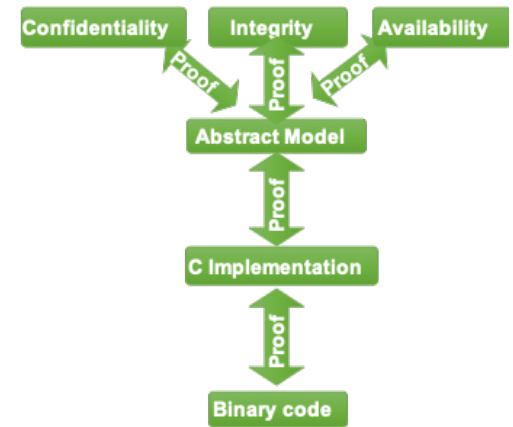School of Computer Science & Engineering

**COMP9242 Advanced Operating Systems**

2019 T2 Week 09b
**Local OS Research**
@GernotHeiser

# Copyright Notice

**These slides are distributed under the
Creative Commons Attribution 3.0 License**

- You are free:
    - to share—to copy, distribute and transmit the work
    - to remix—to adapt the work

- under the following conditions:
    - **Attribution:** You must attribute the work (but not in any way that suggests that the author endorses you or your use of the work) as follows:

        *"Courtesy of Gernot Heiser, UNSW Sydney"*

The complete license text can be found at
http://creativecommons.org/licenses/by/3.0/legalcode

UNSW
SYDNEY

# Quantifying Security Impact of Operating-System Design

# Quantifying OS-Design Security Impact

**Approach:**

- Examine all *critical* Linux CVEs (vulnerabilities & exploits database)

- easy to exploit
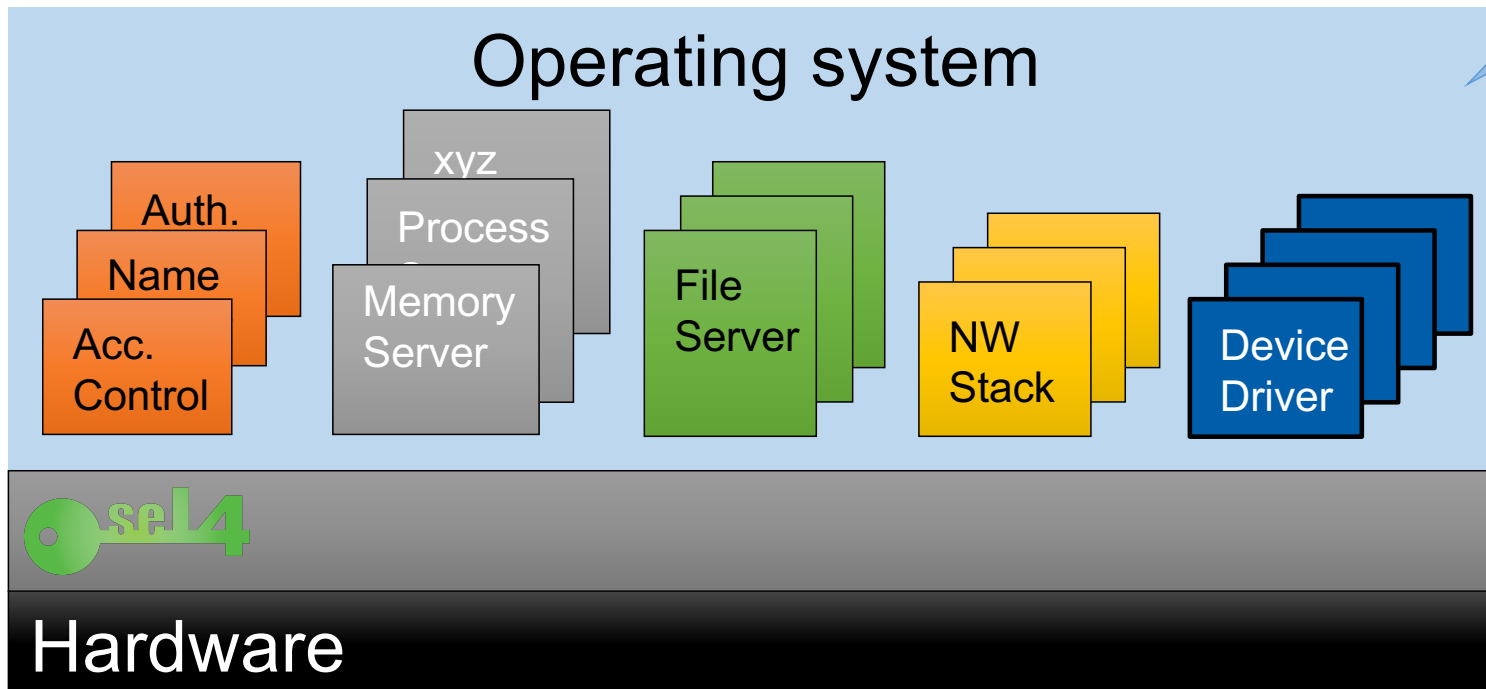- high impact
- no defence available
- confirmed

115 critical
Linux CVEs
to Nov'17

- For each establish how microkernel-based design would change impact
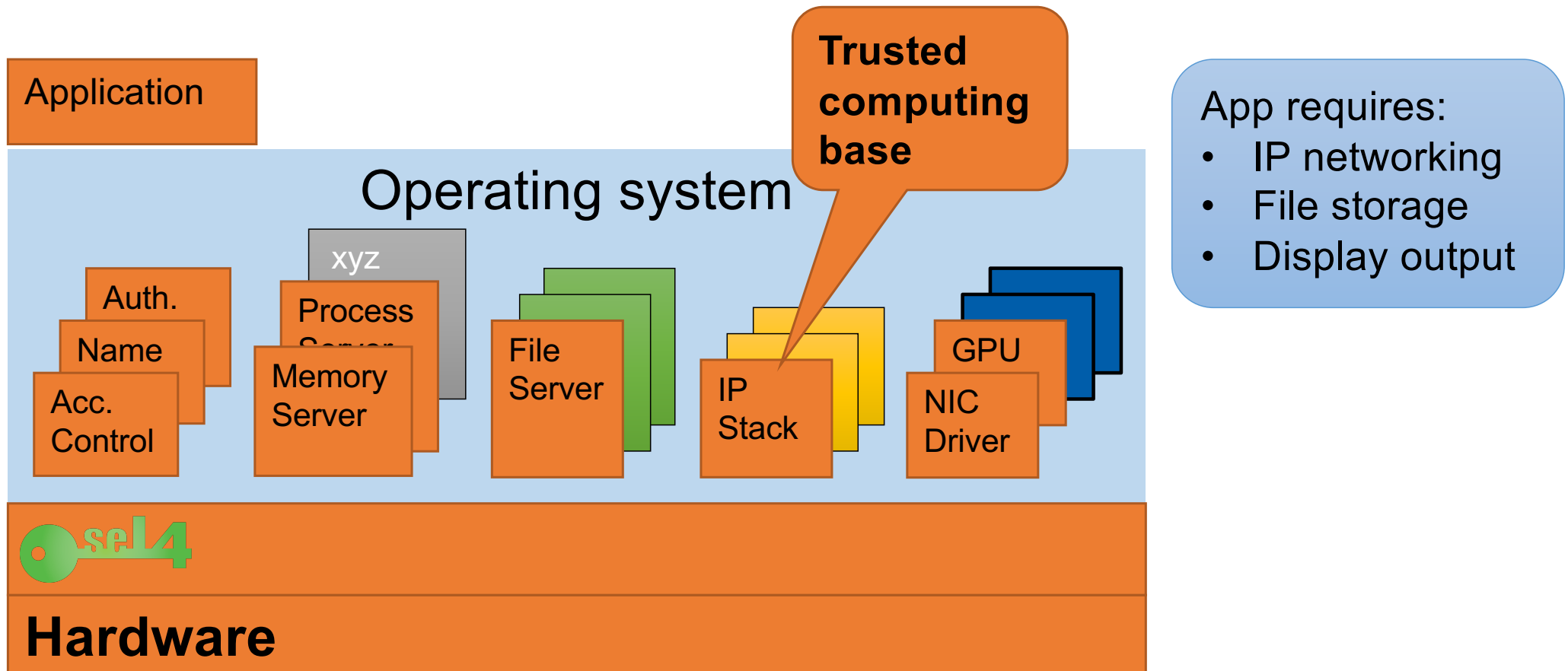
UNSW
SYDNEY

# Hypothetical seL4-based OS

OS structured in *isolated* components, minimal inter-component dependencies, *least privilege*
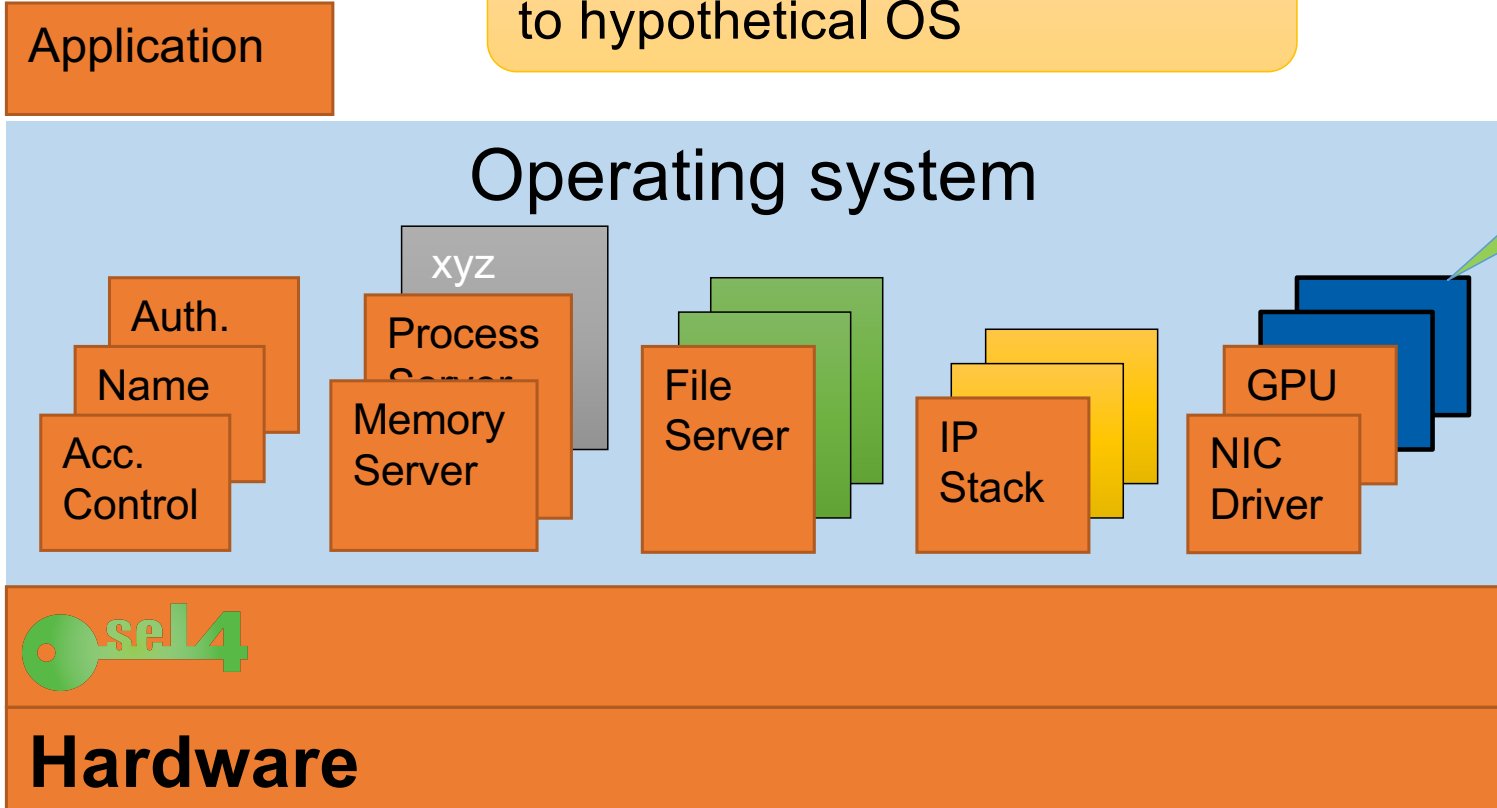
Functionality comparable to Linux

### Operating system

Auth.

Name

Acc. Control

xyz Process

Memory Server

File Server

NW Stack

Device Driver

### seL4

### Hardware

UNSW SYDNEY

# Hypothetical Security-Critical App

**Application**

Operating system

**Trusted computing base**

Auth.

Name

Acc. Control

xyz

Process Server

Memory Server

File Server

IP Stack

GPU

NIC Driver

App requires:
- IP networking
- File storage
- Display output

**sel4**

**Hardware**

# Analysing CVEs

Map compromised component to hypothetical OS

Application

## Operating system

Not in TCB:
**Attack defeated**

Auth.

Name

Acc. Control

xyz

Process Server

Memory Server

File Server

IP Stack

GPU

NIC Driver

**Example:** USB driver bug

**Hardware**

UNSW
SYDNEY

# Analysing CVEs

Map compromised component to hypothetical OS

**Example:** Bug in page-table management

Application

## Operating system

xyz

Auth.

Name

Acc. Control

Process Server

Memory Server

File Server

IP Stack

GPU

NIC Driver

In microkernel: **Attack defeated by verifiation**

sel4

**Hardware**

UNSW SYDNEY

# Analysing CVEs

Map compromised component to hypothetical OS

Only *crash* essential service (DoS): **Strongly mitigated**

Application

## Operating system

xyz

Auth.

Name

Acc. Control

Process Server

Memory Server

File Server

IP Stack

GPU

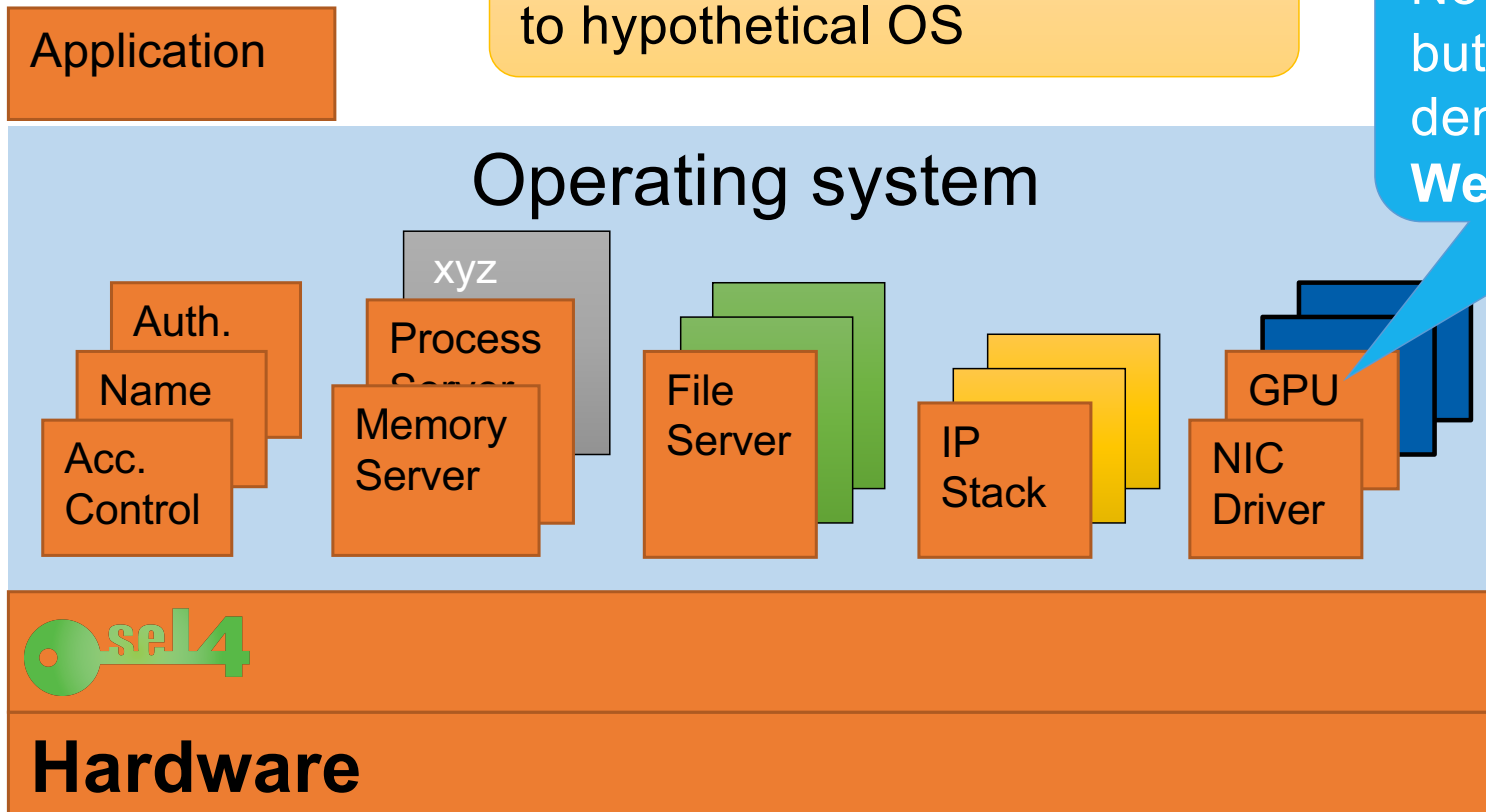NIC Driver

**Example:** File system compromised

**Hardware**

# Analysing CVEs

Map compromised component to hypothetical OS

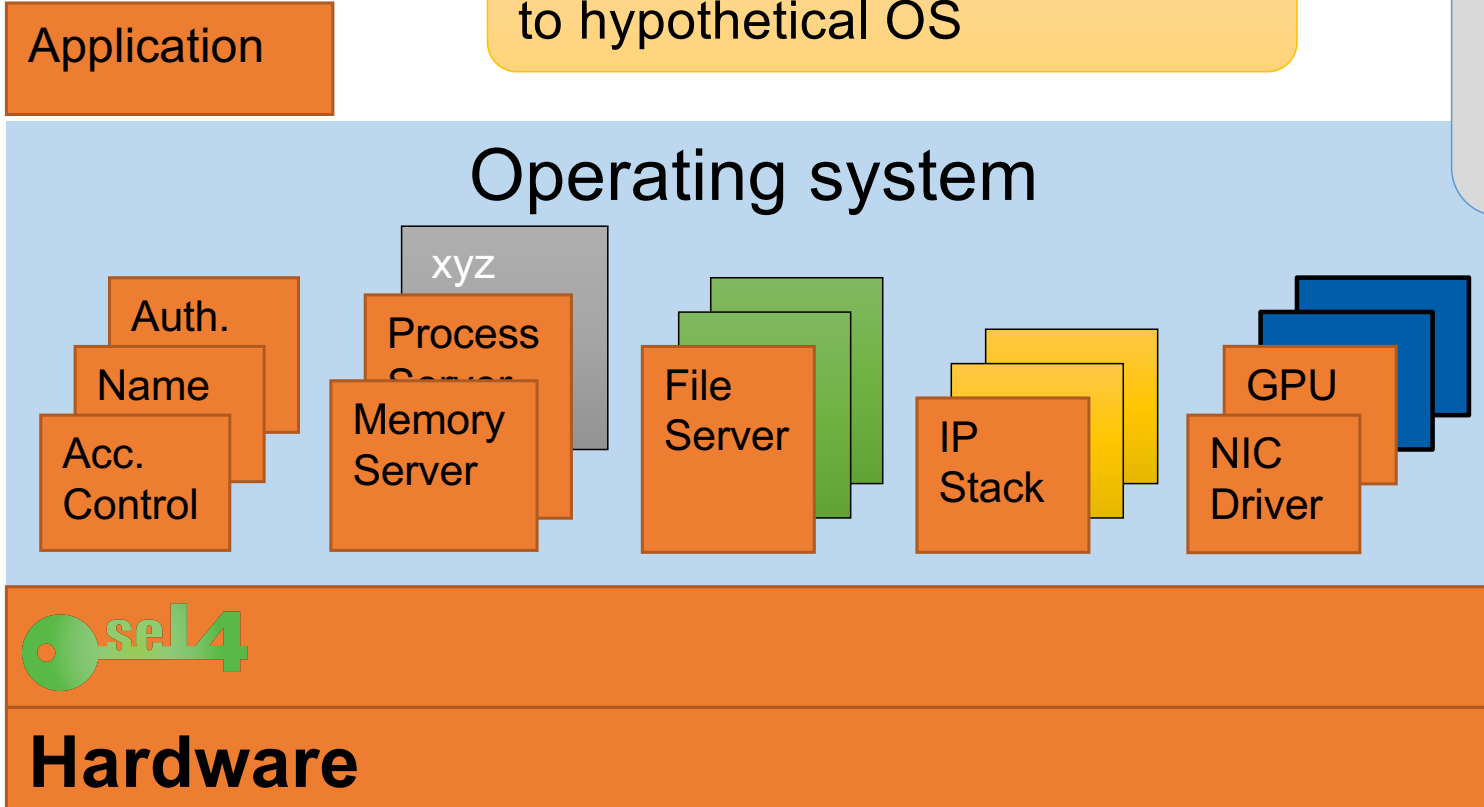No full compromise but integrity or confidentiality violation: **Weakly mitigated**

Application

## Operating system

xyz

Auth.

Name

Acc. Control

Process Server

Memory Server

File Server

IP Stack

GPU

NIC Driver

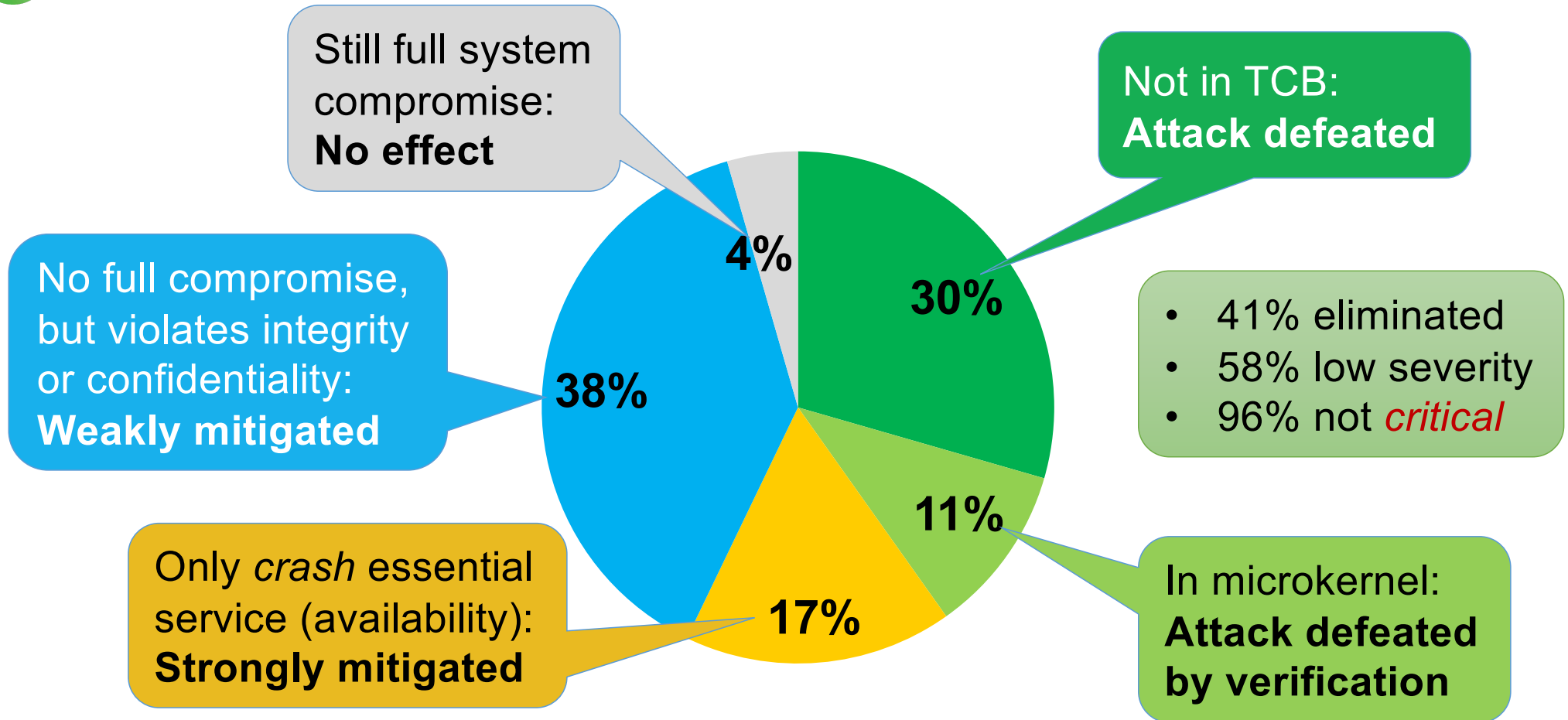**Example:** GPU compromised

**Hardware**

UNSW SYDNEY

# Analysing CVEs

Map compromised component to hypothetical OS

**Example:**
Driver exploit hijacks I2C bus, allowing firmware reflush

Application

## Operating system

xyz

Auth.

Name

Acc. Control

Process Server

Memory Server

File Server

IP Stack

GPU

NIC Driver

Full system compromise: **No effect**

**Hardware**

UNSW SYDNEY

# All Critical Linux CVEs to 2017

Still full system compromise: **No effect**

Not in TCB: **Attack defeated**

No full compromise, but violates integrity or confidentiality: **Weakly mitigated**

- 41% eliminated
- 58% low severity
- 96% not *critical*

Only *crash* essential service (availability): **Strongly mitigated**

In microkernel: **Attack defeated by verification**

4%

30%

38%

17%

11%

UNSW SYDNEY

# Summary

**OS structure matters!**

- Microkernels definitely improve security

- Monolithic OS design is *fundamentally flawed from security point of view*

[Biggs et al., APSys'18]

**Use of a monolithic OS in security- or safety- critical scenarios is professional malpractice!**
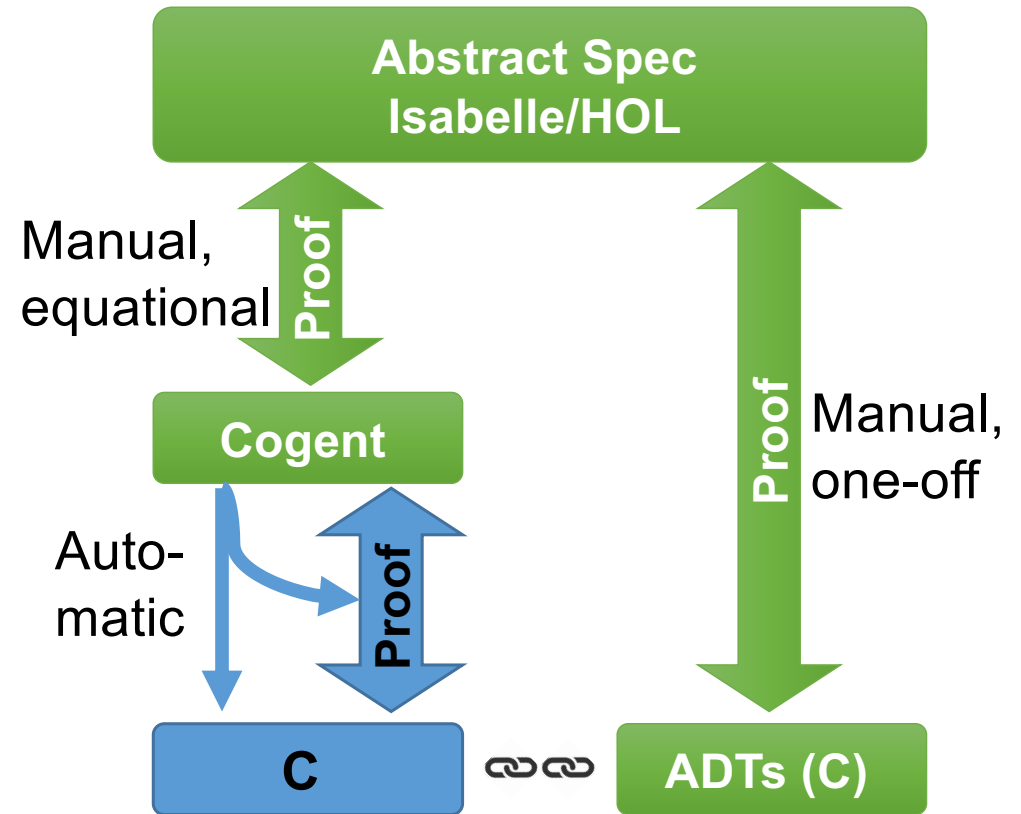
# Cogent

UNSW
SYDNEY

# Beyond the Kernel

5 kLOC?

1 kLOC?

100 kLOC?

10 kLOC?

Control

Device driver

NW stack

File system

Uncritical/ untrusted

Apps

Linux

10 kLOC 11 py

seL4

Aim: Verified TCB at affordable cost!

UNSW
SYDNEY

# Cogent: Code & Proof Co-Generation

Aim: Reduce cost of verified systems code

- Restricted, purely functional *systems* language
- Type- and memory safe, not managed
- Turing incomplete
- File system case-studies: BilbyFs, ext2, F2FS, VFAT

[O'Connor et al, ICFP'16; Amani et al, ASPLOS'16]

**Abstract Spec Isabelle/HOL**

Manual, equational — **Proof**

**Cogent**

Auto-matic — **Proof**

Manual, one-off — **Proof**

**C** 🔗🔗 **ADTs (C)**

UNSW
SYDNEY

# Manual Proof Effort

| BilbyFS functions | Effort | Isabelle LoP | Cogent SLoC | Cost $/SLoC | LoP/ SLOC |
|---|---|---|---|---|---|
| isync()/ iget() library | 9.25 pm | 13,000 | 1,350 | 150 | 10 |
| sync()-specific | 3.75 pm | 5,700 | 300 | 260 | 19 |
| iget()-specific | 1 pm | 1,800 | 200 | 100 | 9 |
| seL4 | 12 py | 180,000 | 8,700 C | 350 | 20 |

BilbyFS: 4,200 LoC Cogent

UNSW SYDNEY

# Addressing Verification Cost



**Dependability-cost tradeoff:**

- Reduced faults through safe language
- Property-based testing (QuickCheck)
- Model checking
- Full functional correctness proof

**Spec reuse!**

**Work in progress:**

- Language expressiveness
- Reduce boiler-plate code
- Network stacks
- Device drivers

UNSW
SYDNEY

# Time Protection

# Refresh: Microarchitectural Timing Channels
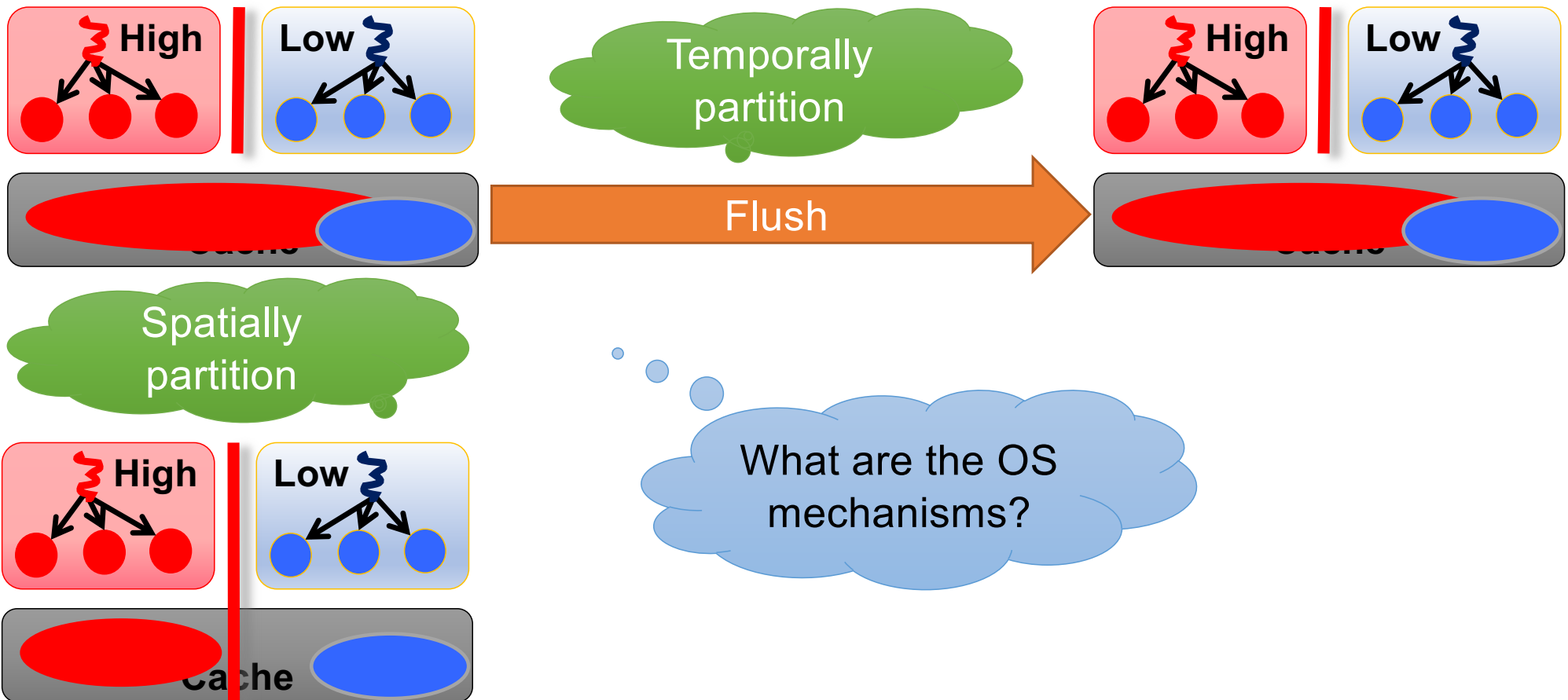


Contention for shared hardware resources affects execution speed, leading to timing channels

# OS Must Enforce *Time Protection*



**Preventing interference is core duty of the OS!**

- *Memory protection* is well established
- *Time protection* is completely absent

UNSW
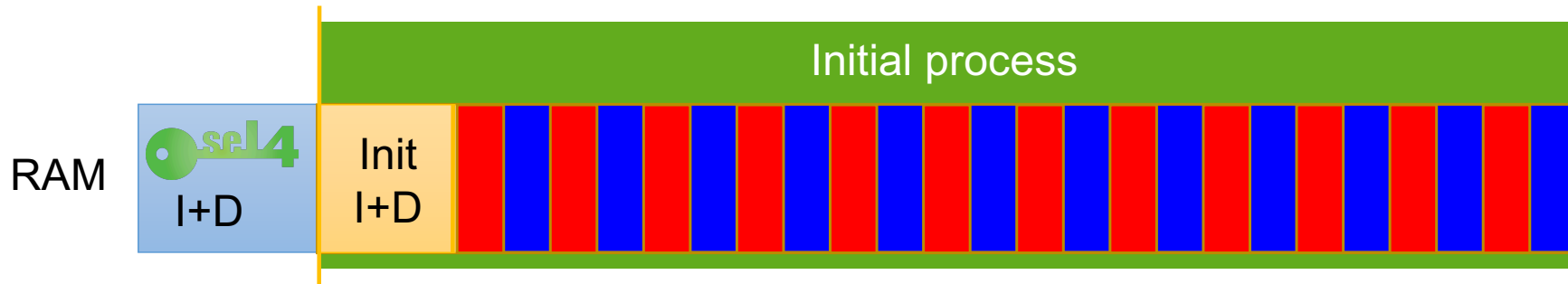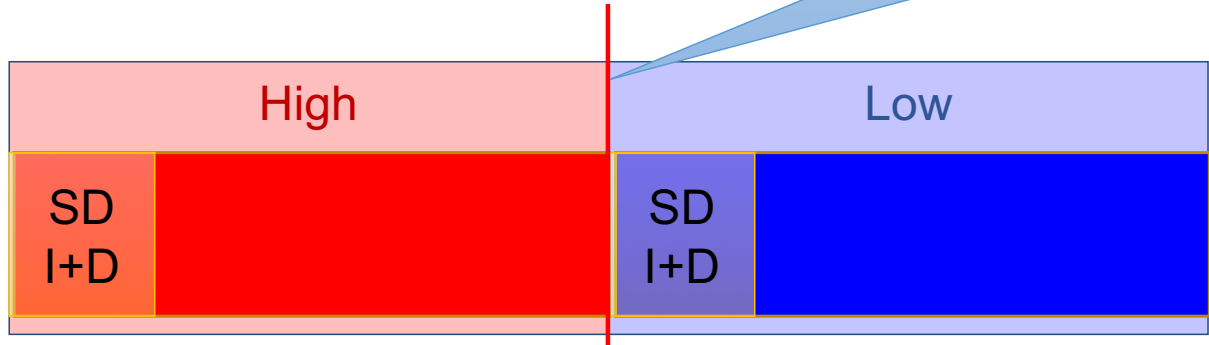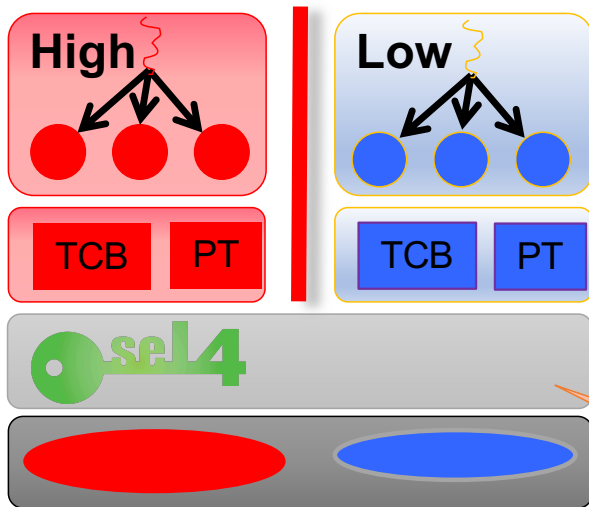SYDNEY

# Time Protection: No Sharing of HW State

# Spatial Partitioning: Cache Colouring

**System permanently coloured**

Partitions restricted to coloured memory

| High | | Low | |
|---|---|---|---|
| SD I+D | | SD I+D | |

Initial process

RAM | sel4 I+D | Init I+D |

COMP9242 2019T2 W09b: Local OS Research

UNSW SYDNEY

# Spatial Partitioning: Cache Colouring

**High**

**Low**

TCB | PT

TCB | PT

- Partitions get frame pools of disjoint colours
- seL4: userland supplies kernel memory
  ⇒ colouring userland colours kernel memory
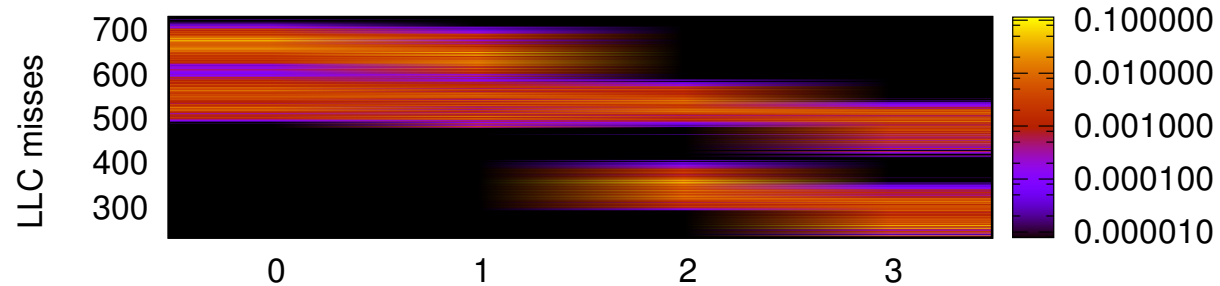
Shared kernel image

# Channel Through Kernel Code

Raw channel



Channel matrix: Conditional probability of observing output signal (time) given input signal (system-call number)
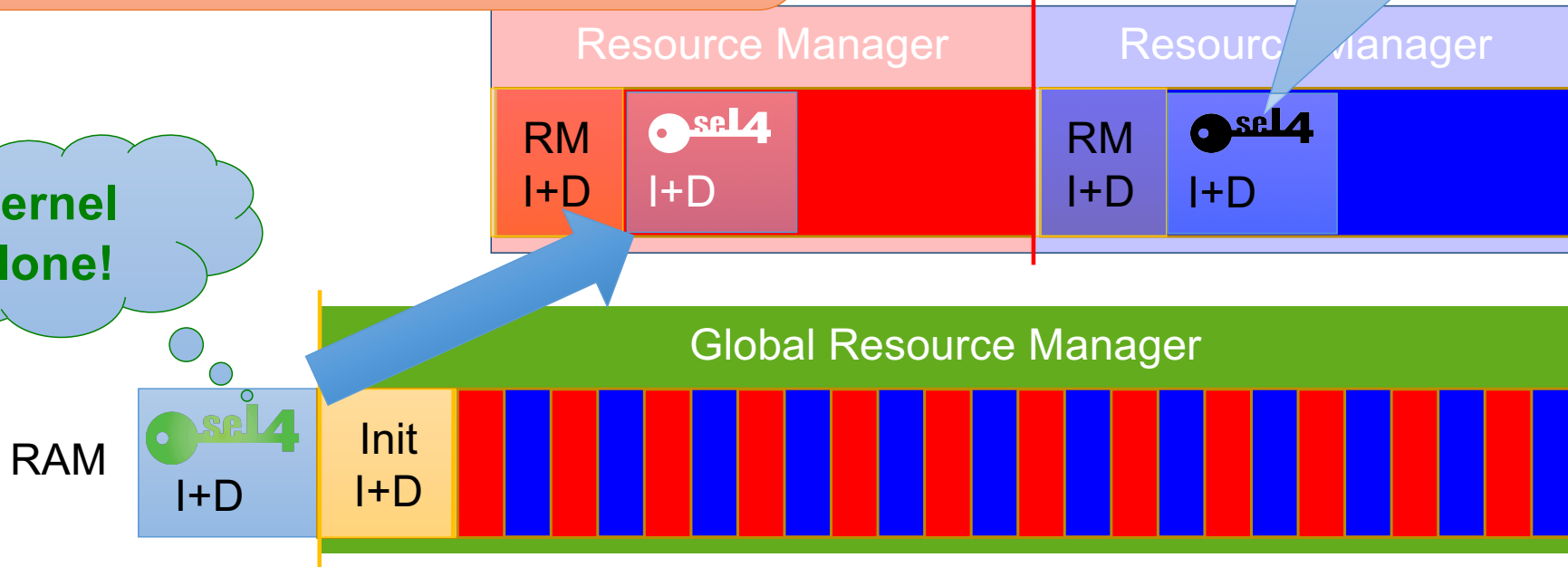
# Colouring the Kernel

Ensure deterministic access!

Each partition has own kernel image

Remaining shared kernel data:

- Scheduler queue array & bitmap
- Few pointers to current thread state

Resource Manager

Resource Manager

RM I+D

seL4 I+D

RM I+D

seL4 I+D

Kernel clone!

Global Resource Manager

RAM

seL4 I+D

Init I+D

UNSW SYDNEY

# Spatial Partitioning: Cache Colouring



- Partitions get frame pools of disjoint colours
- seL4: userland supplies kernel memory
  ⇒ colouring userland colours kernel memory
- Per-partition kernel image to colour kernel
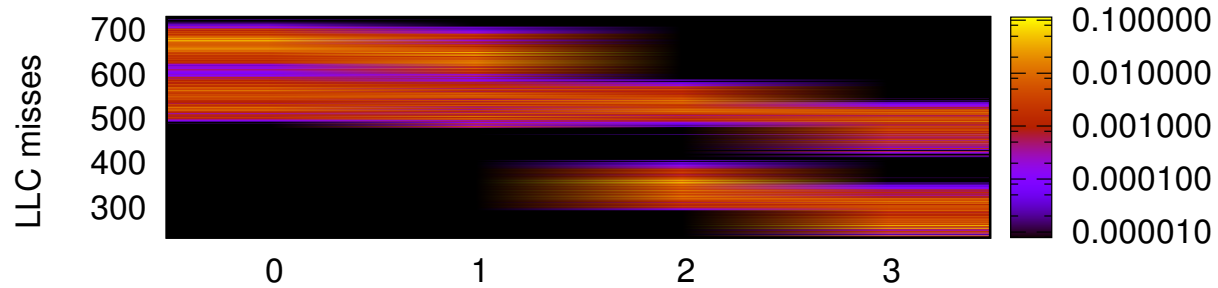
Ensure deterministic access!

Remaining shared kernel data:
- Scheduler queue array & bitmap
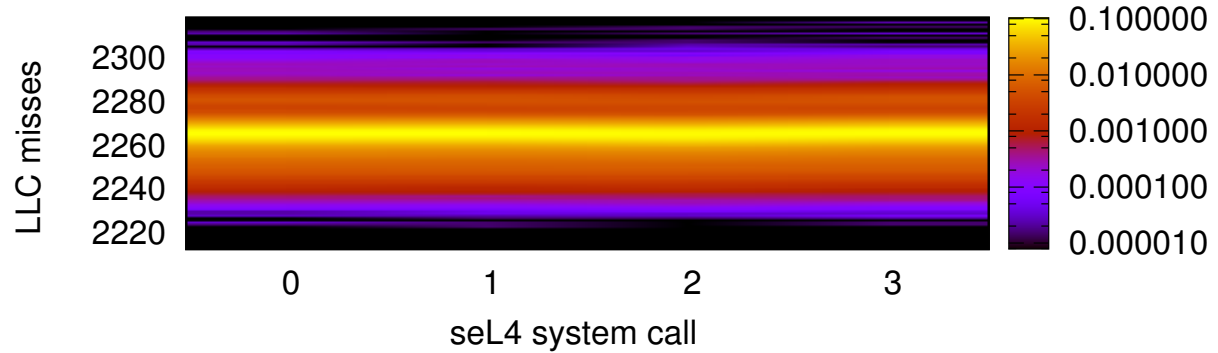- Few pointers to current thread state

UNSW
SYDNEY

# Channel Through Kernel Code



Raw channel

Channel with cloned kernel

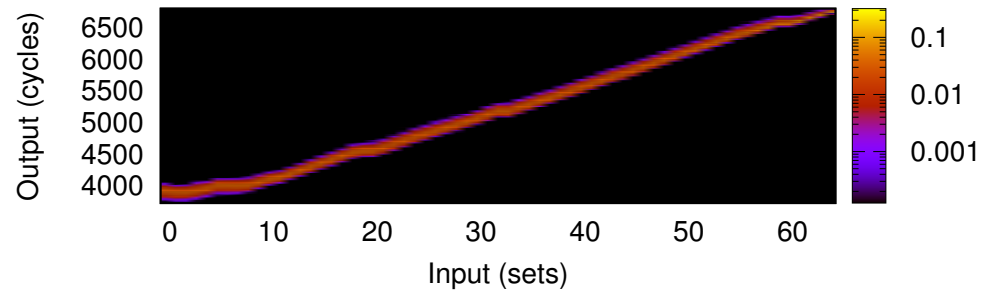COMP9242 2019T2 W09b: Local OS Research

# Temporal Partitioning: Flush on Switch
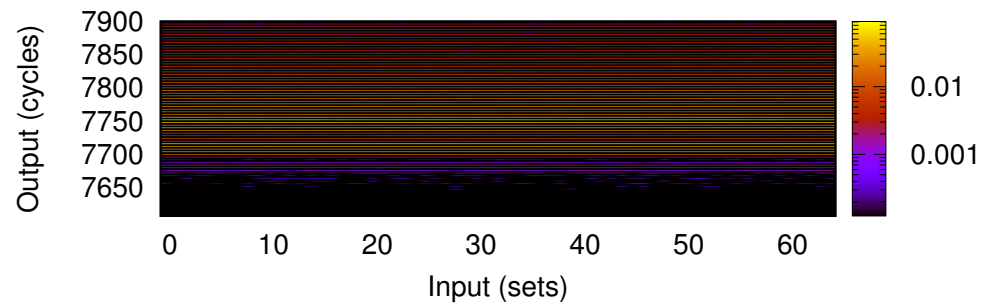
Must remove any history dependence!

2. Switch user context
3. Flush on-core state


6. Reprogram timer
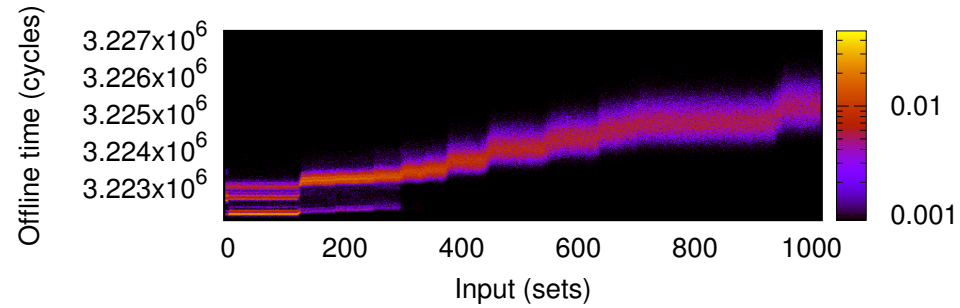7. return

UNSW
SYDNEY

# D-Cache Channel

Raw channel

Channel with flushing

# Flush-Time Channel

Raw channel



COMP9242 2019T2 W09b: Local OS Research

# Temporal Partitioning: Flush on Switch

Must remove any history dependence!

1. $T_0$ = current_time()
2. Switch user context
3. Flush on-core state
4. Touch all shared data needed for return
5. while ($T_0$+WCET < current_time()) ;
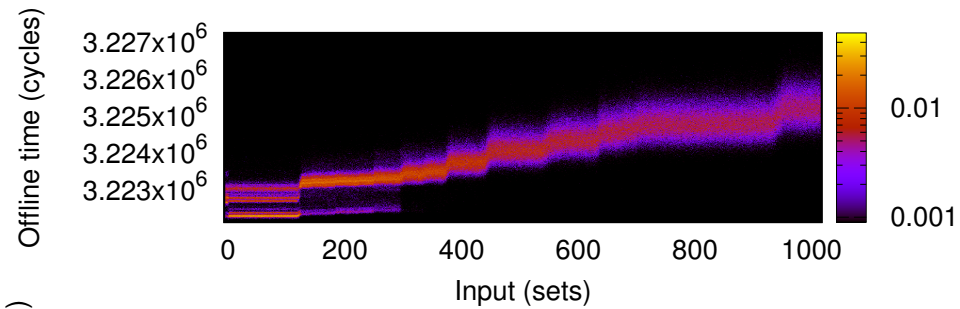6. Reprogram timer
7. return

Latency depends on prior execution!

Time padding to remove dependency

Ensure deterministic execution

UNSW SYDNEY

# Flush-Time Channel

**Raw channel**



**Channel with deterministic flushing**

# Performance Impact of Colouring

## Splash-2 benchmarks on Arm A9



Legend:
- 50% colours base
- 50% colour clone

- Overhead mostly low
- Not evaluated is cost of not using super pages [Ge et al., EuroSys'19]

| Architecture | x86 | Arm |
|---|---|---|
| Mean slowdown | 3.4% | 1.1% |

| Arch | seL4 clone | Linux fork+exec |
|---|---|---|
| x86 | 79 µs | 257 µs |
| Arm | 608 µs | 4,300 µs |

UNSW SYDNEY

# A New HW/SW Contract

For all shared microarchitectural resources:

aISA: augmented ISA

1.  Resource must be spatially partitionable or flushable

2.  Concurrently shared resources must be spatially partitioned

3.  Resource accessed solely by virtual address must be flushed and not concurrently accessed

    Cannot share HW threads across security domains!

4.  Mechanisms must be sufficiently specified for OS to partition or reset

5.  Mechanisms must be constant time, or of specified, bounded latency

6.  Desirable: OS should know if resettable state is derived from data, instructions, data addresses or instruction addresses

[Ge et al., APSys'18]

UNSW SYDNEY

# Can Time Protection Be Verified?

1. Correct treatment of spatially partitioned state:
   - Need hardware model that identifies all such state (augmented ISA)
   - To prove:
     **No two domains can access the same physical state**

   > Functional property!

   > Transforms timing channels into storage channels!

2. Correct flushing of time-shared state
   - Not trivial: eg proving all cleanup code/data are forced into cache after flush
     - Needs an actual cache model
   - Even trickier: need to prove padding is correct
     - … without explicitly reasoning about time!

   > Functional property!

UNSW
SYDNEY

# Verifying Time Padding

- Idea: Minimal formalisation of hardware clocks (abstract time)
  - Monotonically-increasing counter
  - Can add constants to time values
  - Can compare time values

**To prove: padding loop terminates as soon as timer value ≥ $T_0$+WCET**

Functional property

[Heiser et al., HotOS'19]

     UNSW SYDNEY

# Making COTS Hardware Dependable

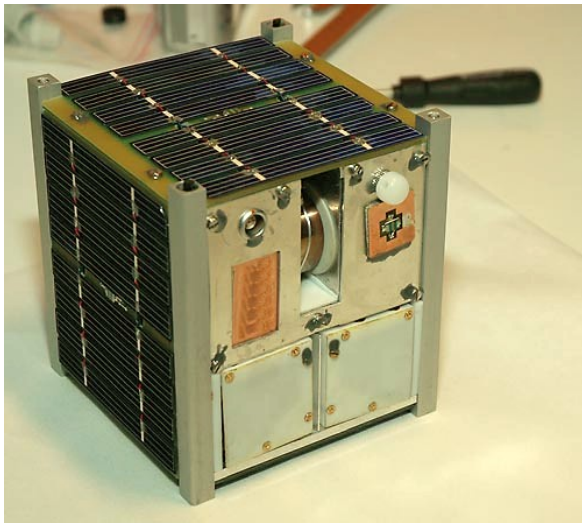COMP9242 2019T2 W09b: Local OS Research © Gernot Heiser 2019 – CC Attribution License

# Satellites: SWaP vs Dependability

Space is becoming commodisized:

- many, small (micro-) satellites
- increasing cost pressure

Harsh evironment for electronics:
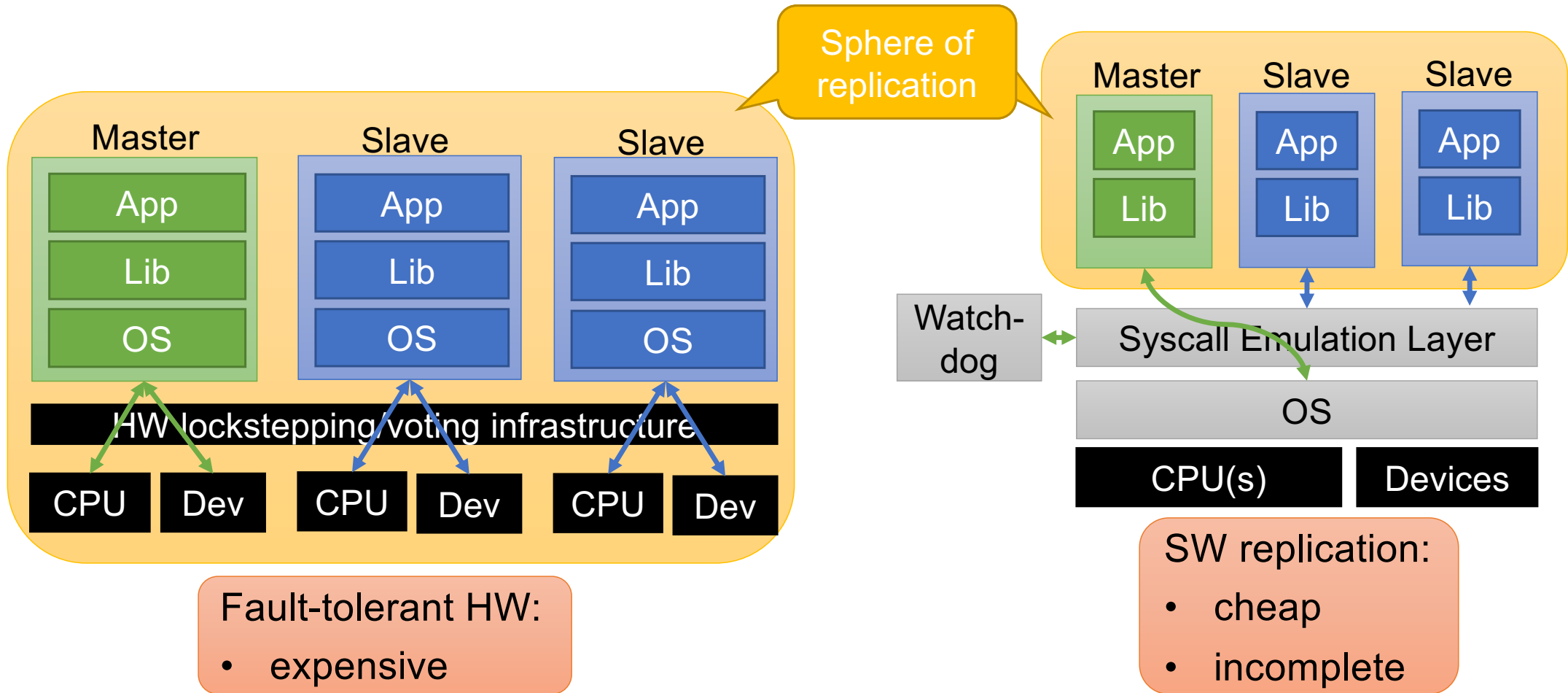
- temperature fluctuations
- ionising radiation



**NCUBE2** by Bjørn Pedersen, NTNU (CC BY 1.0)

Radiation-hardened processors are slow, bulky and expensive

Use redundancy of cheap COTS multicores

UNSW
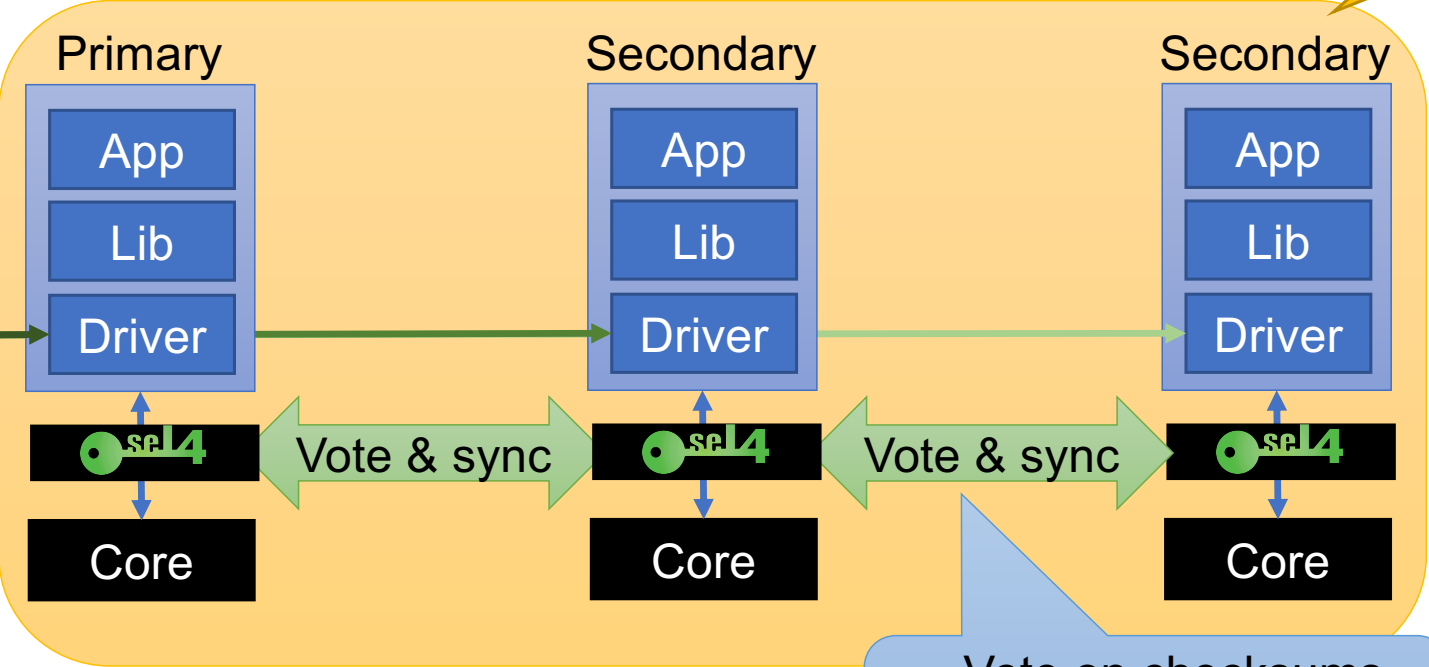SYDNEY

# Traditional Redundancy Approaches

Sphere of replication

Master Slave Slave

| Master | Slave | Slave |
|--------|-------|-------|
| App | App | App |
| Lib | Lib | Lib |
| OS | OS | OS |

HW lockstepping/voting infrastructure

| CPU | Dev | CPU | Dev | CPU | Dev |

Fault-tolerant HW:
• expensive

| Master | Slave | Slave |
|--------|-------|-------|
| App | App | App |
| Lib | Lib | Lib |

Watch-dog

Syscall Emulation Layer

OS

| CPU(s) | Devices |

SW replication:
• cheap
• incomplete

UNSW
SYDNEY

# Redundant Co-Execution (RCoE)

Sphere of replication

Userland transparently replicated

Device access:
- thin shim
- vote outputs
- copy inputs

**Primary**

App

Lib

Driver

**Secondary**

App

Lib

Driver

**Secondary**

App

Lib

Driver

Device interface

Device

Vote & sync

Vote & sync

Core

Core

Core

No master-slave, but peer-to-peer

- Vote on checksums of arguments & state
- Logical time for sync

UNSW SYDNEY

# RCoE: Two Variants

**Loosely-coupled RCoE**

- Sync on syscalls & exceptions

- Preemptions in usermode not further synchronised (imprecise)

- Low overhead
- Cannot support racy apps, threads, virtual machines

**Closely-coupled RCoE**

- Sync on instruction

- Precise preemptions

- High overhead
- Supports all apps
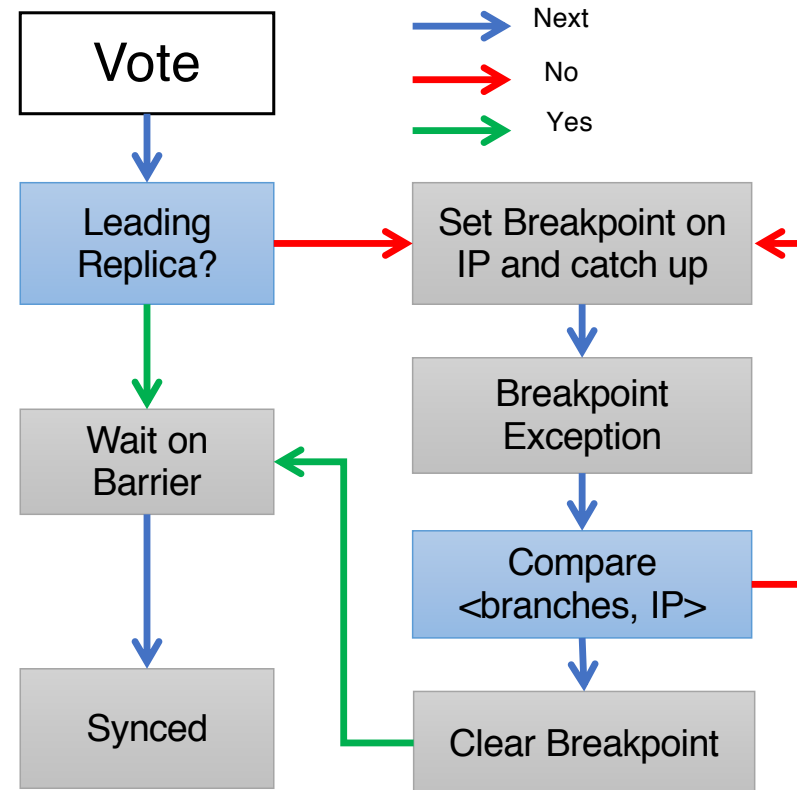- May need re-compile

UNSW
SYDNEY

# Closely-Coupled RCoE Implementation

Precise logical time: Triple of:

- event count
- user-mode branch count
- instruction pointer

x86: Obtained from PMU

Arm v7: Use gcc plugin to count branches

Next

No

Yes

Vote

Leading Replica?

Set Breakpoint on IP and catch up

Breakpoint Exception

Wait on Barrier

Compare <branches, IP>

Synced

Clear Breakpoint

UNSW
SYDNEY

# Performance: Microbenchmarks

| | Dhrystone | | Whetstone | |
|------|-----------|-------|-----------|-------|
| | Arm | x86 | Arm | x86 |
| Base | 146.1 | 108.1 | 108.9 | 120.3 |
| LC | 147.0 | 108.6 | 109.8 | 120.4 |
| CC | 153.4 | 111.9 | 133.5 | 143.0 |

Loosely-coupled

Closely-coupled

LC has low overhead for CPU-bound

LC has usually low inherent overhead for CPU-bound
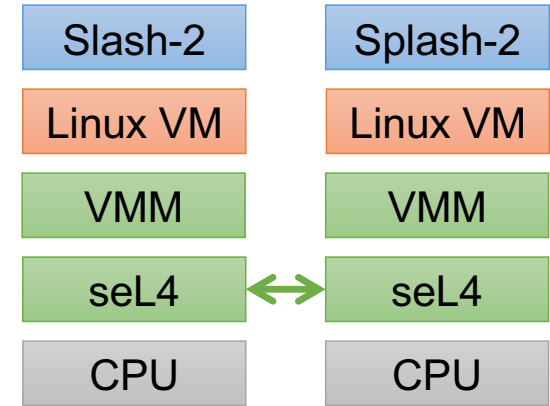
CC has high overhead for tight loops

UNSW SYDNEY

# Performance: SPLASH-2 on x86 VMs

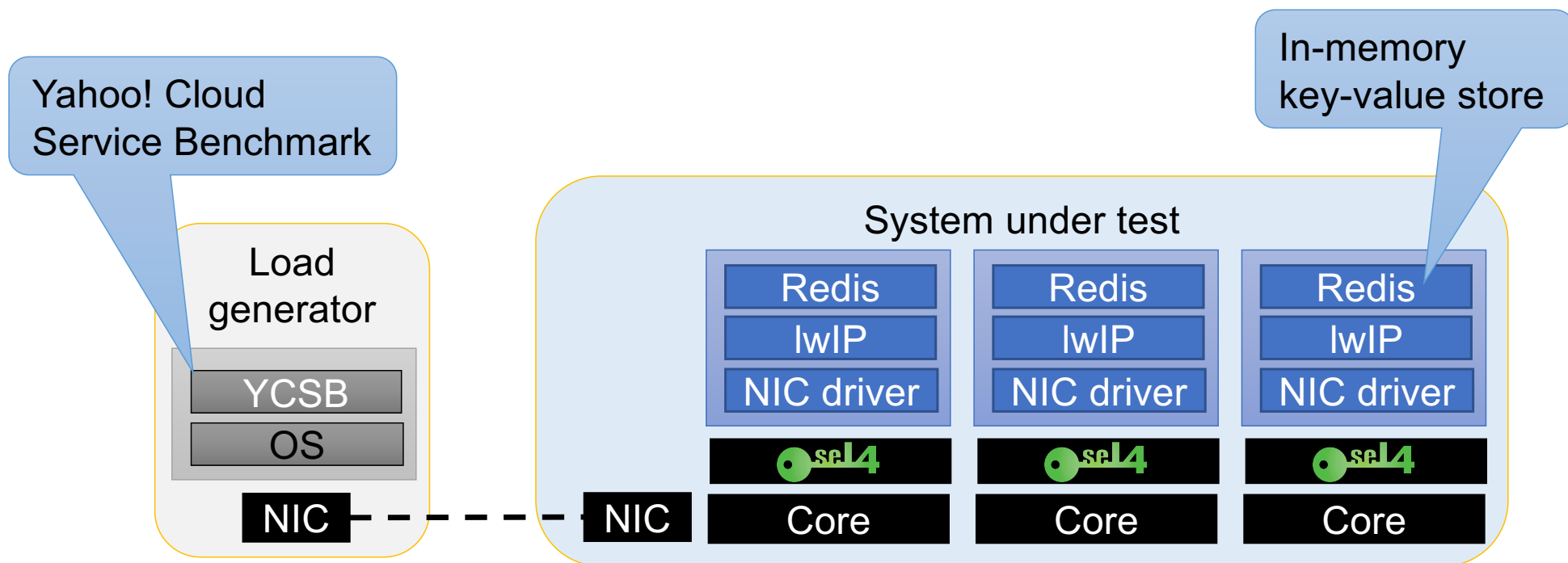| Name | N | Base | CC-D | Factor |
|------|-----|------|------|--------|
| BARNES | 30 | 61 | 93 | 1.52 |
| CHOLESKY | 300 | 66 | 792 | **12.08** |
| FFT | 100 | 64 | 142 | 2.22 |
| FFM | 20 | 76 | 160 | 2.11 |
| LU-C | 30 | 64 | 437 | 6.83 |
| LU-NC | 20 | 62 | 381 | 6.12 |
| OCEAN-C | 1000 | 64 | 173 | 2.71 |
| OCEAN-NC | 1000 | 65 | 171 | 2.65 |
| RADIOSITY | 25 | 66 | 75 | 1.12 |
| RADIX | 20 | 66 | 89 | 1.34 |
| RAYTRACE | 1000 | 60 | 65 | 1.09 |
| VOLREND | 100 | 86 | 133 | 1.54 |
| WATER-NS | 600 | 66 | 92 | 1.41 |
| WATER-S | 600 | 67 | 84 | 1.25 |

- Execution time in sec
- DMR configuration
- Base: unreplicated single-coreVM

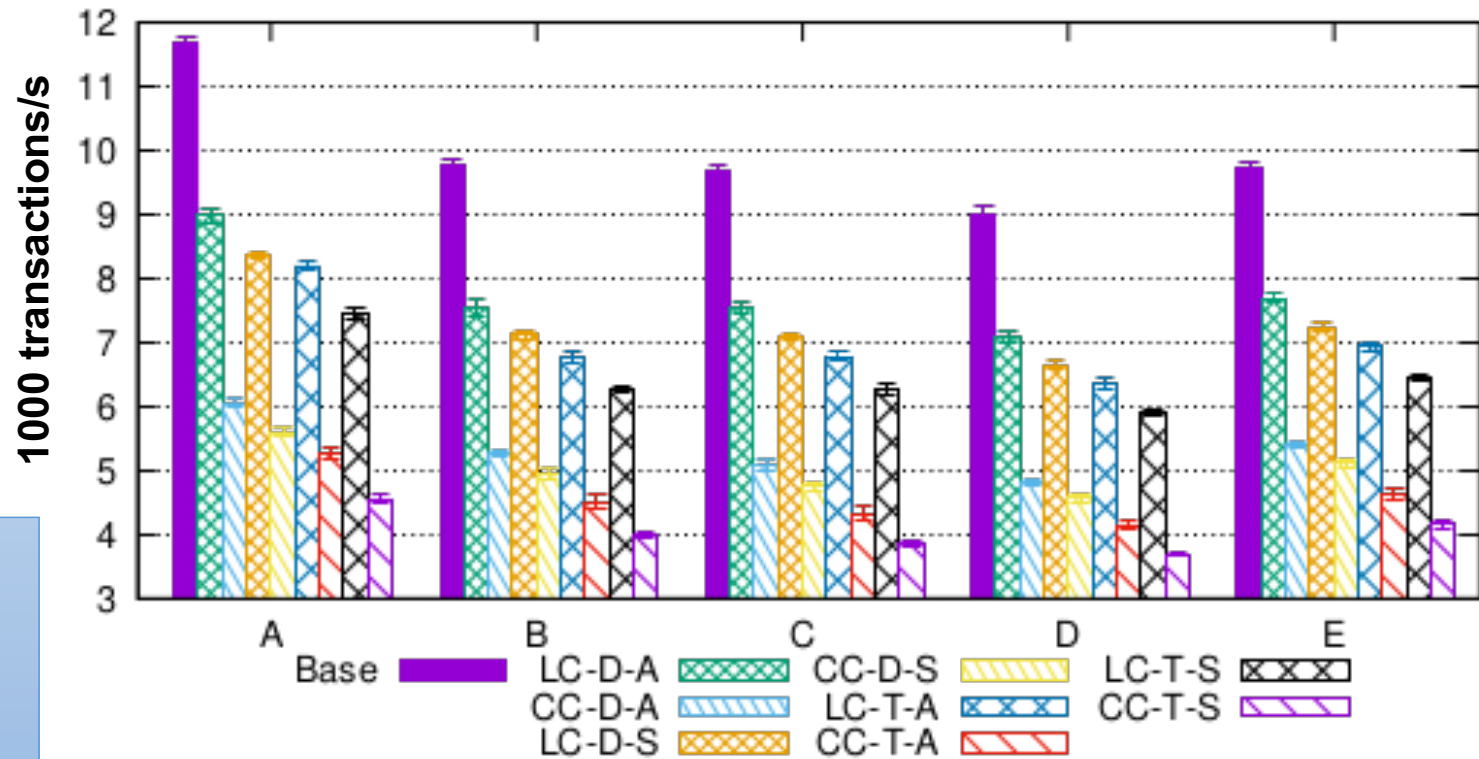Breakpoints in VM are expensive: trigger VM exits

| Slash-2 | Splash-2 |
|---------|----------|
| Linux VM | Linux VM |
| VMM | VMM |
| seL4 ↔ | seL4 |
| CPU | CPU |

Geometric mean overhead: 2.3×

UNSW SYDNEY

# Benchmark: Redis – YCSB

# Performance: Redis on Arm

LC: loosely-coupled

CC: closely-coupled

D: DMR

T: TMR

A: vote on interrupt

S: also vote on syscall



Overhead is 1.2–3 depending on configuration

# Error Detection on Arm

Not checksumming network data

Checksumming NW data

| | Base | LC-D | LC-T | LC-D-N | LC-T-N | CC-D | CC-T |
|---|---|---|---|---|---|---|---|
| Injected faults | 243k | 202k | 184k | 224k | 214k | 205k | 185k |
| YCSB corruptions | 647 | 3 | 1 | **381** | **299** | 3 | 0 |
| YCSB errors | 57 | 1 | 0 | 13 | 10 | 3 | 6 |
| User errors | 296 | 0 | 0 | 0 | 0 | 0 | 0 |
| Kernel exceptions | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Undetected | 1000 | 4 | 1 | 394 | 309 | 6 | 6 |
| RCoE detected | N/A | 996 | 999 | 606 | 691 | 994 | 994 |
| Observed errors | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 |

UNSW SYDNEY

# Comparison to Rad-Hardened Processor

| | Sabre Lite | RAD750 |
|---|---|---|
| Cores @ clock | 4 @ 800 MHz | 1 @ 133 MHz |
| Performance | 4 × 2,000 DMIPS | 240 DMIPS |
| Power | < 5 W | < 6 W |
| Energy Efficiency | 200 DMIPS/W | 40 DMIPS/W |
| Cost | $200 | $200,000 |
| Perf/Cost | 5 DMIPS/$ | 0.0002 DMIPS/$ |

2002 price

Assuming 2× overhead, TMR

[Shen et al., DSN'19]

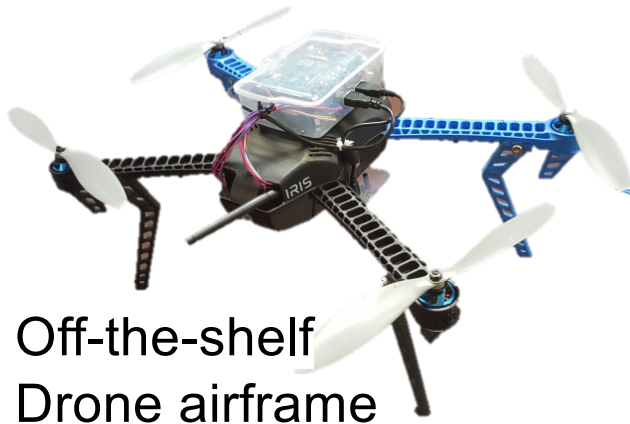UNSW SYDNEY

# Real-World Use
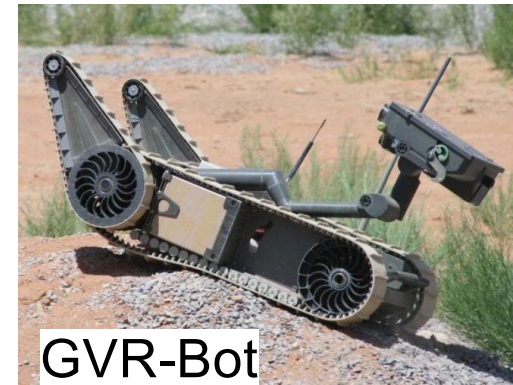
# DARPA HACMS

Unmanned Little Bird (ULB)

Retrofit existing system!

Autonomous trucks

Off-the-shelf Drone airframe

Develop technology

GVR-Bot

UNSW SYDNEY

# ULB Architecture

# Incremental Cyber Retrofit

Original Mission Computer

**Trusted**

Mission Manager

Crypto | Camera

Local NW | GPS

Ground Stn Link

Linux

➤

**Trusted**

Mission Manager

Crypto | Camera

Local NW | GPS

Ground Stn Link

Linux

Virt-Mach Monitor

seL4

➤

**Trusted**

GS Lk

Miss Mgr

Crypto

GPS

Local NW

Linux

VMM

Camera

Linux

VMM

seL4

UNSW SYDNEY

# Incremental Cyber Retrofit



COMP9242 2019T2 W09b: Local OS Research

UNSW SYDNEY

# Incremental Cyber Retrofit

[Klein et al, CACM, Oct'18]

Original Mission Computer

**Trusted**

| Mission Manager |
| Crypto | Camera |
| Local NW | GPS |
| Ground Stn Link |

Linux

Cyber-secure Mission Computer

**Trusted**

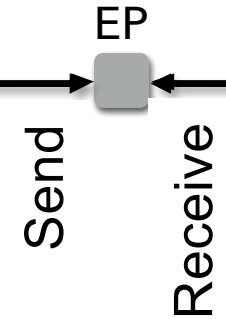| Crypto | Mission Mngr |
| Local NW | Comms |

seL4

Cam-era

Linux

GPS

VMM

UNSW SYDNEY
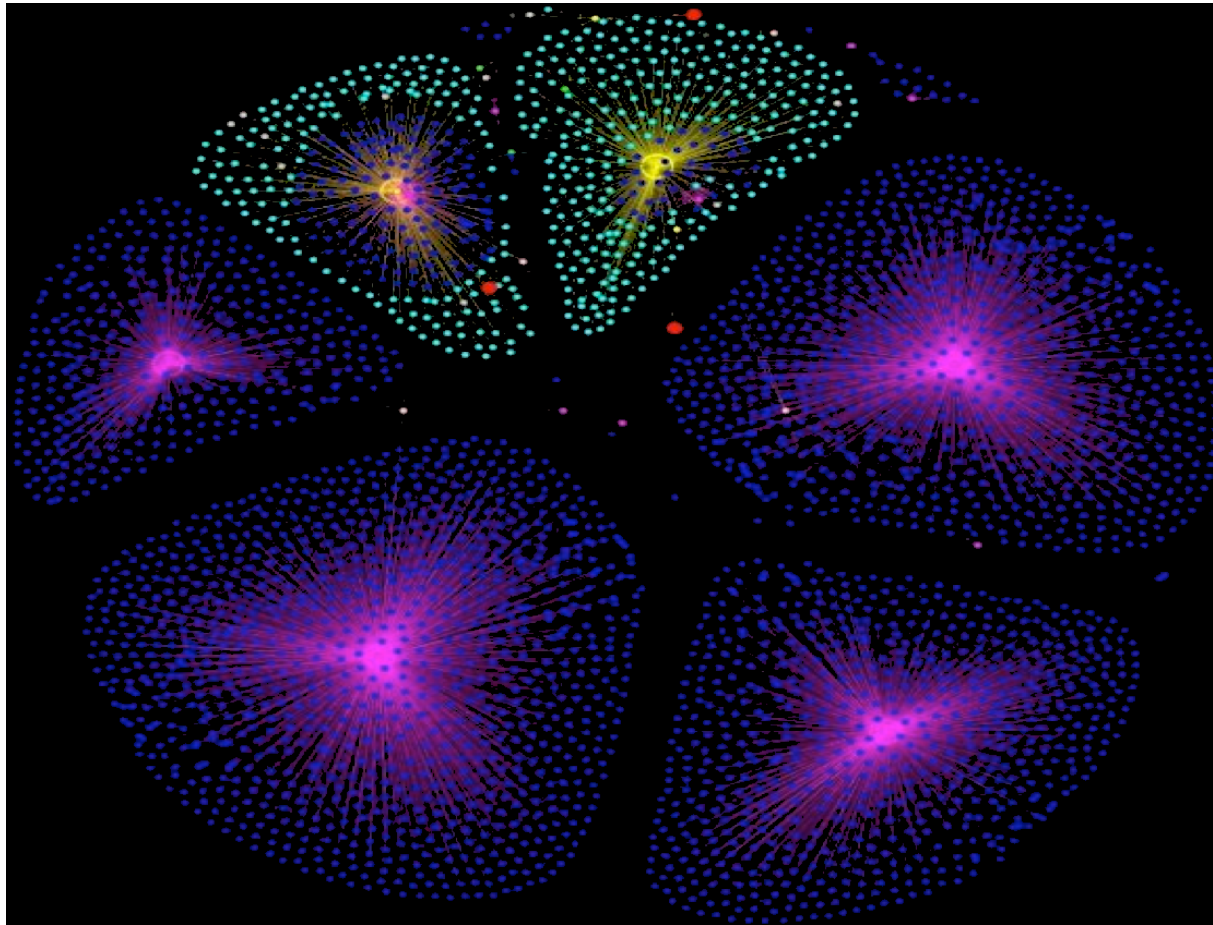
# Issue: Capabilities are Low-Level



>50 capabilities for trivial program!

UNSW
SYDNEY

# Simple But Non-Trivial System

# Component Middleware: CAmkES

Higher-level abstractions of low-level seL4 constructs

# HACMS UAV Architecture



Security enforcement:
Linux only sees
encrypted data

Radio Driver

Data Link

Crypto

CAN Driver

Uncritical/
untrusted,
contained

Wifi

Camera

Linux

UNSW
SYDNEY

# Enforcing the Architecture

UNSW SYDNEY

# Architecture Analysis

Analysis Tools

Safety ✔

Eclipse-based IDE → Design → AADL

Architecture analysis and design language

AADL → Generate → CAmkES → Generate → .h, .c → Compile → Binary

UNSW SYDNEY

# Real-World Use
## Courtesy Boeing, DARPA

UNSW
SYDNEY