### Imperial College London









# Intravisor: Type-3 Hypervisor for Capability-Based Virtual Machines

# Vasily A. Sartakov

Imperial College London

http://lsds.doc.ic.ac.uk <v.sartakov@imperial.ac.uk>

Joint work with Lluís Vilanova<sup>1</sup>, David Eyers<sup>2</sup>, Takahiro Shinagawa<sup>3</sup>, Peter Pietzuch<sup>1</sup> Imperial College London<sup>1</sup>, University of Otago<sup>2</sup>, The University of Tokyo<sup>3</sup> CHERITech – September 2022

### **Clouds: Isolation vs. Sharing**

#### Cloud services must be **isolated** from each other and the cloud stack

Services must **share** data efficiently by crossing isolation boundaries



# VMs: Strong, Heavyweight Isolation

- Strong isolation guarantees
- Small(ish) trusted computing base (TCB)
  - Only hypervisor
- Network communication for sharing → TCP/IP
   Requires data serialisation and copying
- Expensive transitions between services
  - Hypercalls  $\approx 50 \times$  syscalls



Hardware

# **Containers: Weak, Lightweight Isolation**

- Lightweight OS namespace isolation
- Efficient IPC mechanisms
- Large TCB due to shared OS kernel
  - Shared kernel has much unnecessary functionality



 $\rightarrow$  Challenge: efficient data sharing with small TCB

Memory Management Unit (MMU) is a privileged entity

– Intermediary (kernel) is always involved in IPC  $\rightarrow$  Shared TCB, sys/hyper-calls

MMU shares data at page granularity

- Sharing may expose extra data

Can we use another technology for isolation and sharing?

 $\rightarrow$  CHERI: isolation at byte granularity, low dependency on the kernel

## **CHERI** Capabilities

Fat pointers protected by hardware:

- base + length, cursor
- permission, tag
- byte-granularity\*

Fine-grained isolation Limited Dependency on OS kernel Available: Arm CHERI Morello Boards (Armv8)

Capabilities can be created only from capabilities

- Using cap-aware instructions, but not the intermediary



#### Challenges of Cloud Stack with Hardware Capabilities

What would a cloud stack look like if hardware provided <u>efficient</u> mechanisms to share <u>arbitrary-sized</u> memory regions between otherwise <u>isolated</u> entities?

Challenges:

- C1. Support capability-unaware software
- C2. Provide small-TCB OS functionality
- C3. Enable efficient capability-based IPC interfaces

## **cVM: Intra-Process VM-like Abstraction**

- 1. Support cap-unaware software
- $\rightarrow$  Isolated execution of native applications
- 2. Small shared TCB
- $\rightarrow$  Private namespaces by library OSs
- 3. Cap-based IPC interfaces
  - CP\_File: efficient data sharing
  - CP\_Call: remote code invocation



Isolation

cVM

Native ABI: cap-unaware code Pure-capability ABI: requires porting Hybrid-capability ABI: native + cap-aware code

Fine-grained compartmentalisation:

- Cap-unaware instructions are constrained by *default* caps
- Hybrid code can use capability-aware instructions

 $\rightarrow$  Can be used for isolation and IPC primitives



## **Support for Native Software**

Goals:

- POSIX environment
- Cloud deployment model (e.g. Docker or VMs)
- $\rightarrow$  Service for cVM is shipped as disk image
  - Native cap-unaware PIE binaries
  - Compatibility: C standard library (musl libc)
- $\rightarrow$  Intravisor allocates cVM, loads Init and disk

	Program/lib	Dependencies (.so libs)
	LibC(musl)	
		Init
Ē		
	$\bigcirc$	Intravisor

Kernel

# **Small-TCB OS Functionality**

Goals:

- Necessary OS components
- Small attack surface

- → Private LibraryOSs provide OS functionality
- $\rightarrow$  Intravisor provides time/net/disk
- $\rightarrow$  Nested Isolation layers



# **IPC-Like Interfaces Using Capabilities**

#### Data sharing primitives efficient if:

- Non-shared and without intermediary on critical path
- Well-known API (POSIX)
- Usable by cap-unaware code



- CP\_File read/write remote memory at byte granularity using caps
- CP\_Call call a function in a cVM

## **IPC-Like Interfaces Using Capabilities**

- CP File asynchronous data access
- CP\_Call notification mechanism

 $\rightarrow$  Base for complex IPC interfaces

CP Stream - stream-oriented IPC interface

- Destination buffer unknown to the sender
- Recipient pre-registers input buffers and uses <code>cp\_poll</code> to get a notification

# Attacker free to run any code, even cap-aware CHERI instructions

They can get access to entire cVM including compromised libOS

They may try to abuse CP\_Files, CP\_Calls, hostcall, and the kernel



# **Security Analysis**

Kernel ignores syscalls

 $\rightarrow$  cannot abuse the kernel

CP\_Files are data caps

 $\rightarrow$  cannot be used for CInvoke

cVMs cannot store or export caps

 $\rightarrow$  cannot access revoked data

Hostcall caps are sealed

 $\rightarrow$  cannot be changed



 $\rightarrow$  Attacker cannot escape the cVM or gain unauthorised access to data

# Prototype

#### Platforms:

- CHERI RISC-V64, QEMU, AWS F1 (agfi-026d853003d6c433a)
- CheriBSD (host), LKL v4.17 with musl v1.2.1 (cVMs)
- SiFive HiFive Unmatched (No CHERI, but multi-core)

Application and services (In the paper):

- Redis, Data-processing utilities, Python3 with modules, SQLite, benchmarks
- Multi-tier microservice (NGINX with API gate, Redis SiFive only)

#### Evaluation question: Performance of cVM IPC primitives?

- Basic: memcpy, mmap+memcpy
- cVMs: CP\_File, CP\_Stream
- FreeBSD: pipe, Unix, and TCP sockets

## **Comparing with IPC Mechanisms**



CP\_FILE vs. memcpy: - 6% slower

```
CP_Stream faster (1.2 MB+)
- Privileged execution
```

Unix, TCP, mmap+memcpy: – Less than 2.4-3.6 MB/s Processes: 1.6 MB/s max Small-TCB isolation with efficient data sharing in clouds is challenging:

- Containers  $\rightarrow$  large shared TCB with relatively fast IPC mechanisms
- VMs  $\rightarrow$  small TCB with slow IPC mechanisms

CAP-VMs provide VM-like abstraction:

- Secure isolation using memory capabilities
- Controlled shared TCB by private libraryOSs
- Efficient data sharing using capability-based IPC primitives

Available at http://github.com/lsds/intravisor:

- Runtime: musl-LKL, baremetal
- Musl-LKL: cp\_file, cp\_stream, docker-based images (Redis, NGINX), helloworld, python, sqlite
- Baremetal: two pure cVMs, two nested cVMs, hello world
- Arch: RISC-V

Future plans:

- Pure cVMs mostly
- Arm Morello
- Next major update: Jan'23

Thank You — Any Questions? Vasily A. Sartakov v.sartakov@imperial.ac.uk

