# DSbD All Hands November 2023
## Tutorial: New features in CheriBSD 23.11

**Robert N. M. Watson**, Simon W. Moore, Peter Sewell, Peter G. Neumann, Brooks Davis

Hesham Almatary, Ricardo de Oliveira Almeida, Jonathan Anderson, Alasdair Armstrong, Rosie Baish, Peter Blandford-Baker, John Baldwin, Hadrien Barrel, Thomas Bauereiss, Ruslan Bukin, Brian Campbell, David Chisnall, **Jessica Clarke**, Nirav Dave, Lawrence Esswood, Nathaniel W. Filardo, Franz Fuchs, **Dapeng Gao**, Ivan Gomes-Ribeiro, Khilan Gudka, Brett Gutstein, Angus Hammond, Graeme Jenkinson, Alexandre Joannou, Mark Johnston, Robert Kovacsics, Ben Laurie, A. Theo Markettos, J. Edward Maste, **Alfredo Mazzinghi**, Alan Mujumdar, Prashanth Mundkur, Steven J. Murdoch, Edward Napierala, George Neville-Neil, Kyndylan Nienhuis, Robert Norton-Wright, Philip Paeps, Lucian Paul-Trifu, Allison Randal, Ivan Ribeiro, Alex Richardson, Michael Roe, Colin Rothwell, Peter Rugg, Hassen Saidi, Thomas Sewell, Stacey Son, Ian Stark, Domagoj Stolfa, Andrew Turner, Munraj Vadera, **Konrad Witaszczyk**, Jonathan Woodruff, Hongyan Xia, Vadim Zaliva, and Bjoern A. Zeeb

University of Cambridge and SRI International
Manchester, November 8, 2023

# Agenda

| Start: 2:00 PM | |
|---|---|
| **Part I: presentation**<br><br>**(~25 min)** | Overview |
| | Benchmark ABI |
| | Heap temporal memory safety |
| | Library compartmentalization |
| | Visualization tooling |
| **Part II: practice**<br><br>**(~35 min)** | Work environment setup (FreeBSD jails on CheriBSD/Morello hosts) |
| | Exercises |
| **End: 3:00 PM** | |
| **Round 2 workshop**: Measuring, analysing, and understanding performance on Morello<br><br>3:15 PM - 4:15 PM in International Suite (Conference end) | |

UNIVERSITY OF CAMBRIDGE

# Part I: presentation
New features in CheriBSD 23.11.

# CHERI prototype software stack on Morello

- **Complete open-source software stack** from bare metal up: compilers, toolchain, debuggers, hypervisor, OS, applications – all demonstrating CHERI

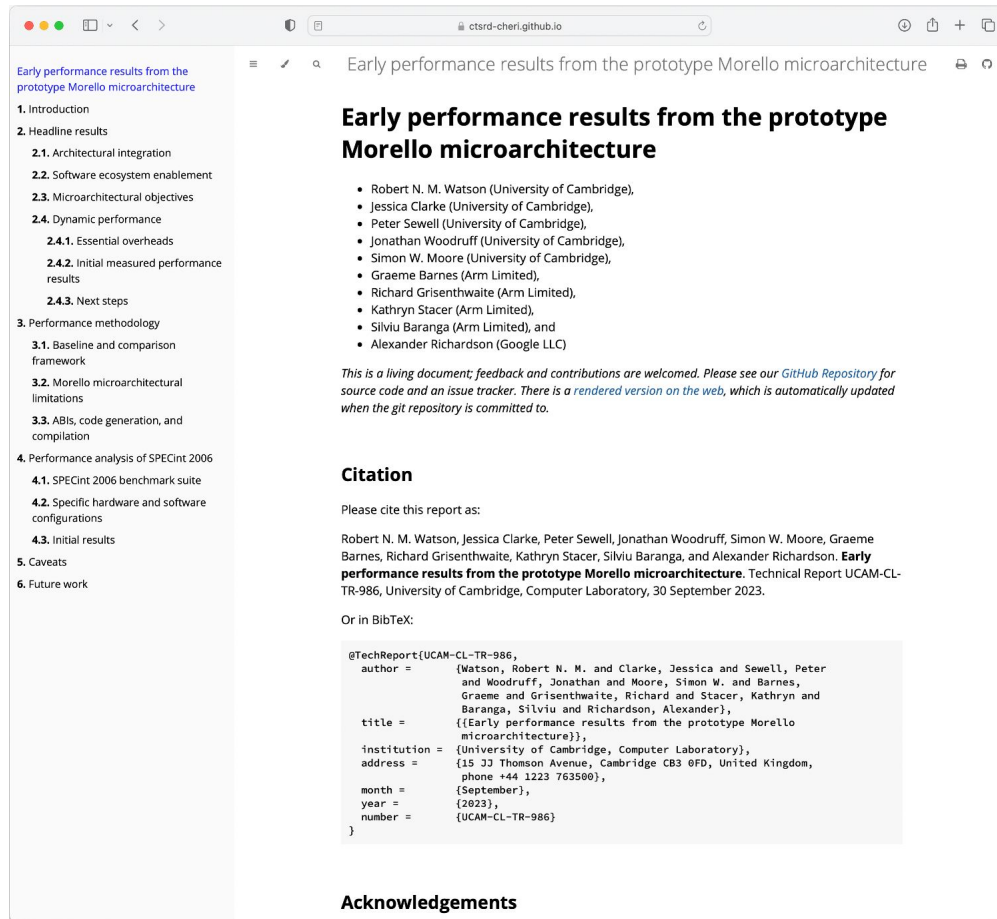- Rich CHERI feature use, but fundamentally incremental/hybridized deployment

**Open-source application suite** (KDE Plasma, Wayland, WebKit, Python, OpenSSH, nginx, …)

**CheriBSD/Morello** (funded by DARPA and UKRI)
(Morello and CHERI-RISC-V)

- FreeBSD kernel + userspace, application stack
- Kernel spatial and referential memory protection
- Userspace spatial, referential, and temporal memory protection
- Co-process compartmentalization (development branch)
- Linker-based compartmentalization
- Morello-enabled bhyve Type-2 hypervisor
- ARMv8-A 64-bit binary compatibility for legacy binaries

**Android** (Arm)
(Morello only)

**Linux** (Arm)
(Morello only)

Baseline CHERI Clang/LLVM from SRI/Cambridge; Morello adaptation by Arm + Linaro

**CHERI Clang/LLVM compiler suite, Morello GCC, LLD, LLDB, GDB**

SRI International

UNIVERSITY OF CAMBRIDGE

# Maturing CHERI software artifacts

| Feature | Status | Availability |
|---|---|---|
| 3rd-party packages (Hybrid) | **23K memory-unsafe** software packages with strong functionality expectations | Since May 2022 (22.05 release) |
| 3rd-party packages (CheriABI) | **11K memory-safe** software packages with mixed functionality expectations | Since May 2022 (22.05 release) **Up from 9k packages in 23.11** |
| Morello GPU device drivers | Memory-safe kernel and user drivers, | Since December 2022 (22.12 release) |
| Benchmark ABI support (+3rd-party packages) | Support for modified code generation addressing Morello bounds prediction | **Shipping in 23.11** (roughly the same packages as CheriABI) |
| Userlevel heap temporal safety | Prototype implements strong temporal safety, developed with Microsoft; testing required | **Shipping in 23.11 (pretty experimental)** |
| Linker-based compartmentalization | Introduces strong encapsulation boundaries around UNIX libraries with no modification | Since 22.12 (very experimental); **Significant improvements in 23.11** |
| bhyve (Type-2) hypervisor | Prototype boots pure-capability guest OS, validation required | **Shipping in 23.11 (very experimental)** |
| Co-process compartmentalization | Prototype runs some compartmentalized software (e.g., OpenSSL); API co-design | Planning to ship in 2024 |

SRI International

UNIVERSITY OF CAMBRIDGE

# CheriBSD 23.11

- The release is expected to be published in the next two weeks – final integration and testing are in progress.

- The release and this workshop are focused on CheriBSD on Arm Morello, but the release is continuously tested not to break existing functionality on CHERI-RISC-V.

- NB: This tutorial is running the 2023-11-01 'dev' snapshot, with a patch to the c18n run-time linker.

# Early Performance Results from the Prototype Morello Microarchitecture



https://ctsrd-cheri.github.io/morello-early-performance-results/

- Presents first detailed performance analysis on the first-generation Morello prototype
- Identifies and mitigates/corrects multiple sources of overhead arising from microarchitectural properties arising from a research design process and timeline
  - "Benchmark ABI" works around lack of bounds prediction in prototype
  - Morello design on FPGA used to validate other incremental improvements
  - Current best estimate of overhead ceiling presented using "P128" code generation
- "Essential reading" if you are doing any performance work on Morello
- Covered in the next workshop session

# Benchmark ABI

- The performance evaluation report shows that capability-relative jumps introduce overhead on Arm Morello in some workloads
  - The overhead is associated with the branch-predictor that was not extended to predict bounds in the Morello prototype
  - This is a Morello prototype performance overhead, not a CHERI overhead – and can be worked around in code generation to better predict future CHERI performance on a mature microarchitecture
- We introduce a new Benchmark ABI that caters to this concern
  - The CheriBSD kernel is aware of the benchmark ABI, and constructs an unbounded PCC for a user-space process
  - A benchmark ABI program is compiled not to use capability-based jumps, leaving PCC unbounded
  - Should not be used for security evaluation; instead use CheriABI compilation that implements bounded PCC

SRI International

UNIVERSITY OF CAMBRIDGE

# Cornucopia Reloaded: Load Barriers for CHERI Heap Temporal Safety (ASPLOS 2024)



- **Cornucopia heap temporal safety** (IEEE SSP 2020), is a GC-inspired, quarantining technique
  - The kernel virtual-memory subsystem tracks "capability dirty" pages
  - A long "stop-the-world" phase - as much as 30 milliseconds measured in practice
- **Cornucopia Reloaded** (ASPLOS 2024) moves to a GC-inspired "load-barrier"
  - VM invariant is that accessible pages have already undergone revocation
  - Depend on 1-bit capability generation added to VM PTEs, implemented by Morello
  - Stop-the-world pauses 10s of microseconds
- Enabled by default in CheriBSD 23.11

# Cornucopia Reloaded in CheriBSD 23.11

- CheriBSD 23.11 will include heap temporal safety enabled by default

  ○ You can enable, disable it at the system, binary file or process level

- The mrs(3) (Malloc Revocation Shim) man page describes available utilities
  ○ May also be usable to extend [some] other heap allocators..?

- The default memory allocator in CheriBSD 23.11 is jemalloc

- We are testing it with the GUI stack and other third-party software

# Ongoing temporal memory-safety deployment

Looking for increased experience, and appreciate any feedback on:

- Semantic impact on any applications vs. bugs/vulnerabilities discovered

- Acceptability of performance behavior, optimization opportunities

- Use in higher-level allocators – e.g., APR, Chromium, etc.

- Support for strong isolation needed for compartmentalization

- Enabling safe inter-compartment communication via shared memory

# Proposed operational models:
# Isolated libraries and UNIX co-processes

## Isolated dynamically linked libraries

- New API loads libraries into in-process sandboxes.
- Calling functions in isolated libraries performs a domain transition, with overheads comparable to function calls.
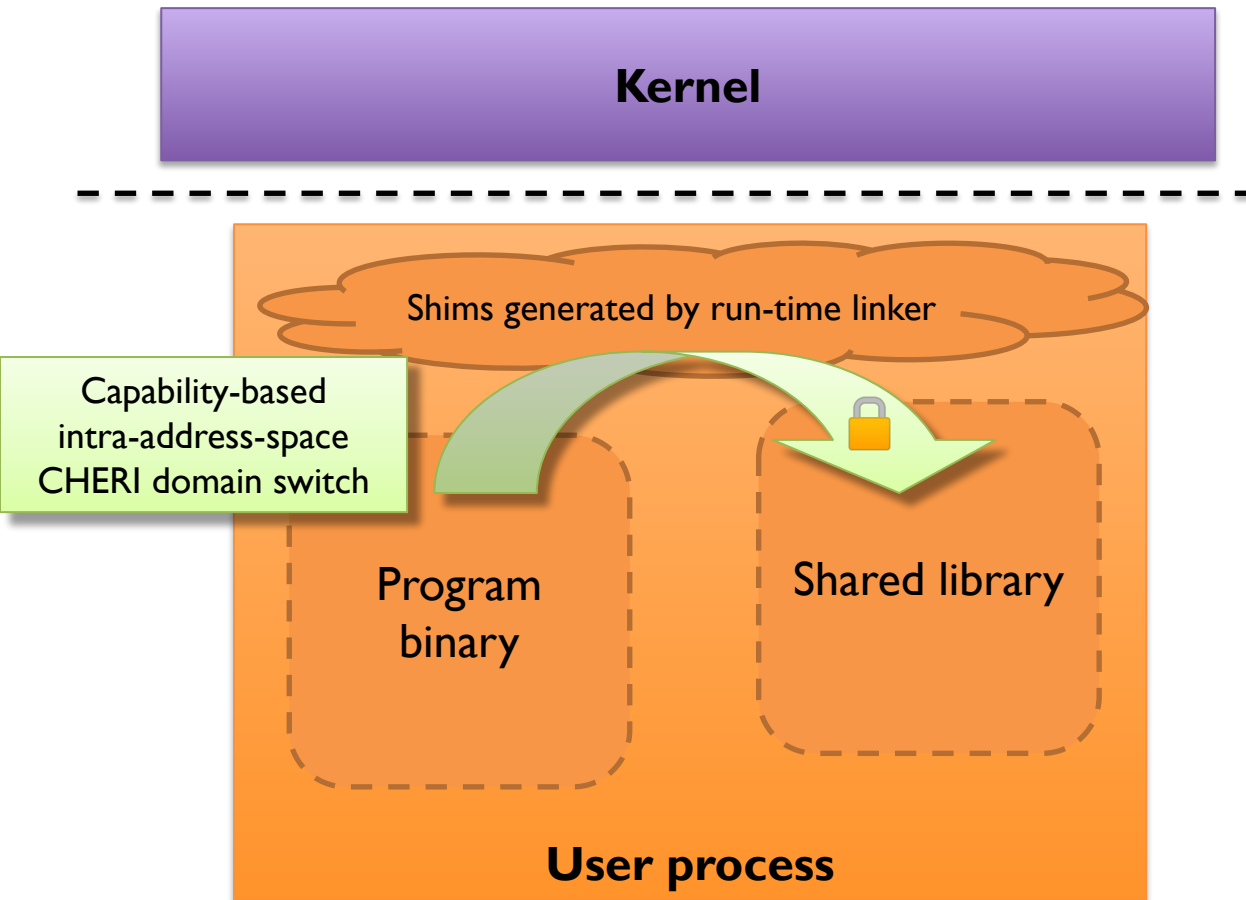- Simple model eschews asynchrony, independent debugging, etc.

Prototype appeared in CheriBSD 22.12; updates in 23.11

## UNIX co-processes

- Multiple processes share a single virtual address space, separated using independent CHERI capability graphs.
- CHERI capabilities enable efficient sharing, domain transition.
- Rich model associates UNIX process with each compartment.

Prototype to appear in future CheriBSD release

UNIVERSITY OF CAMBRIDGE

# Shared library compartmentalization



Kernel

Shims generated by run-time linker

Capability-based intra-address-space CHERI domain switch
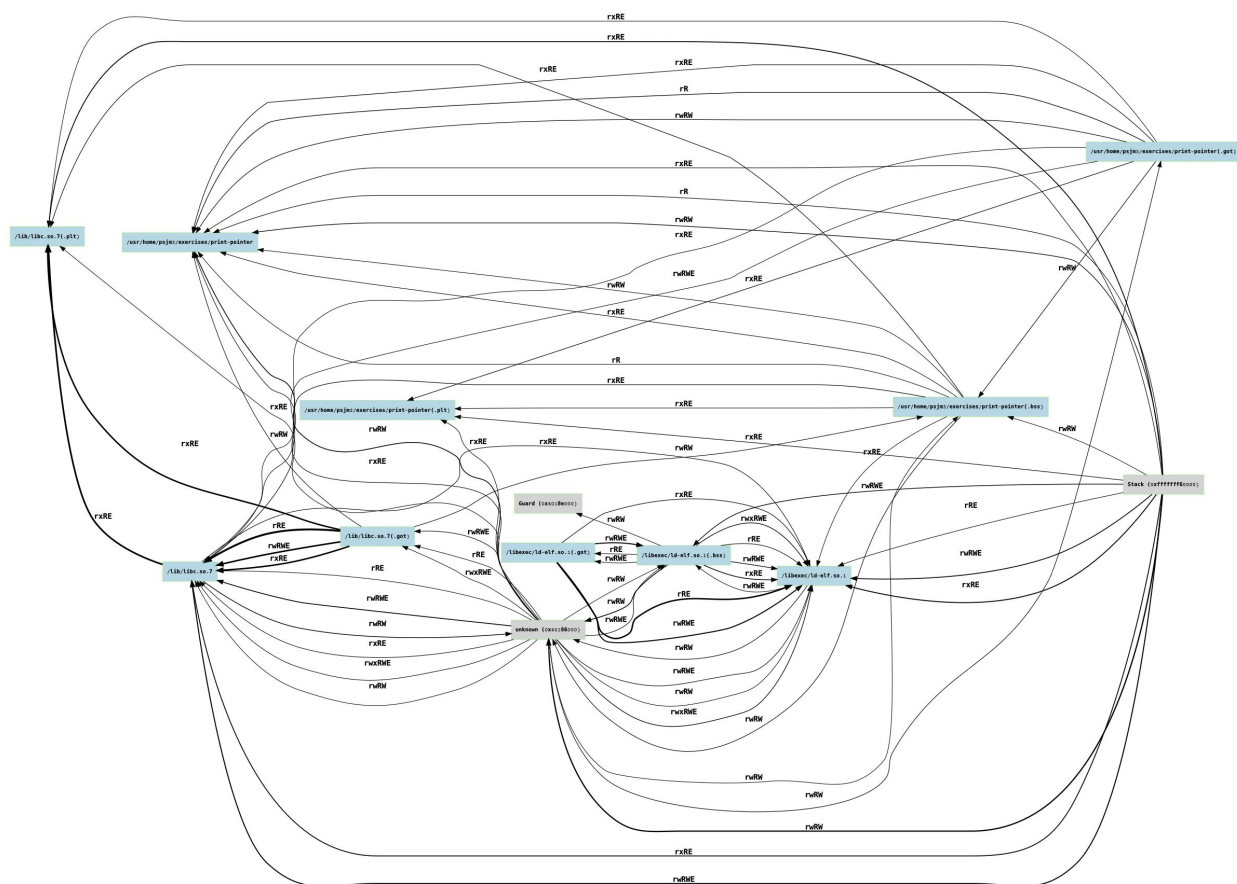
Program binary

Shared library

User process

- Run-time linker limits shared libraries to accesses enabled by ELF

  - Adversary model assumes arbitrary code execution within library

  - Run-time linker delegates capabilities for linked functions, globals via GOT/PLT

  - Domain transitions implemented by trampolines interposed on inter-object calls / returns

- Running prototype on Arm Morello

  - Released in CheriBSD 22.12 in December

  - Low measured overheads in early experiments (e.g., ~1% for image decompression sandboxing)

  - More evaluation results for gRPC and nginx at the Capabilities Limited stand

SRI International

UNIVERSITY OF CAMBRIDGE

# Library-based compartmentalization in CheriBSD 23.11

- (in this workshop) Enabling, tracing, and analysing library compartments

- Changes in 23.11:

    - Support the new ABI for passing memory arguments

    - Trampolines now clear unused argument registers based on function signature information embedded in the ELF. This requires a custom-built toolchain though

    - Allow defining compartments that consist of multiple libraries. Although this functionality is currently only usable via an undocumented interface

    - Support for the benchmark ABI (in testing now – not in today's snapshot)

# Capability graph visualization and analysis



- Pointers are now directly visible in hardware – in memory, ISA-level traces, and so on

    - We can directly analyze capability delegation with CHERI

- New extraction tools scan virtual addresses spaces and binaries to enable:

    - Visualization

    - Validation

    - Debugging and optimization

- **Allows direct analysis of attacker-visible resources and attack surfaces**

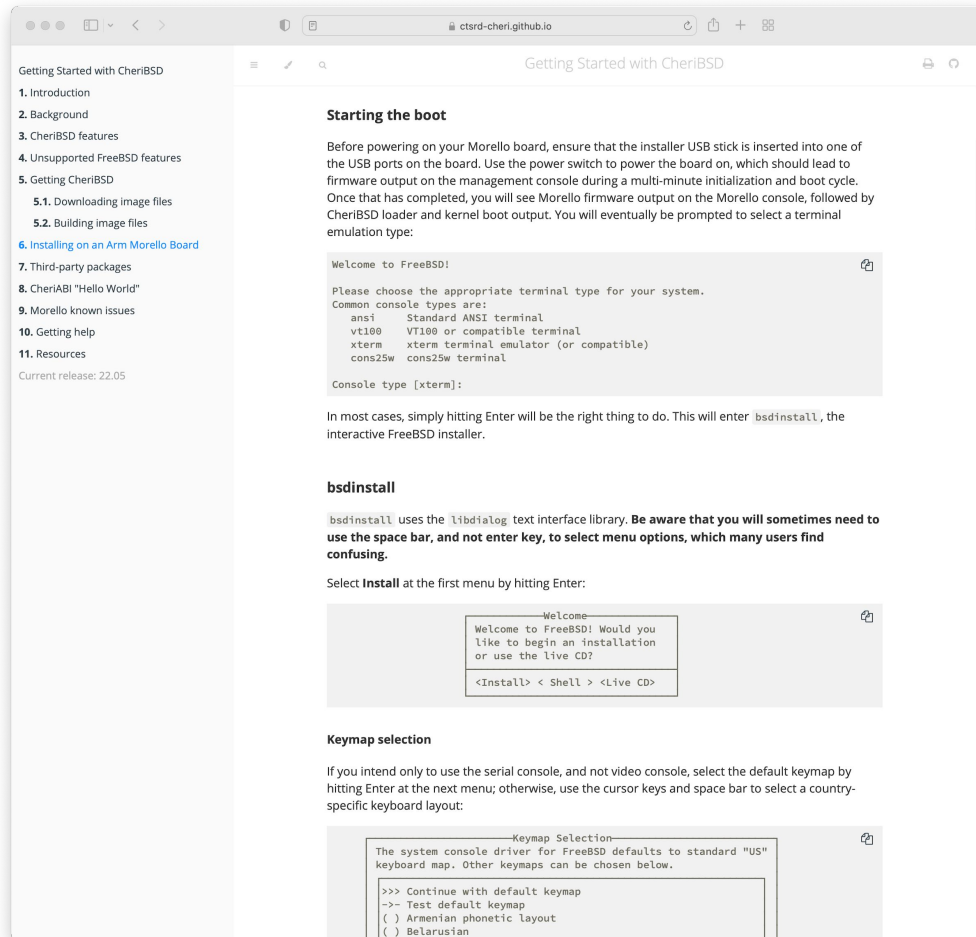# Capability introspection utilities

CheriBSD 23.11 will include the chericat utility that collects information about in-memory capabilities and stores them in an SQLite database for analysis and visualization. This tool will improve over time, growing new visualization, analysis, and testing abilities.

CheriTree from rtegrity allows to print a capability derivation tree to assist a developer in applying the library-based compartmentalization model to their application.

# Other improvements

- Upstream FreeBSD ports changes from April 2022 up to August, 2023 have been merged into CheriBSD ports

- The number of CheriABI packages increased from 9,104 to 10,307

- Added package manager pkg64cb for Benchmark ABI packages (9,741)

- The experimental co-process compartmentalization model now works with for CheriBSD/Morello – but not yet recommended for use

- bhyve hypervisor changes are being reviewed to be included in FreeBSD (arm64) and then merged into CheriBSD 23.11 (w/CHERI support) this week

# Getting Started with CheriBSD 23.11 (WIP)



https://www.cheribsd.org/getting-started/23.11/

- Introduces CheriBSD

- High-level documentation of CHERI-specific features, such as heap temporal safety

- Steps you through installation on a Morello board using a USB stick image that you can download

- Describes third-party package system and pkg64/pkg64c/pkg64cb

- Illustrates "hello world" compilation and debugging for different ABIs – now including Benchmark ABI

- Includes benchmarking guidelines

- Describes some known issues

- Explains how to get support

# How to obtain and install the CHERI software stack (1/2)



- One build tool to rule them all: cheribuild

  **https://github.com/CTSRD-CHERI/cheribuild**

- Builds, installs, and/or runs:

  - QEMU CHERI-RISC-V and Morello, Morello FVP

  - CheriBSD/CHERI-RISC-V and Morello disk images

  - Small suite of adapted third-party applications

- Up and running with one command:

  ./cheribuild.py --include-dependencies run-riscv64-purecap

  ./cheribuild.py --include-dependencies run-morello-purecap

# How to obtain and install the CHERI software stack (2/2)

You can also bootstrap a Docker container with CheriBSD for CHERI-RISC-V or Arm Morello using:

docker pull ctsrd/cheribsd-sdk-qemu-riscv64-purecap

docker pull ctsrd/cheribsd-sdk-qemu-morello-purecap

# Part II: practice
Exercises.

The CheriBSD 23.11 tutorial is available at:

**cheribsd.org/tutorial/23.11/**

Instructions and exercises can be found there now

The jails will remain active until November 15.

# Work environment

- You will work with CheriBSD/Morello using a FreeBSD jail (container)

- You only need an SSH client

- There are pre-installed packages:
  - CheriABI: git, nano, chericat, tmux
  - Benchmark ABI: pkg64cb
  - Hybrid ABI: llvm-base, gdb-cheri, vim

- You can switch to **root** with **sudo su -l**
  You can install packages and change the configuration of your jail

- You can install your own SSH key in ~/.ssh/authorized_keys

- We will distribute initial SSH credentials on paper

UNIVERSITY OF CAMBRIDGE

Talk to us on the CHERI-CPU Slack or mailing lists:

**https://cheri-cpu.org/**

(→ CHERI → CHERI-CPU Slack)

To find get started on the practical material:

**cheribsd.org/tutorial/23.11/**