISA extensions. CheriBSD maintains a coprocessor 2 context for each user thread; it implements CCall/CReturn exception handlers, and recovery paths for when a sandbox triggers a hardware exception (e.g., due to an invalid memory reference). CheriBSD includes a new libcheri which implements a sandbox API and a growing set of system services that may be delegated to sandboxed code.

## 6.5 Extended GNU Assembler (gas)

We have extended the GNU assembler (gas) to support the CHERI ISA, by allowing assembly files, and an inline assembler from C to make use of CHERI instructions. Tools such as **objdump** are also able to interpret CHERI instructions. We have not yet extended the linker to support new CHERI-related linkage types; instead, we rely on hybrid behavior to implement programs for the time being.

## 6.6 CHERI LLVM/Clang

LLVM is a framework for implementing compilers that comprises a well-defined intermediate representation (IR), a set of APIs for generating this representation, optimization passes for transforming it, and back ends for generating native code. We have extended the MIPS back end in LLVM to provide support for capability instructions.



We reserve address space 200<sup>1</sup> for capability pointers. Any pointer to address space 200 is assumed to be a capability. We also add explicit integer to pointer and pointer to integer patterns in the back end. These are required because all existing LLVM back ends regard pointers and integers as interchangeable.

The LLVM back end also provides an alternative assembler. This currently lacks some MIPS instructions, so is not yet a replacement for the GNU assembler, but does provide support for both inline assembly and for stand-alone assembly files.

The modifications to LLVM are intended to make experimentation with programming languages easier. Any language front end that can generate LLVM IR can be modified to support capabilities and we are free to experiment with new languages and modifications to others.

In addition to the support for capabilities as pointers, we also provide a number of intrinsics that map closely to instructions. The example below uses the `llvm.cheri.set.cap.length` intrinsic, which sets the length of a capability. This example shows the LLVM IR for a simple function that wraps the C standard `malloc()` in one that returns a capability that will enforce the length.

```
1 define i8 addrspace(200)* @cmalloc(i64 %s) nounwind {  
2 entry:  
3 ; Call malloc()  
4 %call = tail call i8* @malloc(i64 %s) nounwind  
5 ; Convert the C0-relative pointer to a capability  
6 %0 = ptrtoint i8* %call to i64  
7 %1 = inttoptr i64 %0 to i8 addrspace(200)*  
8 ; CSetLen  
9 %2 = tail call i8 addrspace(200)* @llvm.cheri.set.cap.length(i8  
10 addrspace(200)* %1, i64 %s)  
11 ret i8 addrspace(200)* %2  
}
```

Clang is a front end for LLVM that generates LLVM IR from C-family languages (C, C++, Objective-C, and Objective-C++). Our first work on language extensions involves providing capability support to C. Programmers can annotate pointers as being capabilities, which triggers CHERI rather than MIPS code generation for any resulting memory accesses. We have an experimental version of the same compiler code that now supports stack access via capability.

Objective-C provides a late-bound object oriented model on top of C that makes it an interesting test ground for experimentation. Combined with the MIT-licensed GNUstep Objective-C runtime, we can experiment with adding language features to Objective-C based on our extensions to C.

## 6.7 CHERI Unit-Test Suite

Using the CHERI assembler, we have extended our existing MIPS ISA test suite to exercise various aspects of the CHERI ISA. The current test suite validates the behavior of memory capabilities and capability exceptions, including monotonic decrease in rights using capabil-

---

<sup>1</sup>The number 200 is subject to change and should not be relied on. Address spaces under 256 are intended to be reserved for architecture-agnostic uses, so we may either move this to a higher number, or retain a low number as a general fat-pointer address space in LLVM IR.

ity manipulation instructions. More detailed implementation status for the CHERI hardware prototype may be found in the *BERI Hardware Reference*.

## 6.8 Future Plans

Between 2010 and 2014, we completed basic prototyping of the CHERI hardware and software platform. We are now considering potential future directions for further hardware and software experimentation, including bindings to higher-level languages such as Objective-C, OCaml, and Java. We are also exploring how adding capability support to LLDB and the LLVM debugger would allow us to develop capability-aware programs, as well as considering security implications for designing and implementing debuggers.

# Chapter 7

## Future Directions

The CTSRD project, of which CHERI is just one element, has now been in progress for four years. Our focuses to date have been in several areas:

1. Design the CHERI Instruction-Set Architecture based on hybrid memory-capability and object-capability models. As part of this work, develop PVS/SAL formal models of the ISA, and analyze properties about program expressivity; subsequently, develop a more complete L3 model of both MIPS and CHERI ISAs that would allow automated test generation, trace validation, and eventually, software verification.
2. Flesh out the ISA feature set in CHERI to support a real-world operating system – primarily, this has consisted of adding support for the system management coprocessor, CP0, which includes the MMU and exception model, but also features such as a programmable interrupt controller (PIC). We have also spent considerable time refining second and third versions of the ISA intended to better support high levels of C-language compatibility, as well as automatic use by compilers, which are now implemented. This work has incorporated ideas from the C-language fat-pointer research literature.
3. Prototype, test, and refine CHERI ISA extensions, which are incorporated via a new capability coprocessor, CP2. We have open sourced the BERI and CHERI processor designs in order to allow reproducible experimentation with our approach, as well as to act as an open-source platform for other future hardware-software research projects.
4. Port the FreeBSD operating system first to a capability-free version of CHERI, known as BERI. This is known as FreeBSD/BERI, and this support has been upstreamed such that new releases of FreeBSD support the BERI processor and its peripheral devices.
5. Adapt FreeBSD to make use of CHERI features – first by adapting the kernel to maintain new state and provide object invocation, and then low-level system runtime elements, such as the system library and runtime linker. This is known as CheriBSD. While open sourced, these changes remain outside of the upstream FreeBSD repository due to their experimental nature.
6. Adapt the Clang/LLVM compiler suite to be able to generate CHERI ISA instructions as directed by C-language annotations. We have upstreamed substantial improvements to Clang/LLVM MIPS support, as well as changes making it easier to support ISA features such as extra-wide pointers utilized in the CHERI ISA.

7. Begin to develop semi-automated techniques to assist software developers in compartmentalizing applications using Capsicum and CHERI features. This is a subproject known as Security-Oriented Analysis of Application Programs (SOAAP), and performed in collaboration with Google.
8. Develop FPGA-based demonstration platforms, including an early prototype on the Terasic tPad, and more mature server-style and tablet-style prototypes based on the Terasic DE4 board. We have also made use of CHERI on the NetFGPA 10G board.
9. Develop techniques for translating Bluespec hardware designs into PVS representations so that they can be used for formal analysis purposes. The SRI PVS tool suite has been embedded in the BluespecVerilog compiler chain to enable formal verification, model checking (SAL), and SMT solving (Yices) inline with the compilation.
10. In addition, we are working on a new project task to consider extensions of the CHERI hardware-software architecture that would enable trustworthy systems to be developed despite the presence of numerous potentially untrustworthy microprocessors and devices, some which might otherwise have overly permissive unmonitored direct access to processor memory.

We have made a strong beginning, but clearly there is still much to do in our last year of CTSRD. From this vantage point, we see a number of tasks ahead, which we detail in the next few sections.

In addition, our companion project for DARPA's Mission-oriented Resilient Cloud program has been extending some of the CHERI concepts, notably pursuing multicore and multithreaded CHERI in anticipation of its use in network controllers and cloud datacenters. In addition, we have been experimenting with applications of the CHERI hardware architecture for trustworthy switching and other purposes, and considering efficient operating systems (e.g., DIOS and MirageOS) that might be well adapted to CHERI. This is still work in progress, as that project is a year behind our CRASH project, although both run until September 2015.

## 7.1 An Open-Source Research Processor

One of our goals for the CHERI processor is to produce a reference Bluespec processor implementation, which can then be used as a foundation not only for CHERI, but also for other research projects in the hardware-software interface. Capability processor extensions to the MIPS ISA would then be a core research result from this project, but also the first example of research conducted on the reference processor.

We have spent a considerable amount of effort in open-sourcing CHERI, including enhancing our test suite, updating documentation, and creating a new open-source license intended for hardware-software projects (derived from the Apache software license).

## 7.2 Formal Methods for Bluespec

We have created prototype descriptions of the CHERI ISA in PVS and SAL, and are collaborating with the REMS project at Cambridge to develop an L3 model of the MIPS ISA, with the

intent of also applying it to CHERI. We have used our formal models to automatically generate test suites, and to prove higher-level properties about what the ISA can represent. We have created new tools to automatically process Bluespec designs for use in theorem proving and model checking, and developed new tools to improve SMT performance and to extract higher-level properties from hardware designs. Our longer-term goal has been to link formal models of the hardware itself with the ISA specification and software compiled to that ISA. To that end, we have developed an effective front-end tool (Smten) that embeds the SRI formal analysis tools (PVS, SAL, Yices) into the Bluespec hardware build chain, considerably simplifying the analysis effort. We hope that, by the completion of the CTSRD project, we will also be able to prove a number of basic but interesting properties about the hardware design, such as correctness of pipelining and the capability coprocessor. Perhaps we may even be able to extend the formal analysis into the lower-layer system software – such as properties relating to capability-based protection in sandboxing and compiling.

### **7.3 ABI and Compiler Development**

We have targeted our CHERI ISA extensions at compiler writers, rather than for direct use by application authors. This has required us to design new Application Binary Interfaces (ABIs), and to extend the C programming language to allow specification of protection properties by programmers. We have extended the GNU assembler and the Clang/LLVM compiler suite to generate CHERI instructions, and begun to experiment with modifications to applications. We anticipate significant future work in this area to validate our current approach, but also to extend these ideas both in C and other programming languages, such as Objective C. We are also interested in CHERI instructions as a target for just-in-time compilation by systems such as Dalvik.

### **7.4 Hardware Capability Support for FreeBSD**

With a capability processor prototype complete, and a FreeBSD/BERI port up and running, we have begun an investigation into adding CHERI capability support to the operating system. Currently, the CheriBSD kernel is able to maintain additional per-thread CHERI state for user processes via minor extensions to the process and thread structures, as well as exception-handling code. We have also prototyped object-capability invocation, which we are in the process of integrating with the operating system. A number of further tasks remain, including continuing to explore userspace memory models and their relationship with the userspace object-capability model, and also continued development of the userspace address-space executive that manages memory, code linking, and so on. This work depends heavily on continuing work on Clang/LLVM support for capabilities.

We need to explore security semantics for the kernel to limit access to kernel services (especially system calls) from sandboxed userspace code. This will require developing our notions of privilege described in Chapter 2; the userspace runtime and kernel must agree on which services (if any) are available without passing through a trusted protected subsystem, such as the runtime linker.

Ideally, the kernel should make use of capabilities, initially for bounded memory buffers (offering protection against kernel buffer overflows, for example), but later protected subsys-

tems. An iterative refinement of hardware and software privilege models will be required: for example, a sandboxed kernel subsystem should not be able to modify the TLB without going through a kernel protected subsystem, meaning that simple ring-based notions of privilege for MMU access are insufficient.

## 7.5 Evaluating Performance and Programmability

This report describes a fundamental premise: that through an in-address space capability model, performance and programmability for compartmentalized applications can be dramatically improved. Once the capability coprocessor and initial programming language, toolchain, and operating system support come together, validating this claim will be critical. We anticipate making early efforts to apply compartmentalization to base system components: elements of the operating system kernel, critical userspace libraries, and critical userspace applications.

Our hybrid capability architecture will ease this experimentation, making it possible to apply, for example, capabilities within `zlib` without modifying an application as a whole. Similarly, capability-aware applications should be able to invoke existing library services, even filtering their access to OS services – a similarly desirable hypothesis to test.

We are concerned not only with whether we can express the desired security properties, but also compare their performance with MMU-based compartmentalization, such as that developed in the Capsicum project. An early element of this work will certainly include testing of security context-switch speed as the number of security domains increases, in order to confirm our hypothesis regarding TLB size and highly compartmentalized software, but also that capability context switching can be made orders of magnitude faster as software size scales.

# Bibliography

- [1] M. Accetta, R. Baron, D. Golub, R. Rashid, A. Tevanian, and M. Young. Mach: A New Kernel Foundation for UNIX Development. Technical report, Computer Science Department, Carnegie Mellon University, August 1986.
- [2] W. B. Ackerman and W. W. Plummer. An implementation of a multiprocessing computer system. In *SOSP '67: Proceedings of the First ACM Symposium on Operating System Principles*, pages 5.1–5.10, New York, NY, USA, 1967. ACM.
- [3] J. Anderson. Computer security technology planning study. Technical Report ESD-TR-73-51, U.S. Air Force Electronic Systems Division, October 1972. (Two volumes).
- [4] G. R. Andrews. Partitions and principles for secure operating systems. Technical report, Cornell University, Ithaca, NY, USA, 1975.
- [5] Apple Inc. Mac OS X Snow Leopard. <http://www.apple.com/macosx/>, 2010.
- [6] B. N. Bershad, S. Savage, P. Pardyak, E. G. Sirer, M. E. Fiuczynski, D. Becker, C. Chambers, and S. Eggers. Extensibility safety and performance in the SPIN operating system. In *SOSP '95: Proceedings of the Fifteenth ACM Symposium on Operating Systems Principles*, pages 267–283, New York, NY, USA, 1995. ACM.
- [7] R. Bisbey II and D. Hollingworth. Protection Analysis: Project final report. Technical report, USC Information Sciences Institute (ISI), Marina Del Rey, California, 1978.
- [8] A. Bittau, P. Marchenko, M. Handley, and B. Karp. Wedge: Splitting Applications into Reduced-Privilege Compartments. In *Proceedings of the 5th USENIX Symposium on Networked Systems Design and Implementation*, pages 309–322. USENIX Association, 2008.
- [9] Bluespec, Inc., Waltham, MA. *Bluespec SystemVerilog Version 3.8 Reference Guide*, November 2004.
- [10] M. Branstad and J. Landauer. Assurance for the Trusted Mach operating system. In *Proceedings of the Fourth Annual Conference on Computer Assurance COMPASS '89*, pages 9–13. IEEE, June 1989.
- [11] B. M. Cantrill, M. W. Shapiro, and A. H. Leventhal. Dynamic instrumentation of production systems. In *ATEC '04: Proceedings of the USENIX Annual Technical Conference*, Berkeley, CA, USA, 2004. USENIX Association.

- [12] E. Cohen and D. Jefferson. Protection of the Hydra operating system. In *Proceedings of the Fifth ACM Symposium on Operating Systems Principles*, pages 141–160, 1975.
- [13] F. J. Corbató, M. Merwin-Daggett, and R. C. Daley. An experimental time-sharing system. In *AIEE-IRE '62 (Spring): Proceedings of the May 1–3, 1962, Spring Joint Computer Conference*, pages 335–344, New York, NY, USA, 1962. ACM.
- [14] F. J. Corbató and V. A. Vyssotsky. Introduction and overview of the Multics system. In *AFIPS '65 (Fall, part I): Proceedings of the November 30–December 1, 1965, Fall Joint Computer Conference, part I*, pages 185–196, New York, NY, USA, 1965. ACM.
- [15] R. C. Daley and P. G. Neumann. A general-purpose file system for secondary storage. In *AFIPS Conference Proceedings, Fall Joint Computer Conference*, pages 213–229. Spartan Books, November 1965.
- [16] P. J. Denning. Fault tolerant operating systems. *ACM Computing Surveys*, 8(4):359–389, 1976.
- [17] J. B. Dennis and E. C. Van Horn. Programming semantics for multiprogrammed computations. *Communications of the ACM*, 9(3):143–155, 1966.
- [18] J. Devietti, C. Blundell, M. M. K. Martin, and S. Zdancewic. Hardbound: Architectural support for spatial safety of the C programming language. *SIGPLAN Not.*, 43(3):103–114, Mar. 2008.
- [19] P. Efstathopoulos, M. Krohn, S. VanDeBogart, C. Frey, D. Ziegler, E. Kohler, D. Mazières, F. Kaashoek, and R. Morris. Labels and event processes in the asbestos operating system. *SIGOPS Oper. Syst. Rev.*, 39:17–30, October 2005.
- [20] R. S. Fabry. The case for capability based computers (extended abstract). In *SOSP '73: Proceedings of the Fourth ACM Symposium on Operating System Principles*, page 120, New York, NY, USA, 1973. ACM.
- [21] R. J. Feiertag and P. G. Neumann. The foundations of a Provably Secure Operating System (PSOS). In *Proceedings of the National Computer Conference*, pages 329–334. AFIPS Press, 1979. <http://www.csl.sri.com/neumann/psos.pdf>.
- [22] P. G. Neumann, R. Boyer, R. Feiertag, K. Levitt, and L. Robinson. A Provably Secure Operating System: The system, its applications, and proofs. Technical Report CSL-116, Second edition, Computer Science Laboratory, SRI International, Menlo Park, California, May 1980.
- [23] L. Gong. *Inside Java(TM) 2 Platform Security: Architecture, API Design, and Implementation*. Addison-Wesley, Reading, Massachusetts, 1999.
- [24] L. Gong, M. Mueller, H. Prafullchandra, and R. Schemers. Going beyond the sandbox: An overview of the new security architecture in the Java Development Kit 1.2. In *Proceedings of the USENIX Symposium on Internet Technologies and Systems*, Monterey, California, December 1997.



- [25] J. Gosling, B. Joy, and G. L. Steele. *The Java Language Specification*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 1996.
- [26] R. Graham. Protection in an information processing utility. *Communications of the ACM*, 11(5), May 1968.
- [27] N. Hardy. KeyKOS architecture. *SIGOPS Operating Systems Review*, 19(4):8–25, 1985.
- [28] T. Jim, J. G. Morrisett, D. Grossman, M. W. Hicks, J. Cheney, and Y. Wang. Cyclone: A safe dialect of C. In *ATEC '02: Proceedings of the USENIX Annual Technical Conference*, pages 275–288, Berkeley, CA, USA, 2002. USENIX Association.
- [29] A. Jones and W. Wulf. Towards the design of secure systems. In *Protection in Operating Systems, Proceedings of the International Workshop on Protection in Operating Systems*, pages 121–135, Rocquencourt, Le Chesnay, France, 13–14 August 1974. Institut de Recherche d’Informatique.
- [30] P. Karger. Limiting the damage potential of discretionary Trojan horses. In *Proceedings of the 1987 Symposium on Security and Privacy*, pages 32–37, Oakland, California, April 1987. IEEE Computer Society.
- [31] P. Karger and R. Schell. Multics security evaluation: Vulnerability analysis. In *Proceedings of the 18th Annual Computer Security Applications Conference (ACSAC), Classic Papers section*, Las Vegas, Nevada, December 2002. Originally available as U.S. Air Force report ESD-TR-74-193, Vol. II, Hanscomb Air Force Base, Massachusetts.
- [32] P. A. Karger. Using registers to optimize cross-domain call performance. *SIGARCH Computer Architecture News*, 17(2):194–204, 1989.
- [33] G. Klein, J. Andronick, K. Elphinstone, G. Heiser, D. Cock, P. Derrin, D. Elkaduwe, K. Engelhardt, R. Kolanski, M. Norrish, T. Sewell, H. Tuch, and S. Winwood. sel4: formal verification of an operating-system kernel. *Communications of the ACM*, 53:107–115, June 2009.
- [34] A. Kwon, U. Dhawan, J. M. Smith, T. F. Knight, Jr., and A. DeHon. Low-fat pointers: Compact encoding and efficient gate-level implementation of fat pointers for spatial safety and capability-based security. In *20th ACM Conference on Computer and Communications Security*, November 2013.
- [35] B. Lampson. Redundancy and robustness in memory protection. In *Information Processing 74 (Proceedings of the IFIP Congress 1974)*, volume Hardware II, pages 128–132. North-Holland, Amsterdam, 1974.
- [36] B. W. Lampson. Dynamic protection structures. In *AFIPS '69 (Fall): Proceedings of the November 18-20, 1969, Fall Joint Computer Conference*, pages 27–38, New York, NY, USA, 1969. ACM.
- [37] B. W. Lampson. Protection. *SIGOPS Operating Systems Review*, 8(1):18–24, 1974.

- [38] C. Lattner and V. Adve. LLVM: A compilation framework for lifelong program analysis and transformation. In *Proceedings of the international symposium on code generation and optimization: feedback-directed and runtime optimization*, CGO '04, pages 75–, Washington, DC, USA, 2004. IEEE Computer Society.
- [39] N. G. Leveson and W. Young. An integrated approach to safety and security based on system theory. *Communications of the ACM*, 57(2):31–35, February 2014. url-<http://www.csl.sri.com/neumann/insiderisks.html>.
- [40] H. M. Levy. *Capability-Based Computer Systems*. Butterworth-Heinemann, Newton, MA, USA, 1984.
- [41] J. Liedtke. On microkernel construction. In *SOSP'95: Proceedings of the 15th ACM Symposium on Operating System Principles*, Copper Mountain Resort, CO, Dec. 1995.
- [42] S. B. Lipner, W. A. Wulf, R. R. Schell, G. J. Popek, P. G. Neumann, C. Weissman, and T. A. Linden. Security kernels. In *AFIPS '74: Proceedings of the May 6-10, 1974, National Computer Conference and Exposition*, pages 973–980, New York, NY, USA, 1974. ACM.
- [43] S. McCanne and V. Jacobson. The BSD packet filter: a new architecture for user-level packet capture. In *USENIX'93: Proceedings of the USENIX Winter 1993 Conference*, Berkeley, CA, USA, 1993. USENIX Association.
- [44] E. McCauley and P. Drongowski. KSOS: The design of a secure operating system. In *National Computer Conference*, pages 345–353. AFIPS Conference Proceedings, 1979. Vol. 48.
- [45] M. K. McKusick and G. V. Neville-Neil. *The Design and Implementation of the FreeBSD Operating System*. Pearson Education, 2004.
- [46] A. Mettler and D. Wagner. Class properties for security review in an object-capability subset of Java. In *PLAS '10: Proceedings of the 5th ACM SIGPLAN Workshop on Programming Languages and Analysis for Security*, pages 1–7, New York, NY, USA, 2010. ACM.
- [47] M. S. Miller. *Robust composition: towards a unified approach to access control and concurrency control*. PhD thesis, Johns Hopkins University, Baltimore, MD, USA, 2006.
- [48] M. S. Miller, M. Samuel, B. Laurie, I. Awad, and M. Stay. Caja: Safe active content in sanitized javascript, May 2008. <http://google-caja.googlecode.com/files/caja-spec-2008-06-07.pdf>.
- [49] J. H. Morris, Jr. Protection in programming languages. *Communications of the ACM*, 16(1):15–21, 1973.
- [50] A. C. Myers and B. Liskov. A decentralized model for information flow control. *SIGOPS Oper. Syst. Rev.*, 31:129–142, October 1997.

- [51] S. Nagarakatte, J. Zhao, M. M. Martin, and S. Zdancewic. Softbound: Highly compatible and complete spatial memory safety for C. In *Proceedings of the 2009 ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI '09*, pages 245–258, New York, NY, USA, 2009. ACM.
- [52] G. C. Necula and P. Lee. Safe kernel extensions without run-time checking. In *OSDI '96: Proceedings of the Second USENIX symposium on Operating Systems Design and Implementation*, pages 229–243, New York, NY, USA, 1996. ACM.
- [53] G. C. Necula, S. McPeak, and W. Weimer. Ccured: Type-safe retrofitting of legacy code. In *Proceedings of the 29th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL '02*, pages 128–139, New York, NY, USA, 2002. ACM.
- [54] P. G. Neumann. Holistic systems. *ACM Software Engineering Notes*, 31(6):4–5, November 2006.
- [55] P. G. Neumann, R. Boyer, R. Feiertag, K. Levitt, and L. Robinson. A Provably Secure Operating System: The system, its applications, and proofs. Technical report, Computer Science Laboratory, SRI International, Menlo Park, California, May 1980. 2nd edition, Report CSL-116.
- [56] P. G. Neumann and R. J. Feiertag. PSOS revisited. In *Proceedings of the 19th Annual Computer Security Applications Conference (ACSAC 2003), Classic Papers section*, pages 208–216, Las Vegas, Nevada, December 2003. IEEE Computer Society. <http://www.acsac.org/> and <http://www.csl.sri.com/neumann/psos03.pdf>.
- [57] P. G. Neumann and R. N. Watson. Capabilities Revisited: A Holistic Approach to Bottom-to-Top Assurance of Trustworthy Systems. In *Proceedings of the Fourth Annual Layered Assurance Workshop*, December 2010.
- [58] E. Organick. *The Multics System: An Examination of Its Structure*. MIT Press, Cambridge, Massachusetts, 1972.
- [59] D. A. Patterson and C. H. Sequin. RISC I: A Reduced Instruction Set VLSI Computer. In *ISCA '81: Proceedings of the 8th Annual Symposium on Computer Architecture*, pages 443–457, Los Alamitos, CA, USA, 1981. IEEE Computer Society Press.
- [60] N. Provos, M. Friedl, and P. Honeyman. Preventing Privilege Escalation. In *Proceedings of the 12th USENIX Security Symposium*. USENIX Association, 2003.
- [61] R. Rashid and G. Robertson. Accent: A communications oriented network operating system kernel. In *Proceedings of the Eighth ACM Symposium on Operating System Principles*, pages 64–75, Asilomar, California, December 1981. (ACM Operating Systems Review, Vol. 15, No. 5).
- [62] C. Reis and S. D. Gribble. Isolating web programs in modern browser architectures. In *EuroSys '09: Proceedings of the 4th ACM European Conference on Computer Systems*, pages 219–232, New York, NY, USA, 2009. ACM.

- [63] D. Rémy and J. Vouillon. Objective ML: a simple object-oriented extension of ML. In *POPL '97: Proceedings of the 24th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 40–53, New York, NY, USA, 1997. ACM.
- [64] D. M. Ritchie and K. Thompson. The UNIX time-sharing system. *Communications of the ACM*, 17(7):365–375, 1974.
- [65] Ruby Users Group. Ruby Programming Language. <http://www.ruby-lang.org/>, October 2010.
- [66] J. Rushby. The design and verification of secure systems. In *Proceedings of the Eighth ACM Symposium on Operating System Principles*, pages 12–21, Asilomar, California, December 1981. (ACM Operating Systems Review, 15(5)).
- [67] J. Saltzer. Protection and the control of information sharing in Multics. *Communications of the ACM*, 17(7):388–402, July 1974.
- [68] J. Saltzer and M. Schroeder. The protection of information in computer systems. *Proceedings of the IEEE*, 63(9):1278–1308, September 1975.
- [69] W. Schiller. The design and specification of a security kernel for the PDP-11/45. Technical Report MTR-2934, Mitre Corporation, Bedford, Massachusetts, March 1975.
- [70] M. D. Schroeder. Engineering a security kernel for Multics. In *SOSP '75: Proceedings of the Fifth ACM Symposium on Operating Systems Principles*, pages 25–32, New York, NY, USA, 1975. ACM.
- [71] E. J. Sebes. Overview of the architecture of Distributed Trusted Mach. In *Proceedings of the USENIX Mach Symposium*, pages 20–22. USENIX Association, November 1991.
- [72] J. Shapiro, J. Smith, and D. Farber. EROS: a fast capability system. In *SOSP '99: Proceedings of the seventeenth ACM Symposium on Operating Systems Principles*, Dec 1999.
- [73] R. Spencer, S. Smalley, P. Loscocco, M. Hibler, D. Andersen, and J. Lepreau. The Flask security architecture: System support for diverse security policies. In *Proceedings of the 8th USENIX Security Symposium*, pages 123–139, Washington, D.C., USA, Aug. 1999. USENIX Association.
- [74] R. Wahbe, S. Lucco, T. E. Anderson, and S. u. L. Graham. Efficient software-based fault isolation. In *SOSP '93: Proceedings of the Fourteenth ACM Symposium on Operating Systems Principles*, pages 203–216, New York, NY, USA, 1993. ACM.
- [75] B. J. Walker, R. A. Kemmerer, and G. J. Popek. Specification and verification of the UCLA Unix security kernel. *Communications of the ACM*, 23(2):118–131, 1980.
- [76] S. T. Walker. The advent of trusted computer operating systems. In *AFIPS '80: Proceedings of the May 19-22, 1980, national computer conference*, pages 655–665, New York, NY, USA, 1980. ACM.

- [77] H. J. Wang, C. Grier, A. Moshchuk, S. T. King, P. Choudhury, and H. Venter. The multi-principal OS construction of the Gazelle web browser. In *Proceedings of the 18th USENIX Security Symposium*, pages 417–432, Berkeley, CA, USA, 2009. USENIX Association.
- [78] R. N. Watson, P. G. N. J. Woodruff, J. Anderson, R. Anderson, N. Dave, B. Laurie, S. W. Moore, S. J. Murdoch, P. Paeps, M. Roe, and H. Saidi. CHERI: a research platform deconflating hardware virtualization and protection. In *Runtime Environments, Systems, Layering and Virtualized Environments (RESoLVE 2012)*, March 2012.
- [79] R. N. M. Watson. New Approaches to Operating System Security Extensibility. Technical report, Ph.D. Thesis, University of Cambridge, Cambridge, UK, October 2010.
- [80] R. N. M. Watson. A decade of OS access-control extensibility. *Commun. ACM*, 56(2), Feb. 2013.
- [81] R. N. M. Watson, J. Anderson, B. Laurie, and K. Kennaway. Capsicum: Practical capabilities for Unix. In *Proceedings of the 19th USENIX Security Symposium*. USENIX, August 2010.
- [82] R. N. M. Watson, D. Chisnall, B. Davis, W. Koszek, S. W. Moore, S. J. Murdoch, P. G. Neumann, and J. Woodruff. Bluespec Extensible RISC Implementation (BERI): Software Reference. Technical Report UCAM-CL-TR-853, University of Cambridge, Computer Laboratory, 15 JJ Thomson Avenue, Cambridge CB3 0FD, United Kingdom, phone +44 1223 763500, June 2014.
- [83] R. N. M. Watson, D. Chisnall, B. Davis, W. Koszek, S. W. Moore, S. J. Murdoch, P. G. Neumann, and J. Woodruff. Capability Hardware Enhanced RISC Instructions (CHERI): User’s guide. Technical Report UCAM-CL-TR-851, University of Cambridge, Computer Laboratory, 15 JJ Thomson Avenue, Cambridge CB3 0FD, United Kingdom, phone +44 1223 763500, June 2014.
- [84] R. N. M. Watson, P. G. Neumann, J. Woodruff, J. Anderson, D. Chisnall, B. Davis, B. Laurie, S. W. Moore, S. J. Murdoch, and M. Roe. Capability Hardware Enhanced RISC Instructions (CHERI): Instruction-Set Architecture. Technical Report UCAM-CL-TR-850, University of Cambridge, Computer Laboratory, 15 JJ Thomson Avenue, Cambridge CB3 0FD, United Kingdom, phone +44 1223 763500, June 2014.
- [85] R. N. M. Watson, J. Woodruff, D. Chisnall, B. Davis, W. Koszek, A. T. Marketos, S. W. Moore, S. J. Murdoch, P. G. Neumann, R. Norton, and M. Roe. Bluespec Extensible RISC Implementation (BERI): Hardware Reference. Technical Report UCAM-CL-TR-852, University of Cambridge, Computer Laboratory, 15 JJ Thomson Avenue, Cambridge CB3 0FD, United Kingdom, phone +44 1223 763500, June 2014.
- [86] M. Wilkes and R. Needham. *The Cambridge CAP Computer and Its Operating System*. Elsevier North Holland, New York, 1979.
- [87] J. Woodruff, R. N. M. Watson, D. Chisnall, S. W. Moore, J. Anderson, B. Davis, B. Laurie, P. G. Neumann, R. Norton, and M. Roe. The CHERI capability model: Revisiting RISC in an age of risk. In *Proceedings of the 41st International Symposium on Computer Architecture (ISCA 2014)*, June 2014.

- [88] W. Wulf, E. Cohen, W. Corwin, A. Jones, R. Levin, C. Pierson, and F. Pollack. HYDRA: the kernel of a multiprocessor operating system. *Communications of the ACM*, 17(6):337–345, 1974.
- [89] W. Wulf, R. Levin, and S. Harbison. *Hydra/C.mmp: An Experimental Computer System*. McGraw-Hill, New York, 1981.
- [90] B. Yee, D. Sehr, G. Dardyk, J. B. Chen, R. Muth, T. Ormandy, S. Okasaka, N. Narula, and N. Fullagar. Native client: A sandbox for portable, untrusted x86 native code. In *Proceedings of the 2009 30th IEEE Symposium on Security and Privacy*, pages 79–93, Washington, DC, USA, 2009. IEEE Computer Society.
- [91] B. Yee, D. Sehr, G. Dardyk, J. B. Chen, R. Muth, T. Ormandy, S. Okasaka, N. Narula, and N. Fullagar. Native client: A sandbox for portable, untrusted x86 native code. In *SP '09: Proceedings of the 2009 30th IEEE Symposium on Security and Privacy*, pages 79–93, Washington, DC, USA, 2009. IEEE Computer Society.
- [92] N. Zeldovich, S. Boyd-Wickizer, E. Kohler, and D. Mazières. Making information flow explicit in histar. In *Proceedings of the 7th symposium on Operating systems design and implementation, OSDI '06*, pages 263–278, Berkeley, CA, USA, 2006. USENIX Association.