



Capability Hardware Enhanced RISC Instructions: CHERI User's guide

Robert N.M. Watson, David Chisnall,
Brooks Davis, Wojciech Koszek,
Simon W. Moore, Steven J. Murdoch,
Peter G. Neumann, Jonathan Woodruff

April 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, David Chisnall, Brooks Davis,
Wojciech Koszek, Simon W. Moore, Steven J. Murdoch,
Peter G. Neumann, Jonathan Woodruff

SRI International is acknowledged as an additional copyright
holder

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

The *CHERI User's Guide* documents the software environment for the Capability Hardware Enhanced RISC Instructions (CHERI) prototype developed by SRI International and the University of Cambridge. The *User's Guide* is targeted at hardware and software developers working with capability-enhanced software. It describes the CheriBSD operating system, a version of the FreeBSD operating system that has been adapted to support userspace capability systems via the CHERI ISA, and the CHERI Clang/LLVM compiler suite. It also describes the earlier Deimos demonstration microkernel.

Acknowledgments

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

Ross J. Anderson	Jonathan Anderson	Gregory Chadwick	Nirav Dave
Khilan Gudka	Jong Hun Han	Alex Horsman	Alexandre Joannou
Asif Khan	Myron King	Ben Laurie	Patrick Lincoln
Anil Madhavapeddy	Ilias Marinos	A. Theodore Markettos	Ed Maste
Andrew Moore	Will Morland	Alan Mujumdar	Prashanth Mundkur
Robert Norton	Philip Paeps	Michael Roe	Colin Rothwell
John Rushby	Hassen Saidi	Hans Petter Selasky	Muhammad Shahbaz
Stacey Son	Richard Uhler	Philip Withnall	Bjoern Zeeb

The CHERI team wishes thank 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			

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

1	Introduction	7
1.1	Background	7
1.2	Getting BERIC	8
1.3	Licensing	8
1.4	Version History	8
1.5	Document Structure	9
2	CheriBSD	10
2.1	CheriBSD Modifications	10
2.2	Obtaining FreeBSD/BERI and CheriBSD Source Code	11
2.3	About FreeBSD/BERI and CheriBSD	11
2.4	Building CheriBSD	11
3	CHERI Clang/LLVM	14
3.1	Cross-Compiling for CHERI	14
3.2	Building the Assembler	14
3.3	Building the Compiler	14
3.4	Building a Complete SDK	15
3.5	Using Clang	16
3.6	Disassembling CHERI Binaries	16
3.7	Capability Extensions to C	17
3.7.1	Querying Reserved Capability Registers	18
3.7.2	Other Builtins	18
3.7.3	Inline Assembly	18
3.7.4	Const and Capabilities	18
3.7.5	Output-Only Capabilities	19
3.7.6	Capability Implicit Range Checking	19
3.7.7	Opaque Types	19
3.7.8	Object-Capability Invoke Calling Convention	20
3.7.9	Stack Spills and Safety	20
3.8	Assembly Extensions	21
3.8.1	Capability Move	21
3.8.2	Capability-Relative Floating-Point Loads and Stores	21

4	The Deimos Demonstration Operating System	22
4.1	Demonstration Narrative	22
4.2	Deimos Design and Implementation	23
4.2.1	Supervisor	23
4.2.2	Memory	25
4.2.3	CHERI-Aware ABI	25
4.3	Building Deimos	25
4.4	Conclusion	26

Chapter 1

Introduction

This is the *User's Guide* for the Capability Hardware Enhanced RISC Instructions (CHERI) prototype. The *User's Guide* complements the *CHERI Instruction-Set Architecture* (which specifies the CHERI ISA), the *BERI Hardware Reference* (which describes the hardware prototype), and the *BERI Software Reference* (which describes the BERI software development environment).

The *User's Guide* describes the CheriBSD operating-system prototype, a version of the Clang/LLVM/gas toolchain used with CHERI, and the Deimos demonstration microkernel. The Guide is intended to address the needs of hardware and software developers who are prototyping new hardware features or bringing up operating systems, language runtimes, and compilers on CHERI, rather than literal end users. Future iterations will continue to flesh out operational aspects of the CHERI processor.

1.1 Background

Capability Hardware Enhanced RISC Instructions (CHERI), developed by SRI International and the University of Cambridge, are security extensions for the 64-bit MIPS Instruction Set Architecture (ISA). The CHERI ISA provides direct processor support for fine-grained memory protection and scalable compartmentalization of (and within) system software and application software.

Whereas traditional CPU designs impose heavy performance and programmability penalties for employing compartmentalization, CHERI's ISA features support fast and easy application compartmentalization of C-language systems software. These improvements are made possible by integrated processor support for continuous memory protection and enforcement using memory capabilities and the object-capability model. Contemporary software trusted computing bases (TCBs) such as operating-system kernels and language runtimes are particularly interesting targets, as CHERI will allow us to improve their security and reliability – and, therefore, the security and reliability of applications built on top of those services.

Detailed information on the CHERI ISA and possible uses may be found in the *CHERI Instruction-Set Architecture*, including new coprocessor registers and instructions. The CHERI prototype is a reference implementation of the CHERI ISA, intended to help validate the approach through a complete system implementation. CHERI is based on the Bluespec Extensible RISC Implementation (BERI) FPGA soft core, and is implemented as an additional coprocessor. The distinction between BERI and CHERI is evolving; however, our hope is that BERI

will be a reusable platform across multiple research projects relating to the hardware-software interface. We have open sourced both our hardware prototypes and software support for them in the form of FreeBSD device drivers, extensions to the Clang/LLVM compiler suite, and CheriBSD, a version of FreeBSD enhanced to utilize the CHERI ISA for fine-grained memory protection and compartmentalization.

1.2 Getting BERI

We distribute the BERI prototype and software stack as open source via the BERI website:

<http://www.beri-cpu.org/>

1.3 Licensing

The BERI hardware design, simulated peripherals, and software tools are available under the BERI Hardware-Software License, a lightly modified version of the Apache Software License that takes into account hardware requirements.

We have released our extensions to the FreeBSD operating system to support BERI under a BSD license; initial support for BERI was included in FreeBSD 10.0, but further features will appear in FreeBSD 10.1. We have also released versions of FreeBSD and Clang/LLVM that support the CHERI ISA under a BSD license; these are distributed via GitHub.

We welcome contributions to the BERI project; however, we are only able to accept non-trivial changes when an individual or corporate contribution agreement has been signed. The BERI hardware-software license and contribution agreement may be found at:

<http://www.beri-open-systems.org/>

1.4 Version History

- 1.0 An initial version of the *CHERI User's Guide* documented the implementation status of the CHERI prototype, including the CHERI ISA and processor implementation, as well as user information on how to build, simulate, debug, test, and synthesize the prototype.
- 1.1 Minor refinements were made to the text and presentation of the document, with incremental updates to its descriptions of the SRI/Cambridge development and testing environments.
- 1.2 This version of the *CHERI User's Guide* followed an initial demonstration of CHERI synthesized for the Terasic tPad FPGA platform. The Guide contained significant updates on the usability of CHERI features, the build process, and debugging features such as CHERI's debug unit. A chapter was added on Deimos, a demonstration microkernel for the CHERI architecture.
- 1.3 The document was restructured into hardware prototype and software reference material. Information on the status of MIPS ISA implementation was updated and expanded, especially with respect to the MMU. Build dependencies were updated, as was information on the CHERI simulation environment. The distinction between BERI and CHERI was discussed in detail. The Altera development environment was described in its own chapter. A new chapter was added that detailed bus and device configuration and use of the Terasic tPad and DE4 boards, including the Terasic/Cambridge MTL touch screen display. New chapters were added on building and using

CheriBSD, as well as a chapter on FreeBSD device drivers on BERI/CHERI. A new chapter was added on cross-building and using the CHERI-modified Clang/LLVM suite, including C-language extensions for capabilities.

- 1.4 This version introduced improved Altera build and Bluespec simulation instructions. A number of additional C-language extensions that can be mapped into capability protections were introduced. FreeBSD build instructions were updated for changes to the FreeBSD cross-build system. Information on the CHERI2 prototype was added.
- 1.5 In this version of the *CHERI User's Guide*, several chapters describe the CHERI hardware prototype have been moved into a separate document, the *CHERI Platform Reference Manual*, leaving the User's Guide focused on software-facing activities.
- 1.6 This version updated the *CHERI User's Guide* for changes in the CheriBSD build including support for the CFI driver, incorporation of Subversion into the FreeBSD base tree, and non-root cross builds. It also added information on the `quartus_pgm` command, and made a number of minor clarifications and corrections throughout the document.
- 1.7 In this version, information on building and using FreeBSD/BERI was moved to the *BERI Software Reference*. A short chapter describing CheriBSD was retained, and updated to reflect a migration from Perforce to Github. Information on the CHERI Clang/LLVM compiler was updated to include new C-language extensions. CheriBSD build instructions were extended to support the the CHERI Clang/LLVM compiler. Information on CheriBSD extensions to FreeBSD was expanded.
- 1.8 - **UCAM-CL-TR-851** This version of the *CHERI User's Guide* was made available as a University of Cambridge Technical Report. Information on CheriBSD was updated to reflect enhanced support for sandboxing, signal handling, and testing. A variety of changes were made to reflect open sourcing of the BERI/CHERI processors and associated software stacks. New CHERI Clang built-in functions for clearing capability-pointer tags and querying the capability-program counter are now documented.

1.5 Document Structure

This document is an introduction and user manual for Version 1 of the Capability Hardware Enhanced RISC Instructions CPU prototype:

Chapter 2 describes CheriBSD.

Chapter 3 describes CHERI extensions to Clang/LLVM and the C programming language.

Chapter 4 describes Deimos, a prototype operating-system software stack for CHERI. This microkernel, which runs on CHERI in simulated and synthesized forms, illustrates some of CHERI's protection and compatibility properties.

Chapter 2

CheriBSD

FreeBSD/BERI is a port of the open-source FreeBSD operating system that extends support for the Bluespec Extensible RISC implementation (BERI). CheriBSD extends FreeBSD/BERI to implement memory protection and software compartmentalization features supported by the CHERI ISA. General crossbuild and use instructions for FreeBSD/BERI may be found in the *BERI Software Reference*. Procedures for building and using FreeBSD/BERI should entirely apply to CheriBSD, except as documented in this chapter.

2.1 CheriBSD Modifications

FreeBSD/BERI has been modified in the following ways to support CHERI's security features:

- Platform boot code has been extended to enable the capability co-processor.
- The per-thread PCB context structure has been extended to hold a saved capability register file, as well as a per-thread 'trusted stack' that tracks object-capability invocations to provide a reliable return path.
- Kernel context-switching code has been extended to save and restore the capability register file for userspace.
- The kernel's handling of user exceptions has been extended to provide additional debugging information when userspace protection faults occur.
- The kernel rejects attempts to perform system calls from user threads whose **PCC** (program-counter capability) register does not have the `CHERI_PERM_SYSCALL` user-defined permission bit, preventing sandboxes from directly invoking system services. They must instead invoke a system class that is authorized to invoke system calls. We also hope to introduce new system calls that are safe within sandboxes and are authorized using special user capabilities – e.g., user capabilities that represent kernel file descriptors directly, avoiding the need for interposition – similar to the behavior of Capsicum.
- The kernel implements `CCall` and `CReturn` fast exception handlers that unseal invoked object capabilities, push the caller state onto the trusted stack, and restore it on return. If a fault occurs in the invoked object, control is returned to the caller.
- The kernel delivers capability-coprocessor faults in userspace processes as signals, extending the signal trap frame to include capability registers. This allows userspace software (and, in particular, language runtimes) to catch and handle software protection faults.

- The kernel is extended to allow processes that implement sandboxing to export class, method, and object statistics.
- The `libprocstat(3)` library and `procstat(1)` command have been extended to inspect exported sandbox statistics.
- A new library, `libcheri(3)`, has been added to provide a sandbox API, and to implement a set of system-class objects that can be delegated to sandboxes. Currently, this consists of a singleton system object that provides the ability to print to `stdout`, and a file-descriptor class that allows delegation of individual kernel-provided file descriptors to sandboxes.
- A new library, `libc_cheri(3)`, has been added to provide core C-language APIs and services within sandboxes. This library is able to use the system and file-descriptor classes to provide access to APIs such as `printf()`.
- A new command-line tool, `cheritest`, implements test cases for a variety of capability-related functions including sandboxing; `cheritest` relies on `cheritest-helper.bin` to provide sandboxed code.
- A new command-line tool, `cheri_tcpdump`, implements sandboxed packet sniffing and parsing; `cheri_tcpdump` relies on `tcpdump-helper` to provide sandboxed code.
- A new library, `libz-cheri(3)`, implements compression routines with fine-grained memory protection.

2.2 Obtaining FreeBSD/BERI and CheriBSD Source Code

Source code for CheriBSD is maintained on GitHub in the following repository:

```
https://github.com/CTSRD-CHERI/cheribsd
```

The CheriBSD development tree is branched from the FreeBSD GitHub repository at:

```
https://github.com/freebsd/freebsd
```

Development takes place on the master branch, which will eventually become FreeBSD 11.x. CheriBSD may be retrieved from GitHub as follows:

```
$ git clone https://github.com/CTSRD-CHERI/cheribsd
```

2.3 About FreeBSD/BERI and CheriBSD

CheriBSD contains additions to FreeBSD/BERI to support the CHERI capability coprocessor. Table 2.1 contains a list of directories affected by CHERI-specific behavior. Table 2.2 lists CHERI-specific files in the common MIPS configuration directory.

2.4 Building CheriBSD

CheriBSD follows the same build instructions as those found in the *BERI Software Reference* chapter on building FreeBSD/BERI, substituting source code from the above Git repository, as well as pathnames and kernel names in build commands.

Filename	Description
bin/cheritest/ ctsrds/ lib/libc_cheri/ lib/libcheri/	Command-line utility exercising CHERI and CheriBSD features, including sandboxing CTSRD-project demo code In-sandbox C library/runtime Library implementing the CHERI sandbox API; the CHERI system class implementation
libexec/cheritest-helper/ libexec/tcpdump-helper/ sys/mips/cheri/	Sandboxed components for cheritest Sandboxed components for cheri_tcpdump CHERI-specific code: coprocessor 2 initialization and context management
lib/libz-cheri	Version of libz compiled with CHERI memory protection
usr.sbin/tcpdump/cheri_tcpdump	Version of tcpdump able to use CHERI sandboxing
lib/libprocstat/ usr.bin/procstat/	Extensions to this library allow procstat(1) to monitor libcheri sandboxes procstat(1) command extended to monitor libcheri sandboxes

Table 2.1: CheriBSD directories

Filename	Description
CHERI_DE4_MDROOT	CheriBSD kernel configuration to use a memory root file system on the Terasic DE4
CHERI_DE4_SDROOT	CheriBSD kernel configuration to use an SD Card root file system on the Terasic DE4
CHERI_SIM_MDROOT	CheriBSD kernel configuration to use a memory root file system while in simulation
CHERI_SIM_SDROOT	CheriBSD kernel configuration to use a simulated SD Card root file system

Table 2.2: CheriBSD files in `src/sys/mips/conf`.

Portions of the CheriBSD tree rely on our CHERI-aware Clang/LLVM extensions (described in Chapter 3). You must tell `make buildworld` where to find the extended Clang/LLVM by adding this argument to the build command line:

```
CHERI_CC=/path/to/cheri-unknown-freebsd-clang
```

Typically, this argument will be a pointer to the `Build/bin` directory in your CHERI Clang/LLVM build, or `bin` in the CHERI SDK.

Some utility and demonstration software is stored in the `ctsrds` and `tools/tools/atsect1` directories. They can be built with the world by adding the following to the `make buildworld` command line:

```
LOCAL_DIRS="ctsrds tools/tools/atsect1" \  
LOCAL_LIB_DIRS=ctsrds/lib \  
LOCAL_MTREE=ctsrds/ctsrds.mtree
```

You will similarly need to include the following lines for the `make installworld` target:

```
LOCAL_DIRS="ctsrds tools/tools/atsect1" \  
LOCAL_MTREE=ctsrds/ctsrds.mtree
```

Support for the capability coprocessor is an optionally compiled extension enabled using `cpu CPU_CHERI` in CheriBSD. Ensure that you have replaced a BERIC kernel configuration-file name with a similar CHERI name to ensure that `nocpu CPU_BERIC` and `cpu CPU_CHERI` lines have been used.

Chapter 3

CHERI Clang/LLVM

This chapter describes CHERI-specific modifications to the Clang/LLVM compiler suite and the GNU assembler, as well as our extensions to the C programming language to support explicit capability use.

3.1 Cross-Compiling for CHERI

For cross-compiling code that targets CHERI, we provide a modified LLVM back end and Clang front end for [Objective-]C[++]. The back end can generate CHERI assembly and object code from LLVM's intermediate representation (IR). The front end generates the IR from C family languages and supports some capability extensions to C.

For assembly language programming, we also provide a modified version of the GNU binutils (including the GNU assembler `gas`) that has support for the capability instructions. This approach is gradually being deprecated in favor of the LLVM integrated assembler, but is still used in a number of places including the CHERI test suite and CheriBSD kernel build.

3.2 Building the Assembler

To build the assembler, you will need to have Git installed. Check out the source code and build it like this:

```
$ git clone git://github.com/CTSRD-CHERI/binutils.git
$ cd binutils
$ ./configure --target=mips64 --disable-werror
$ make
```

3.3 Building the Compiler

To build the compiler, you will need to have Git, CMake, and Ninja installed. Check out the code and build like this:

```
$ git clone git://github.com/CTSRD-CHERI/llvm.git
$ cd llvm/tools
$ git clone git://github.com/CTSRD-CHERI/clang.git
$ cd ..
$ mkdir Build
```

```
$ cd Build
$ cmake -G Ninja -DCMAKE_BUILD_TYPE:STRING=Debug \
  -DBUILD_SHARED_LIBS:BOOL=ON \
  -DLLVM_DEFAULT_TARGET_TRIPLE:STRING=cheri-unknown-freebsd ..
$ ninja
```

Note: A recent version of CMake (at least 2.8.8) is required. Also, the `std*.h` files included with Clang/LLVM appear to be incompatible with those shipped in the FreeBSD base; they cannot be used for a CheriBSD crossbuild. The SDK build will automatically trim them, but manual builds of Clang/LLVM require those files to be deleted from your build tree:

```
$ rm lib/clang/3.4/include/std*
```

If these files are not deleted, the CheriBSD build may experience compiler errors relating to variable-argument functions and type. These errors will hopefully be resolved in the future by upstreamed improvements to the integrated headers so that they are appropriate for use in FreeBSD.

By default, Ninja will select a number of processes to run in parallel on the build based on your number of processors. You can increase or decrease this number with the `-j` flag. Building LLVM is somewhat memory intensive, with compilation steps taking around 300MB of RAM and linking steps taking 1-2GB, so you may wish to reduce this number if you have less than 1GB of RAM per core. If you are on a 32-bit system, you may want to pass the following option to CMake to build a release build with asserts, rather than a debug build:

```
-DLLVM_ENABLE_ASSERTIONS:BOOL:=ON
```

Attempting to Link a debug build of LLVM can run out of address space in the linker in a 32-bit system.

Building Clang also requires a recent version of `gcc`. To compile Clang with itself (or to compile the CTSRD-modified version of Clang with an unmodified clang) pass the following options to `cmake`:

```
-DCMAKE_C_COMPILER=clang -DCMAKE_CXX_COMPILER=clang++
```

You can run the LLVM test suite, including Cheri-related tests, using:

```
$ ninja check
```

3.4 Building a Complete SDK

On a FreeBSD host system, if you have access to the Cheri svn repository, you can run `cherilibs/trunk/tools/build.sdk.sh`. This should be run in a new directory, as it will check out and build both CheriBSD and Cheri/LLVM.

After running this script, you should have subdirectories for the various projects that must be built, and an `sdk` directory containing the SDK. The SDK includes all of the core FreeBSD libraries and headers along with all of the tools required to cross build C/C++ programs for Cheri.

3.5 Using Clang

Once you have built LLVM, you (mostly) have a working cross-compiler. You can generate CHERI assembly code from [Objective-C][++] source code with this command:

```
$ clang -S {source file} -target cheri-unknown-freebsd \
    -msoft-float
```

You can also generate native code directly:

```
$ clang --sysroot={cheribsd sysroot} {source file(s)} \
    -B{cheribsd sdk directory} \
    -target cheri-unknown-freebsd -msoft-float \
    -o {output executable}
```

The **-target** flag specifies CHERI as the architecture and FreeBSD as the platform. If you have built the CHERI SDK, you will have a `cheri-unknown-freebsd-clang` which you can use instead of the `clang` and the **-target** flag.

The **--sysroot** and **-B** flags tell the compiler where to look for various things. The first specifies where to search for headers (when compiling) and libraries (when linking). The second specifies where to search for other parts of the toolchain, specifically the linker and (if you're not using the integrated assembler) the assembler.

The **-msoft-float** flag ensures that, if your version of CHERI has no FPU, we will emit calls to emulated FPU functions rather than causing illegal instruction traps.

If you do a debug build of LLVM, then Clang will default to using the simple register allocator. To see significantly better code, add the following CFLAGS:

```
-mllvm -regalloc=greedy -O3 -mllvm -enable-mips-delay-filler
```

These flags will enable a better register allocator and will attempt to replace the `nops` in delay slots with instructions from before the branch – which also turns on the full set of LLVM optimizations. Note that not all of these optimizations are well tested with CHERI; a lower optimization level may be required for the generation of correct code.

3.6 Disassembling CHERI Binaries

Disassembly of some instructions is common during debugging. You can do this for individual instructions with the `llvm-mc` tool:

```
$ echo 0x48 0x02 0x08 0x02 | llvm-mc -disassemble - \
    -triple=cheri-unknown-freebsd
.section TEXT,__text,regular,pure_instructions
CGetType $2, $c1
```

This tool expects a string of hex bytes and will write out the corresponding assembly. To disassemble entire object code files, use the `llvm-objdump` tool:

```
$ llvm-objdump -disassemble -triple=cheri-unknown-freebsd \
    {something.o}
```


3.7 Capability Extensions to C

Clang predefines a `__capability` macro. If you want to make your code portable to non-CHERI platforms, then you can begin it with:

```
1 #if !__has_feature(capabilities) && !defined(__capability)
2 #   define __capability
3 #endif
```

You can then use the `__capability` qualifier on any pointer. This pointer will be treated as a capability with the following semantics:

- Casts to integers return the base address.
- Casts from integers are treated as C0-relative¹.
- Relative addressing will use the load/store via capability instructions
- Pointer increments will return a new capability, which may be invalid if it exceeds the upper bound of the original.
- Pointer decrements will always fail.

Clang also provides a number of built-in functions for accessing aspects of these semantics:

```
1 size_t __builtin_cheri_get_cap_base(__capability void *);
2 size_t __builtin_cheri_get_cap_length(__capability void*);
3 size_t __builtin_cheri_get_cap_perms(__capability void*);
4 _Bool __builtin_cheri_get_cap_tag(__capability void*);
5 size_t __builtin_cheri_get_cap_type(__capability void*);
6 _Bool __builtin_cheri_get_cap_unsealed(__capability void*);
7 __capability void* __builtin_cheri_and_cap_perms(__capability void*,
8     size_t);
9 __capability void* __builtin_cheri_clear_cap_tag(__capability void*)
10 ;
11 __capability void* __builtin_cheri_inc_cap_base(__capability void*,
12     size_t);
13 __capability void* __builtin_cheri_set_cap_length(__capability void
14 *, size_t);
15 __capability void* __builtin_cheri_set_cap_type(__capability void*,
16     size_t);
17 __capability void* __builtin_cheri_seal_cap_code(__capability void*)
18 ;
19 __capability void* __builtin_cheri_seal_cap_data(__capability void*,
20     __capability void*);
21 __capability void* __builtin_cheri_unseal_cap(__capability void*,
22     __capability void*);
```

Most of these aspects correspond directly to the relevant capability inspection / modification instructions.

¹This means that a `__capability void*` to `uintptr_t` to `__capability void*` round trip will work only if the capability references something inside the C0 address space

3.7.1 Querying Reserved Capability Registers

Several capability registers are set aside in the ISA and ABI for specific functions. These can be queried using the following builtins:

```
1 __capability void *__builtin_cheri_get_exception_program_counter_cap
  (void);
2 __capability void *__builtin_cheri_get_global_data_cap(void);
3 __capability void *__builtin_cheri_get_invoke_data_cap(void);
4 __capability void *__builtin_cheri_get_kernel_cap1(void);
5 __capability void *__builtin_cheri_get_kernel_cap2(void);
6 __capability void *__builtin_cheri_get_kernel_code_cap(void);
7 __capability void *__builtin_cheri_get_kernel_data_cap(void);
8 __capability void *__builtin_cheri_get_program_counter_cap(void);
```

Only builtins that query the global and invoke data capabilities will be available to regular userspace code.

3.7.2 Other Builtins

Clang provides builtins that correspond to the CGetCause and CSetCause instructions:

```
1 size_t __builtin_cheri_get_cause(void)
2 void __builtin_cheri_set_cause(size_t)
```

These builtins are useful only for writing kernel code, as they correspond to privileged instructions.

3.7.3 Inline Assembly

Clang adds a "C" constraint to inline assembly. This allows capability registers to be used as input and output for inline assembly blocks. For example:

```
1 __capability void* getC0(void)
2 {
3     __capability void *c0;
4     __asm__("cmovl_%0,_%$c0" : "+C" (c0));
5     return c0;
6 }
```

This will return the current C0 value.

3.7.4 Const and Capabilities

In C/C++, you can explicitly cast a **const** pointer to a non-**const** pointer. This casting is not supported for `__capability` pointers. For example, consider the following two `__capability` pointers:

```
1 __capability void *a = something;
2 const __capability void *b;
```

This action will implicitly clear the store and store-capability flags:

```
1 b = a;
```

So far, this casting is almost exactly like the C behavior, except that now you will get a hardware trap instead of just a compiler warning if you try to store through `b`. If you perform the assignment in the other direction, however, the 'const-ness' is still preserved:

```
1 | a = (__capability void)b;
```

Any attempts to store through `a` will cause a trap.

3.7.5 Output-Only Capabilities

Clang also now supports an `__output` qualifier on capabilities. These qualifiers are the twin of `const` capabilities: they can be written only to, not read. You can create a write-only capability by casting any other capability. The compiler will statically check attempts to write through an `__output` capability, and code that does so via an explicit cast will abort at run time.

```
1 | int readFail(__capability __output int *x)
2 | {
3 |     *x = 12;
4 |     (*x)++;
5 |     return *x;
6 | }
```

This code contains two bugs: two attempts to read through the pointer passed as a parameter, even though it is declared solely for output. Attempting to compile it will raise the following warnings:

```
writetest.c:4:2: error: write-only variable is not readable
    (*x)++;
    ~~~~
writetest.c:5:9: error: write-only variable is not readable
    return *x;
           ^^
```

You can call this function transparently, without an explicit cast:

```
1 | int caller(__capability int *x)
2 | {
3 |     return readFail(x);
4 | }
```

In this case, the caller will automatically modify the capability to remove load and load-capability permissions before passing it.

3.7.6 Capability Implicit Range Checking

At optimization level 1 and above, LLVM will attempt to automatically limit the range of capabilities. At `-O0`, the compiler will not transform the code at all. If static code flow analysis can determine that the capability was constructed from a global or a stack allocation, LLVM will initially begin with its range limited to the size of that allocation. This approach will also work for various heap-allocation functions, including any allocations that are the direct result of casting from the result of one of the standard C allocation functions. Note that this analysis is *not* interprocedural, and so will not work in the general case. It is primarily intended as a way to provide a small amount of extra checking.

3.7.7 Opaque Types

CHERI-Clang supports a new pragma, `opaque`, for linking a type to a key. This pragma is intended to provide enforcement of opaque types in C: no code outside of a compilation unit (where the key is

visible) can dereference the opaque pointers, even after casting. This pragma is intended to be used as follows:

```
1 // In a public header (one or more of the following):
2 typedef __capability struct foo* foo_c;
3 typedef struct foo* foo_t;
4
5 // In either the implementation file or a private header:
6
7 // If we're in private header, these should be __attribute__
8 // ((visibility("hidden"))), otherwise they should be static:
9 void *ptrKey;
10 __capability void *capKey;
11
12 struct foo
13 {
14     int a, b;
15 };
16
17 #pragma opaque foo_c capKey
18 #pragma opaque foo_t ptrKey
```

In an implementation file, any function returning a `foo_t` will have its value XORed with the value of `ptrKey` before return and any function returning a `foo_c` will have the value sealed with `capKey` before return. Similarly, any function receiving a `foo_c` or `foo_t` as an argument will have the inverse applied. Within the function, `foo_t` and `foo_c` behave exactly as they would without the pragma.

3.7.8 Object-Capability Invoke Calling Convention

By default, code generated for functions that take capability arguments or return capability values will conform to the normal CHERI calling convention. However, functions prototypes may be tagged with `__attribute__((cheri_ccall))` to indicate to the compiler that the first two capability arguments should be placed in C1 and C2 rather than the normal argument registers. This convention is currently used when the function being called will perform a `CCall` to another protection domain, which requires the object-capability code and data capabilities to be in those registers. Typically, the invoked function will be an assembly function, `cheri_invoke`, which saves and clears registers as required before invoking potentially untrustworthy code, and then properly restores register state following `CReturn`. End-user code will rarely employ this attribute directly; however, it will be of value to developers creating new sandboxing models.

3.7.9 Stack Spills and Safety

Functions declared with the `sensitive` attribute deal with some kind of sensitive data. C already provides tools (`volatile` in older versions and explicit memory operations in C11) to ensure that sensitive data is not left in areas of memory managed by the programmer. C does not provide a mechanism for making the same guarantee for stack spills.

The traditional way of avoiding this problem is to use `volatile` variables everywhere and explicitly zero them. However, this method effectively disables all compiler optimization for a function, purely to avoid the possibility of a stack spill. In general, however, higher optimization levels are more likely to reduce the need to ever spill registers to the stack. The correctness of such code is difficult to determine

in the presence of optimizations such as common subexpression elimination. These optimizations may accidentally store intermediate results on the stack, which could persist beyond the end of the function – for future reuse.

Marking a function as sensitive does not disable *any* optimizations. It does, however, ensure that any values spilled to the stack during the function are zeroed at the end. This zeroing is especially important for CHERI systems, where stack spills include capabilities. The pass does not attempt to destroy the entire capability; it simply writes a 64-bit zero value over the start. This write is enough to destroy the base address (without which the rest of the information in the capability is largely useless) and, more importantly, invalidate the tag. With the tag invalidated, the capability is rendered invalid.

Note that the current implementation may still spill values to the stack if they were left in a callee-save register across function calls and the callee spills them. To avoid this issue, the compiler will warn if you call a function that is not marked as sensitive from one that is. One possible approach would be to store all callee-save registers to the stack and invalidate them before every call, but that would impose a significant overhead. Instead, we aim to provide tools that allow programmers to write secure code and to decide when these trade-offs are appropriate.

In typical use, the overhead of this attribute is a few instructions in the function epilog (or epilogs, if it has multiple return paths). Of course, this overhead depends on the complexity of the function and the size of its working set. In general, it does improve significantly at higher optimization levels where the register allocator eliminates most, if not all, stack spills. Spill slots are often reused, so the number of invalidations required is often lower than the total number of spills over the program.

We do not invalidate the spills in the function prolog, because these spills are unlikely to contain any sensitive data. They contain the initial stack pointer, the globals pointer and the return address. Once we begin using capabilities for these, invalidating them may be desirable.

3.8 Assembly Extensions

The LLVM integrated assembler (used by default by Clang, unless `-no-integrated-as` is passed) provides some mnemonics for ease of assembly programming.

3.8.1 Capability Move

The `cmove` pseudo operation expands to a `CIncBase` instruction with `$zero` as the increment size. This is a shorthand for moving a value between capability registers. Its use is discouraged, as the similarity of this instruction and `cmov` makes code difficult to read.

3.8.2 Capability-Relative Floating-Point Loads and Stores

The assembler provides `clwc1`, `cldc1`, `cswc1`, and `csdc1` pseudoinstructions. These load and store 32- or 64-bit floating point values and are the same format as integer load and store operations. For example, to load a single-precision floating point value from offset 32 inside `C3` you would do:

```
1 | clwc1 $f2, $zero, 32($c2)
```

This will expand to:

```
1 | clw   $1, $zero, 32($c2)
2 | mtc1  $1, $f2
```

A future revision of the CHERI ISA may add instructions for loading and storing floating point values directly, at which point assembly code using this pseudo will generate a real instruction.

Chapter 4

The Deimos Demonstration Operating System

Deimos is a demonstration microkernel operating system that uses the CHERI ISA's capability features to sandbox untrustworthy applications. For the purposes of this demonstration, the CHERI prototype CPU was implemented in the Terasic tPad platform using an Altera FPGA; the tPad includes a VGA touchscreen, which is used for the Deimos user interface.

This chapter is largely historical in nature, reflecting our early experimentation with the CHERI ISA; Deimos as demonstrated in November 2011 is not able to run on current CHERI hardware due to ISA changes. It is also unable to take advantage of a number of contemporary CHERI features, such as access to a CHERI-aware compiler.

4.1 Demonstration Narrative

Figure 4.1 illustrates the demonstration user interface: two untrusted applications, a touchscreen drawing application and a weather status bar, as well as a trustworthy status bar across the top of the display maintained by the operating system. All applications run within the kernel ring and a single address space that is traditionally reserved only for privileged applications. Using CHERI's memory capability features, That address space is partitioned into a supervisor and two sandboxes. Access to I/O devices (such as touchscreen input and portions of the frame buffer) is delegated to the sandboxes using the capability mechanism to provide hardware-enforced control over display access. Attempts to draw outside of the delegated display area trigger a hardware exception, returning control to the supervisor.

In the demonstration, all touchscreen input is sent to the on-screen drawing application – drawing outside of its delegated area using a stylus causes the application to ‘misbehave’; attempting to write outside its screen region triggers an exception. The supervisor displays an exception indicator that shows when a disallowed access has been requested. The sandbox is then restarted on the next instruction, allowing it to continue. Over 80% of code in the Deimos kernel and applications is portable C and stock 64-bit MIPS assembly code, requiring only a very small amount of CHERI-specific assembly in order to manage additional register state, delegate capabilities, and perform I/O access via capabilities.

The demonstration shows off a number of features of the CHERI architecture:

- A hybrid capability model, allowing conventional C and MIPS assembly to coexist, transparently, with a hardware-enforced sandbox model.
- Selective delegation and containment of I/O access within the kernel ring, illustrating how device drivers (for example) can be contained using the capability model.



Figure 4.1: The Deimos touchscreen interface

- The flexibility of CHERI hybridization: the ability to run an operating system and applications using solely capability protection features, without employing the TLB.
- How CHERI features, combined with hardware user interface elements and software UI components, can be used to implement a trusted path.

4.2 Deimos Design and Implementation

Figure 4.2 illustrates the rough structure of the Deimos demonstration:

- 10,500 lines of Bluespec that implement a fully pipelined, 64-bit MIPS CPU, including a capability coprocessor as CP2.
- 2,950 lines of C and assembly code that implement a preemptive, multi-tasking microkernel, which uses capabilities for sandboxing of untrustworthy applications.
- 640 lines of C and assembly code that implement supporting libraries for sandboxed applications, including graphics and string management.
- 140 and 90 lines, respectively, of C and assembly code for two touchscreen applications.

4.2.1 Supervisor

Deimos uses only the capability coprocessor for sandboxing, eschewing use of virtual addressing. The supervisor runs entirely within exception-handling context – that is, with `EXL = 1`; a small bootstrap routine initializes the Deimos kernel stack and data structures, and then triggers an exception to enter the supervisor. The supervisor consists of an exception handler that implements a small microkernel. On an exception, the supervisor saves *userspace*'s general-purpose and capability registers to the *process*'s context structure; when the supervisor has completed execution, user context is restored – possibly after a context switch that selects a new process.

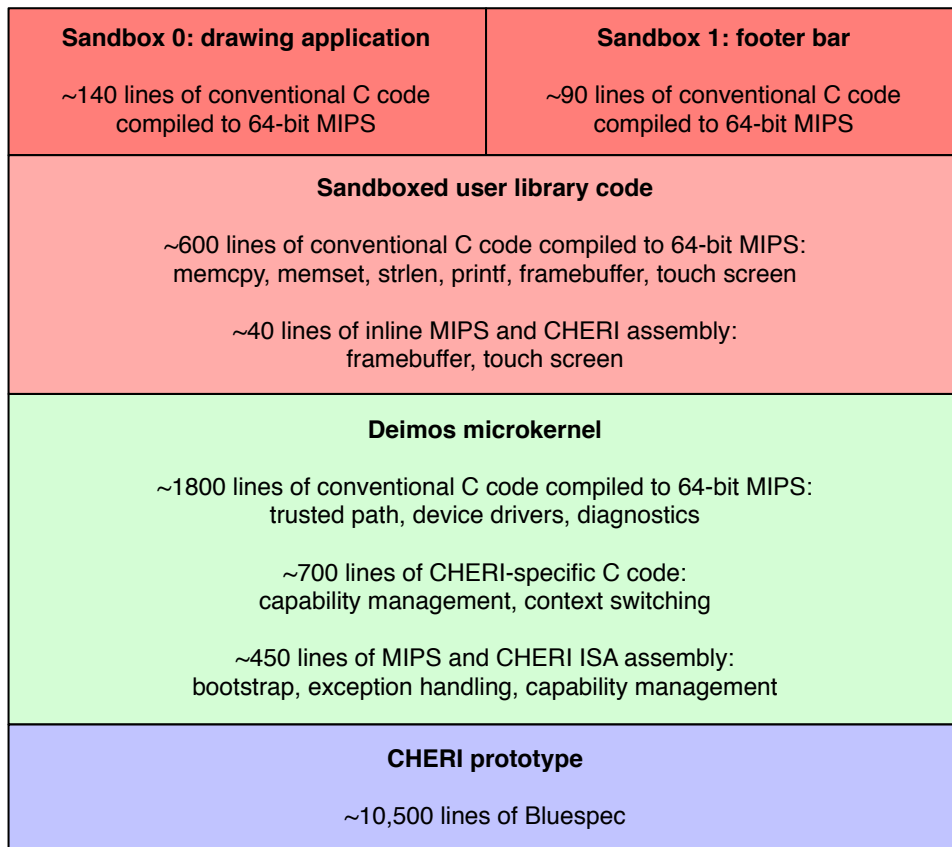


Figure 4.2: The Deimos software architecture: a blend of C code compiled to stock 64-bit MIPS machine code, stock 64-bit MIPS assembly, and a small amount of CHERI-specific assembly.

Target	Description
run	Build Deimos and run in the CHERI simulator
verilog	Build Deimos and generate a memory image for FPGA use
clean	Clean the Deimos build; the same options must be passed to clean as were used to build Deimos.

Table 4.1: Deimos make targets

Deimos handles a number of exception types, including userspace preemption using the MIPS cycle counter, explicit system calls from the user process, and capability exception handling faults. Deimos system calls include simple UART I/O routines, a voluntary yield call, and a call to retrieve capabilities using a string-based namespace; capabilities delegating access to different portions of user-addressable screen space, as well as touchscreen input, are available to processes, subject to a simple access control policy.

4.2.2 Memory

Deimos uses capability registers **PCC** and **EPC** to lay out per-sandbox address space with respect to conventional MIPS memory access instructions, providing a small dead space at address 0, code segment, heap, and stack. Using this technique, unmodified MIPS application code can be isolated transparently. Access to additional memory outside of the compatibility address space (such as to the frame buffer) is possible through delegated capabilities. Application code is linked to run at alternative addresses, and some care must be taken during context switch to ensure that general-purpose and capability registers are saved and restored in appropriate order to allow both user applications and the supervisor to run correctly – especially as most of the supervisor is implemented in C code compiled to stock MIPS machine code.

4.2.3 CHERI-Aware ABI

Each user process is described by context structures that hold preserved general-purpose and capability register context that is saved and restored when the supervisor is entered and exited. Deimos implements an enhanced Application Binary Interface (ABI) that allows capabilities to be passed to and returned from functions, including system calls. As with general-purpose registers, this calling convention designates specific capability registers for use as arguments, returned values, and temporary values. Most code in CHERI applications is not aware of the enhanced ABI and hence preserves capability registers; certain kernel and application code is aware of capability registers and manages them per these conventions. Application frame buffer library code, for example, queries frame buffer capabilities when starting, and then preserves them for later use.

4.3 Building Deimos

The Deimos source code may be found in the `deimos/trunk` tree in the CHERI Subversion repository. Table 4.1 shows a number of useful make targets for Deimos. Table 4.2 shows options that may be passed when building Deimos.

Option	Description
CP2=YES	Enable capability coprocessor support, sandboxing applications
PREEMPTION=YES	Enable cycle timer-based preemption
SIMULATOR=YES	Modify runtime timing parameters to improve performance in a simulated environment
TRACE=YES	Enable UART tracing of Deimos system events
UART_DISABLE_FC=YES	Disable flow control support for the Deimos console on the UART – useful for some FPGA targets where a remote endpoint may not be connected

Table 4.2: Deimos make options

4.4 Conclusion

Deimos is a simplistic operating system that implements an elementary process model and minimalist system services. However, it is sufficiently capable to illustrate a number of key features of the CHERI processor – especially, its support for a hybrid capability model in which some portions are explicitly aware of capabilities, but much is not. Deimos also illustrates how capability delegation can be used even within a privileged ring, to contain software components with delegation of limited I/O access. This ability suggests various applications of CHERI features that would be useful in conventional operating systems. We hope to build on Deimos for further demonstration and testing purposes, and use lessons learned from the experience of building Deimos to inform our larger-scale experiments with CHERI, such as the adaptation of the FreeBSD operating system to exploit capabilities within both its kernel and userspace.