# COMP9242 Advanced OS

S2/2016 W12: **Local Systems 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 Australia"*

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

**UNSW**
AUSTRALIA

# Present Systems are *NOT* Trustworthy!

HIJACKED!

HACKED!

HACKED!

pwned!

HACKED!

**Yet they are expensive:**

- $1,000 per line of code for "high-assurance" software!
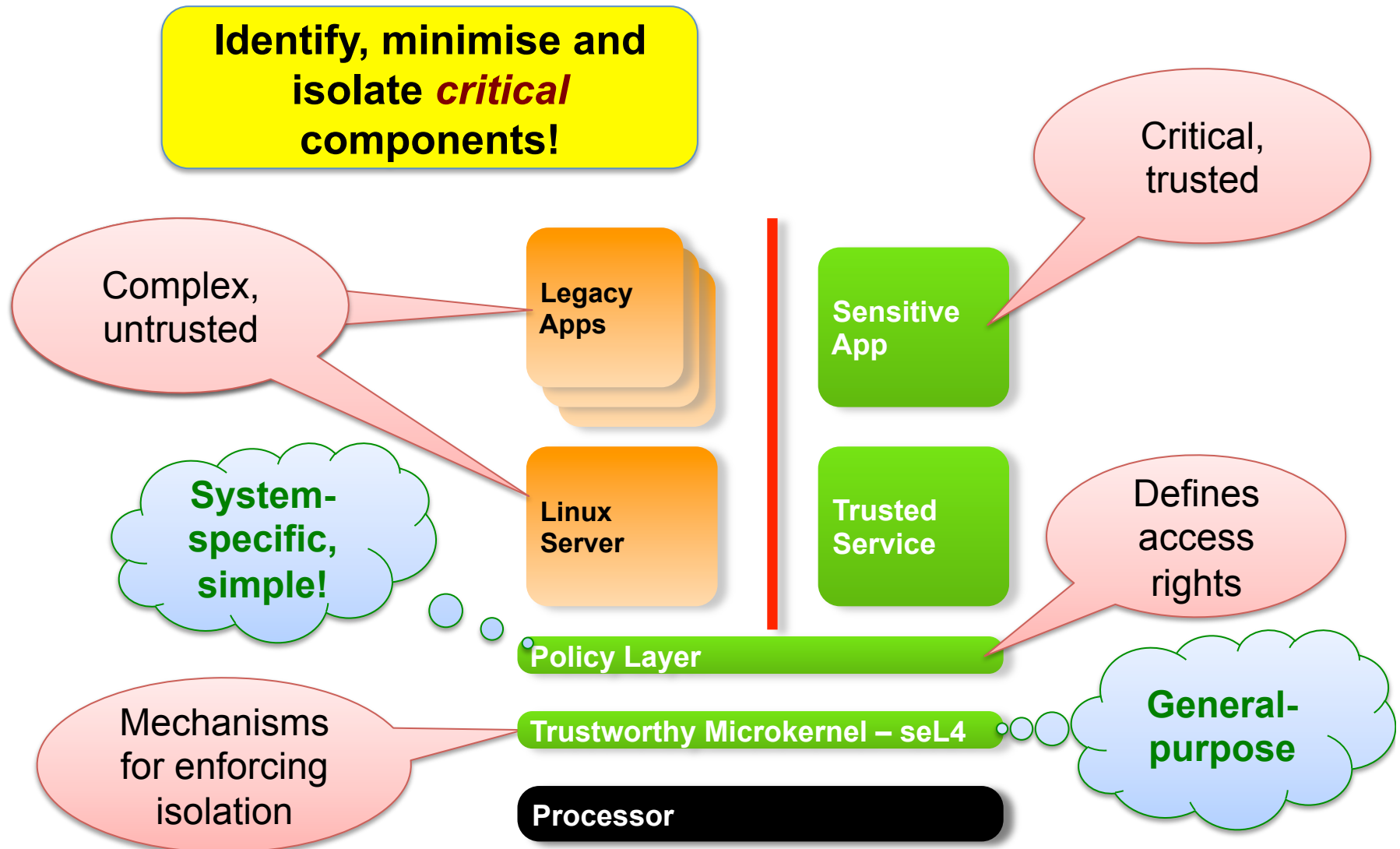
       UNSW
AUSTRALIA

# Trustworthy Systems Vision

Suitable for real-world systems

**We will change the *practice* of designing and implementing critical systems, using rigorous approaches to achieve *true trustworthiness***

Hard *guarantees* on safety/security/reliability

UNSW
AUSTRALIA

# Isolation is Key!

**Identify, minimise and isolate *critical* components!**

Critical, trusted

Complex, untrusted

**Legacy Apps**

**Sensitive App**

**System-specific, simple!**

**Linux Server**

**Trusted Service**

Defines access rights

**Policy Layer**

Mechanisms for enforcing isolation

**Trustworthy Microkernel – seL4**

**General-purpose**

**Processor**

UNSW
AUSTRALIA

# Trustworthy Systems Agenda

1. **Dependable microkernel (seL4) as a rock-solid base**
   - Formal specification of functionality
   - Proof of functional correctness of implementation
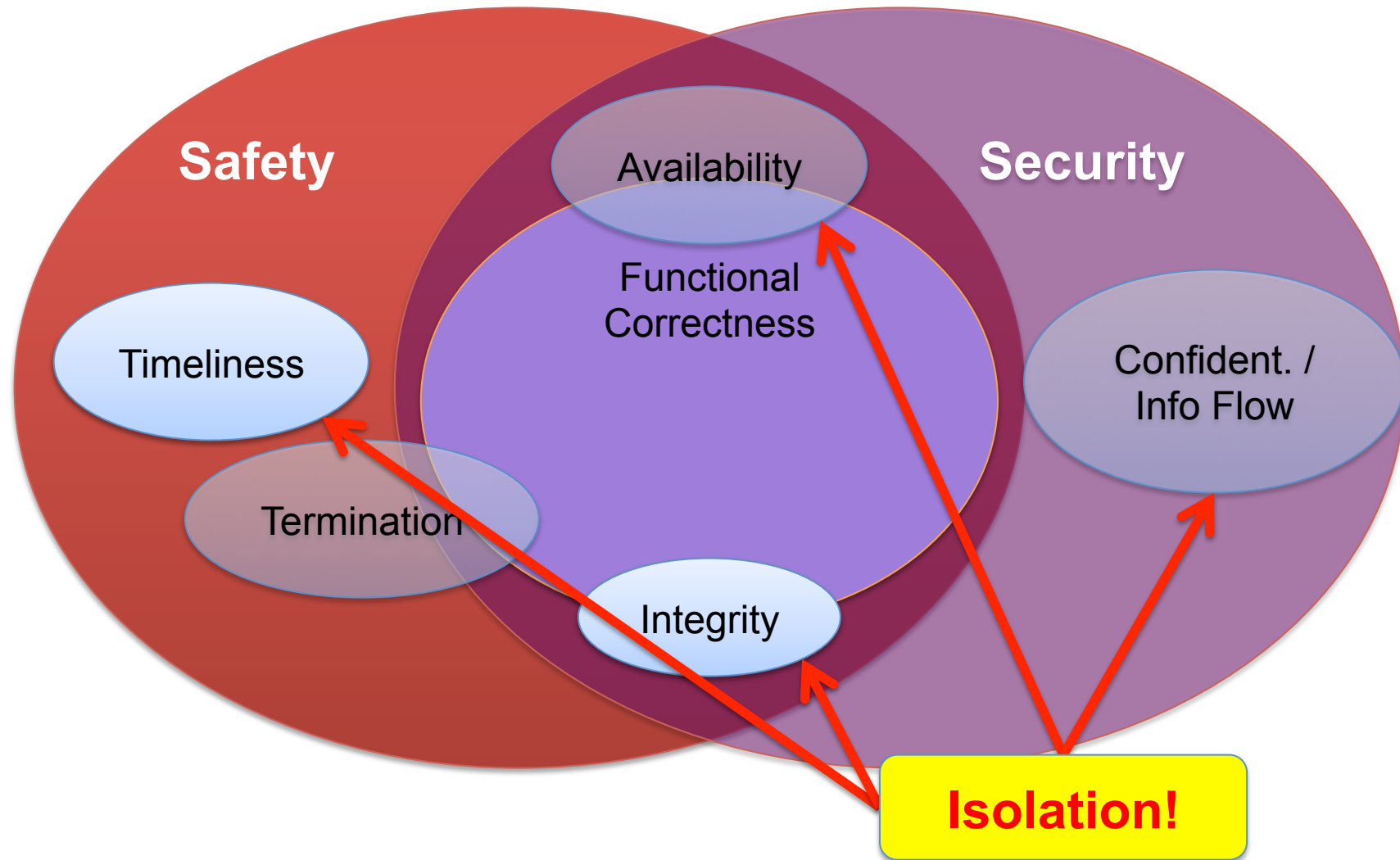   - Proof of safety/security properties

   **DONE**

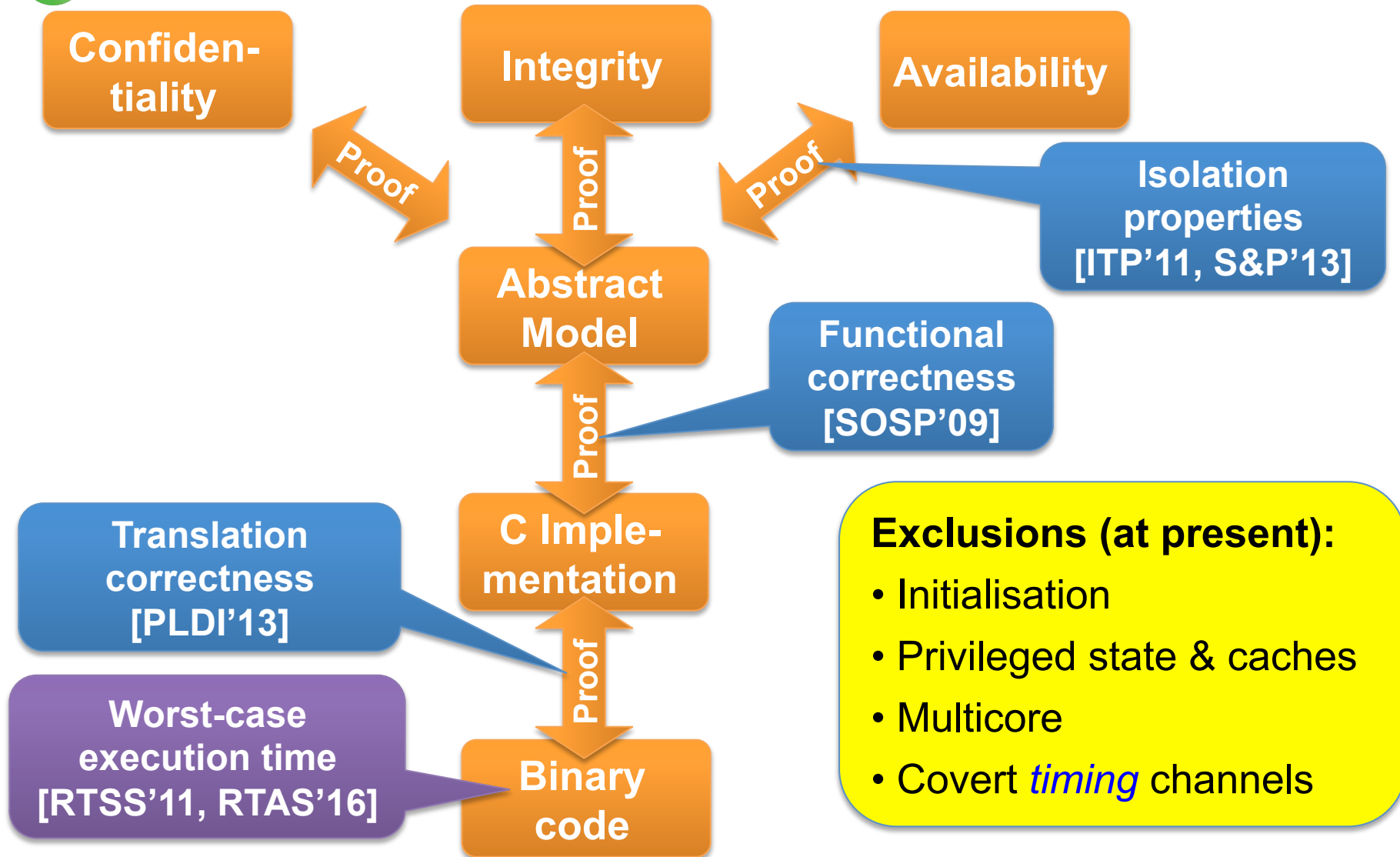2. **Lift microkernel guarantees to whole system**
   - Use kernel correctness and integrity to guarantee critical functionality
   - Ensure correctness of balance of trusted computing base
   - Prove dependability properties of complete system
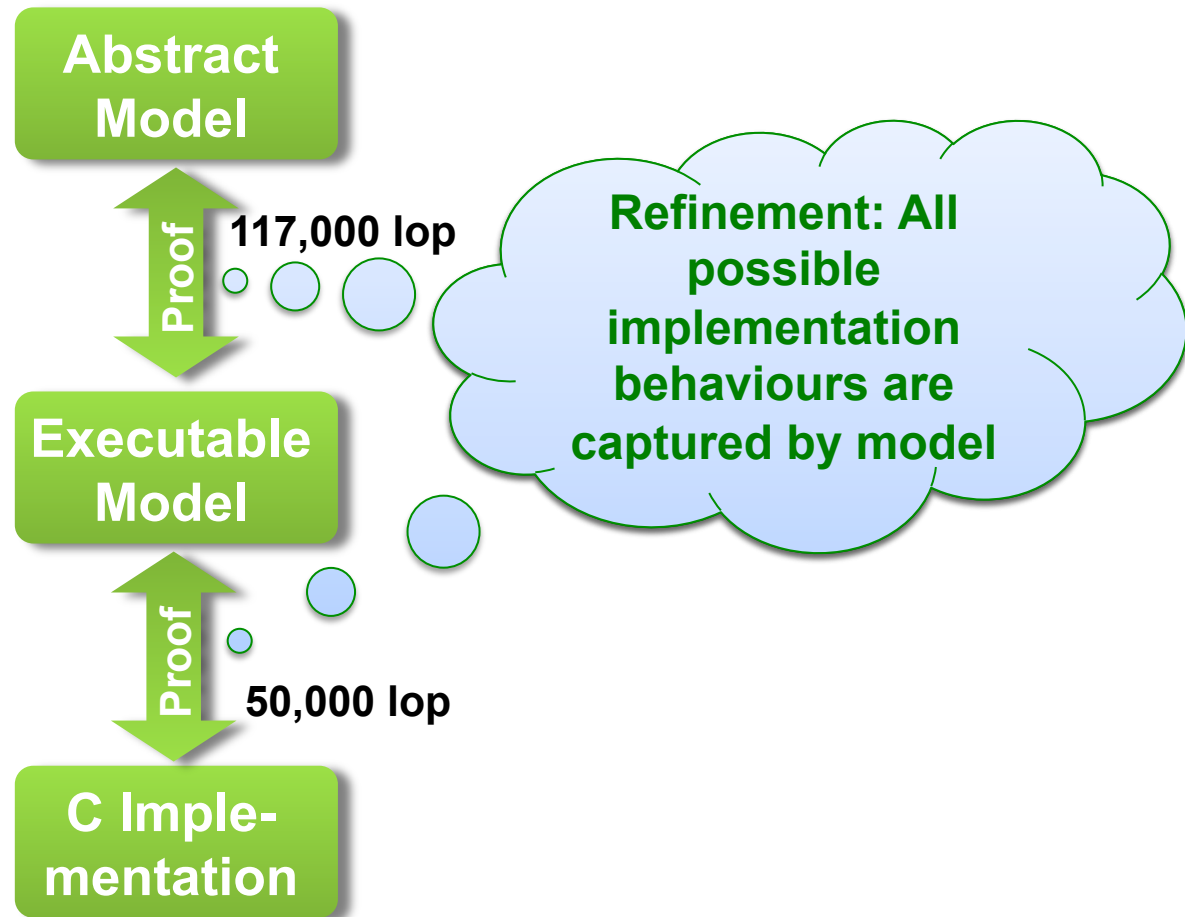     - despite 99 % of code untrusted!

© 2016 Gernot Heiser. Distributed under CC Attribution License

**UNSW** AUSTRALIA

# Requirements for Trustworthy Systems



COMP9242 S2/2016 W12 © 2016 Gernot Heiser. Distributed under CC Attribution License

# Provable Security and Safety

**Confiden-tiality**

**Integrity**

**Availability**

Proof

Proof

Proof

**Isolation properties [ITP'11, S&P'13]**

**Abstract Model**

Proof

**Functional correctness [SOSP'09]**

**Translation correctness [PLDI'13]**

**C Imple-mentation**

Proof

**Exclusions (at present):**
- Initialisation
- Privileged state & caches
- Multicore
- Covert *timing* channels

**Worst-case execution time [RTSS'11, RTAS'16]**

**Binary code**

© 2016 Gernot Heiser. Distributed under CC Attribution License

**UNSW** AUSTRALIA

# Proving Functional Correctness

**Abstract Model**

↑ Proof

**117,000 lop**

**Executable Model**

↑ Proof

**50,000 lop**

**C Imple-mentation**

**Refinement: All possible implementation behaviours are captured by model**

UNSW AUSTRALIA

# Proving Functional Correctness

```
constdefs
   schedule :: "unit s_monad"
   "schedule ≡ do
       threads ← allActiveTCBs;
       thread ← select threads;
       do_machine_op flushCaches OR return ();
       modify (\s . s ⦇ cur_thread := thread ⦈)
```

```
schedule :: Kernel ()
schedule = do
        action <- getSchedulerAction
        case action of
```

```
void
setPriority(tcb_t *tptr, prio_t prio) {
    prio_t oldprio;

    if(thread_state_get_tcbQueued(tptr->tcbState)) {
        oldprio = tptr->tcbPriority;
        ksReadyQueues[oldprio] = tcbSchedDequeue(tptr, ksReadyQueues[
        if(isRunnable(tptr)) {
            ksReadyQueues[prio] = tcbSchedEnqueue(tptr, ksReadyQueues
        }
        else {
            thread_state_ptr_set_tcbQueued(&tptr->tcbState, false);
        }
    }

    tptr->tcbPriority = prio;
}

void
yieldTo(tcb_t *target) {
    target->tcbTimeSlice += ksCurThread->tcbTimeSlice;
```

```
rThread
Slice curThread
e == 0) chooseThread
```

UNSW
AUSTRALIA

**MIT Technology Review**

# 10 BREAKTHROUGH TECHNOLOGIES

2011

Share

---

## Crash-Proof Code

*Making critical software safer*

7 comments
**WILLIAM BULKELEY**
*May/June 2011*

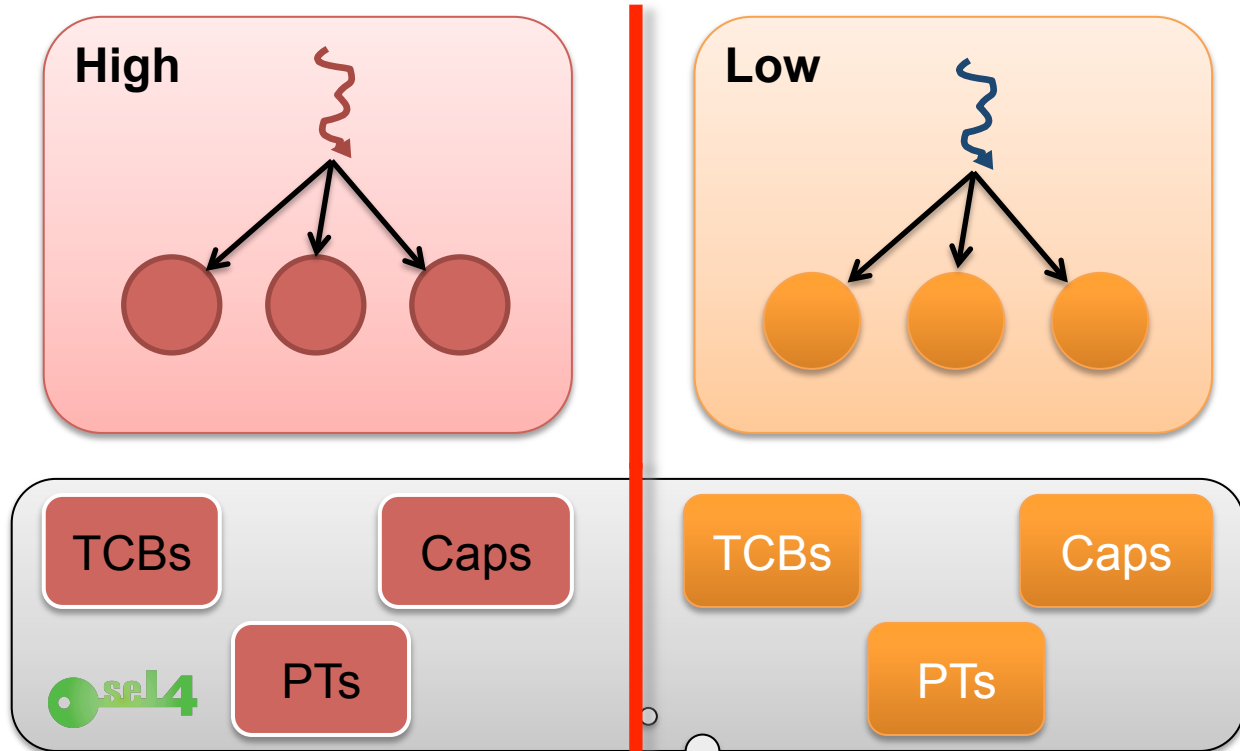# seL4 Formal Verification Summary

**Kinds of properties proved**

- Behaviour of C code is fully captured by abstract model
- Behaviour of C code is fully captured by executable model
- Kernel never fails, behaviour is always well-defined
  - assertions never fail
  - will never de-reference null pointer
  - cannot be subverted by misformed input
- All syscalls terminate, reclaiming memory is safe, ...
- Well typed references, aligned objects, kernel always mapped…
- Access control is decidable

> Can prove further properties on abstract level!

**Did you find bugs?**

- During (very shallow) testing: 16
- During verification: 460
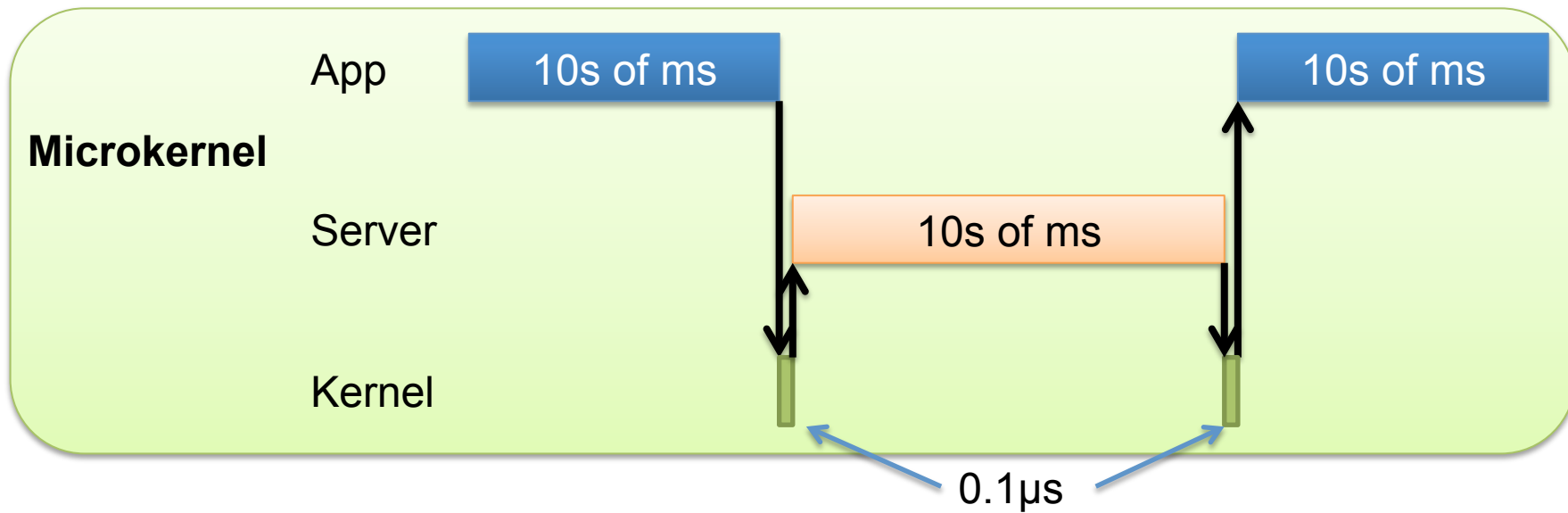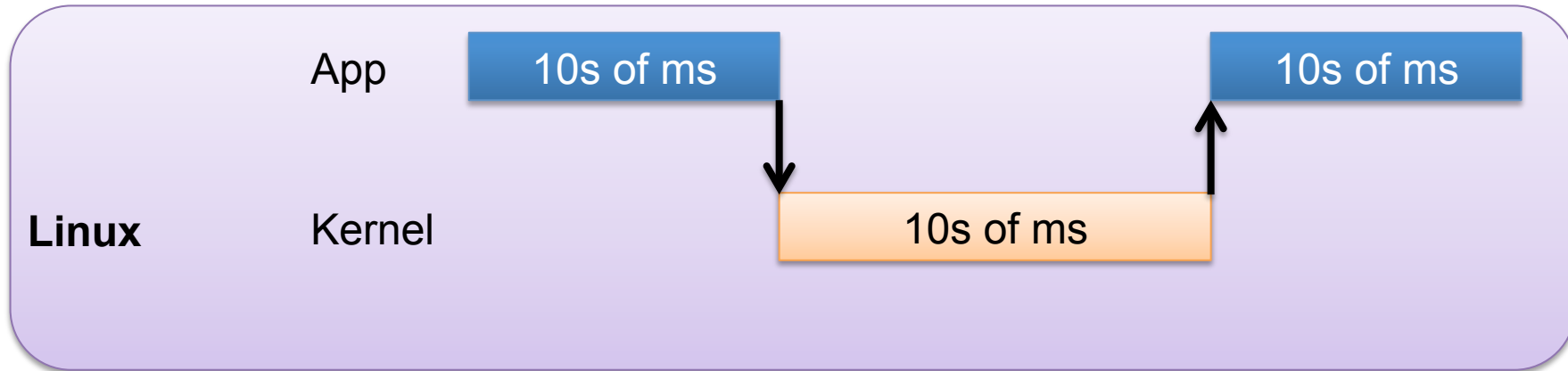  - 160 in C, ~150 in design, ~150 in spec

   UNSW AUSTRALIA

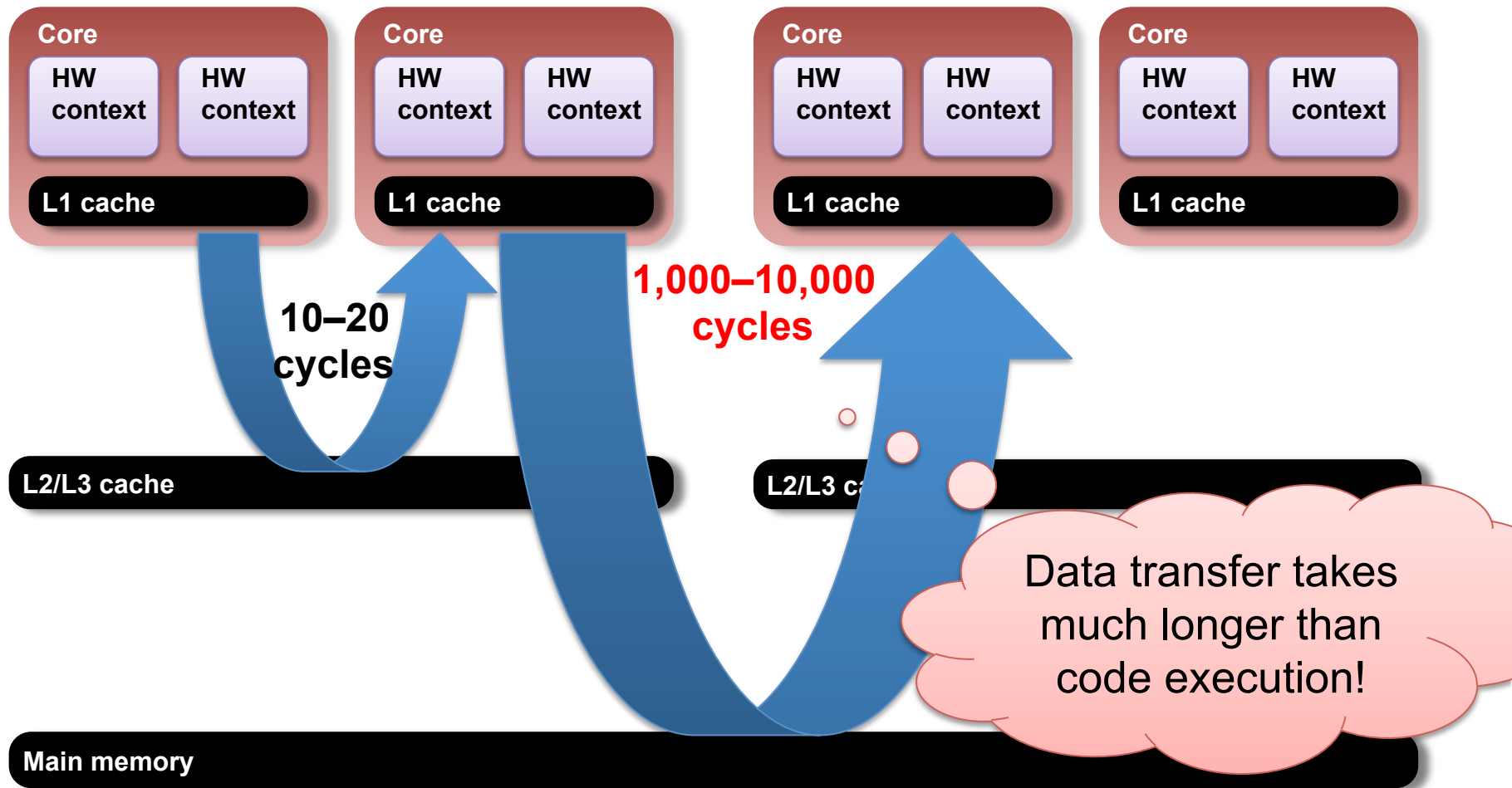# Isolation Goes Deep

**High**

**Low**

TCBs

Caps

PTs

TCBs

Caps

PTs

Kernel data partitioned like user data

© 2016 Gernot Heiser. Distributed under CC Attribution License

# Multicore

# Microkernel vs Linux Execution

**Linux**

App | 10s of ms | | 10s of ms
Kernel | 10s of ms

**Microkernel**

App | 10s of ms | | 10s of ms
Server | | 10s of ms
Kernel

0.1µs

UNSW
AUSTRALIA

# Cache Line Migration Latencies



COMP9242 S2/2016 W12    © 2016 Gernot Heiser. Distributed under CC Attribution License

# Cost of Locking

**X86 (Haswell)**



Cycles

424 | 436 | 508 | 496

No lock | Big lock | Fine-grained locking | Transactions
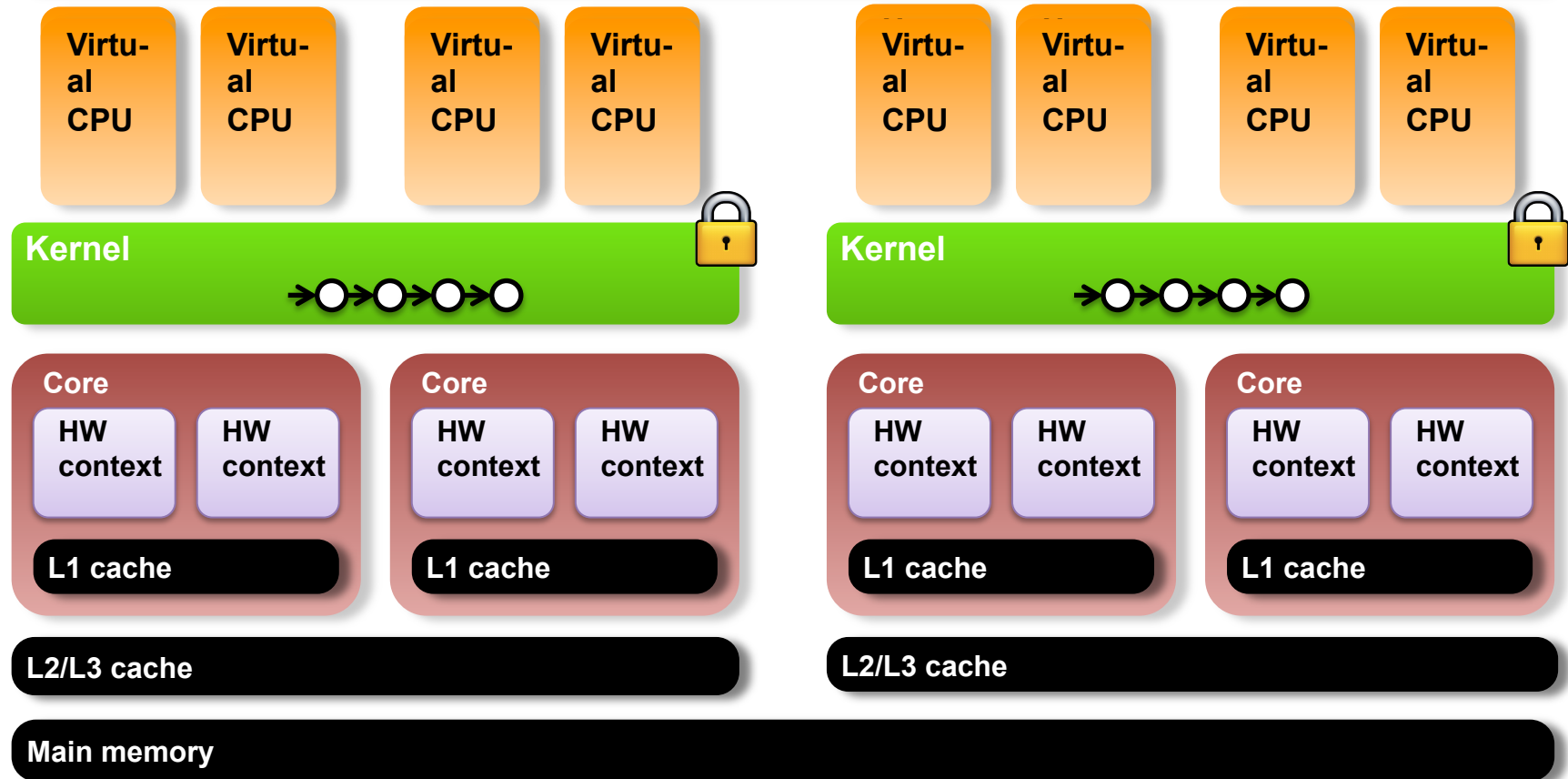
**ARM A9**



316 | 390 | 548

No lock | Big lock | Fine-grained locking

Locks have a cost –
significant in a fast microkernel!

UNSW
AUSTRALIA

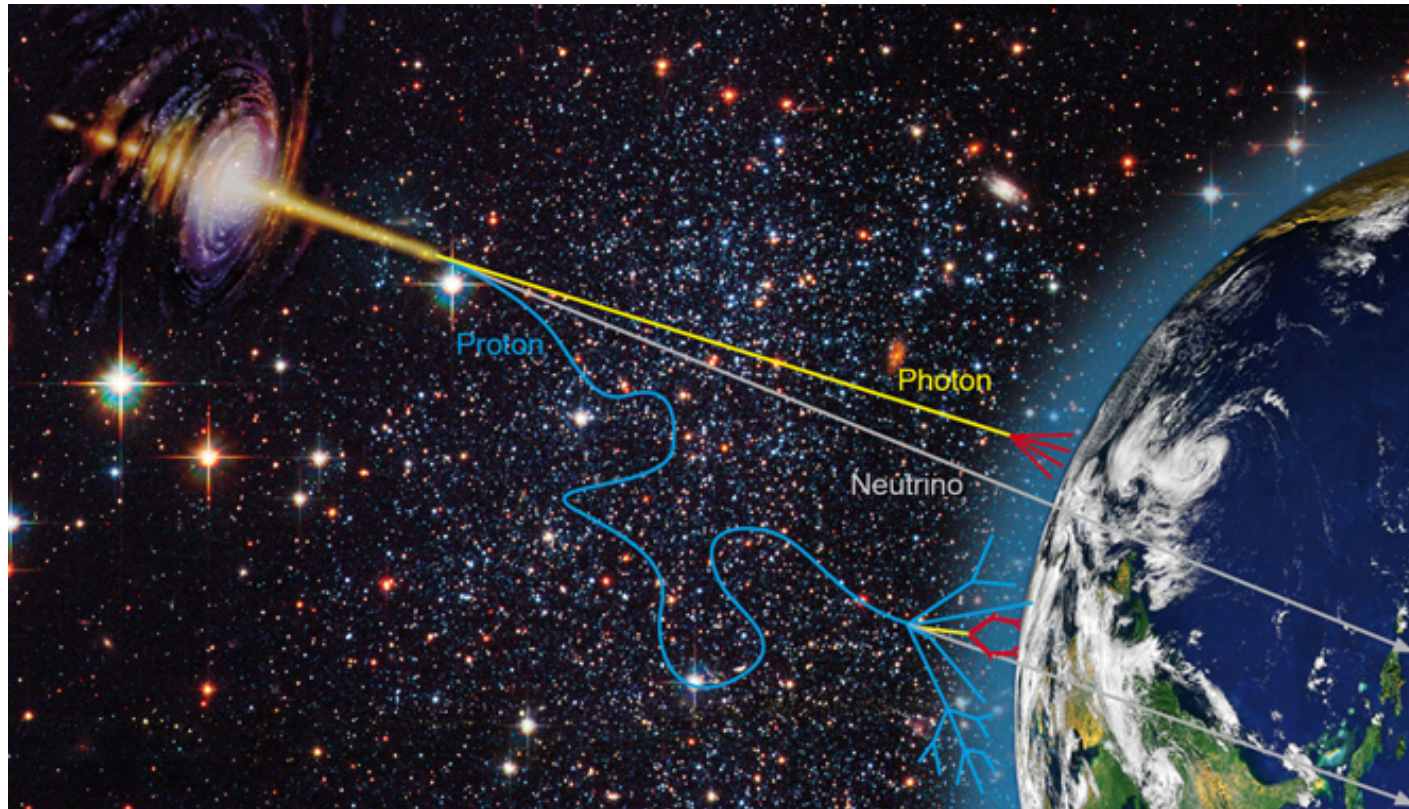# Multicore Design: Clustered Multikernel

**NUMA-aware Linux**

| Virtu-al CPU | Virtu-al CPU | Virtu-al CPU | Virtu-al CPU | | Virtu-al CPU | Virtu-al CPU | Virtu-al CPU | Virtu-al CPU |

**Kernel**

**Kernel**

**Core**
HW context | HW context

**Core**
HW context | HW context

**Core**
HW context | HW context

**Core**
HW context | HW context

**L1 cache** | **L1 cache** | **L1 cache** | **L1 cache**

**L2/L3 cache** | **L2/L3 cache**

**Main memory**

© 2016 Gernot Heiser. Distributed under CC Attribution License

UNSW
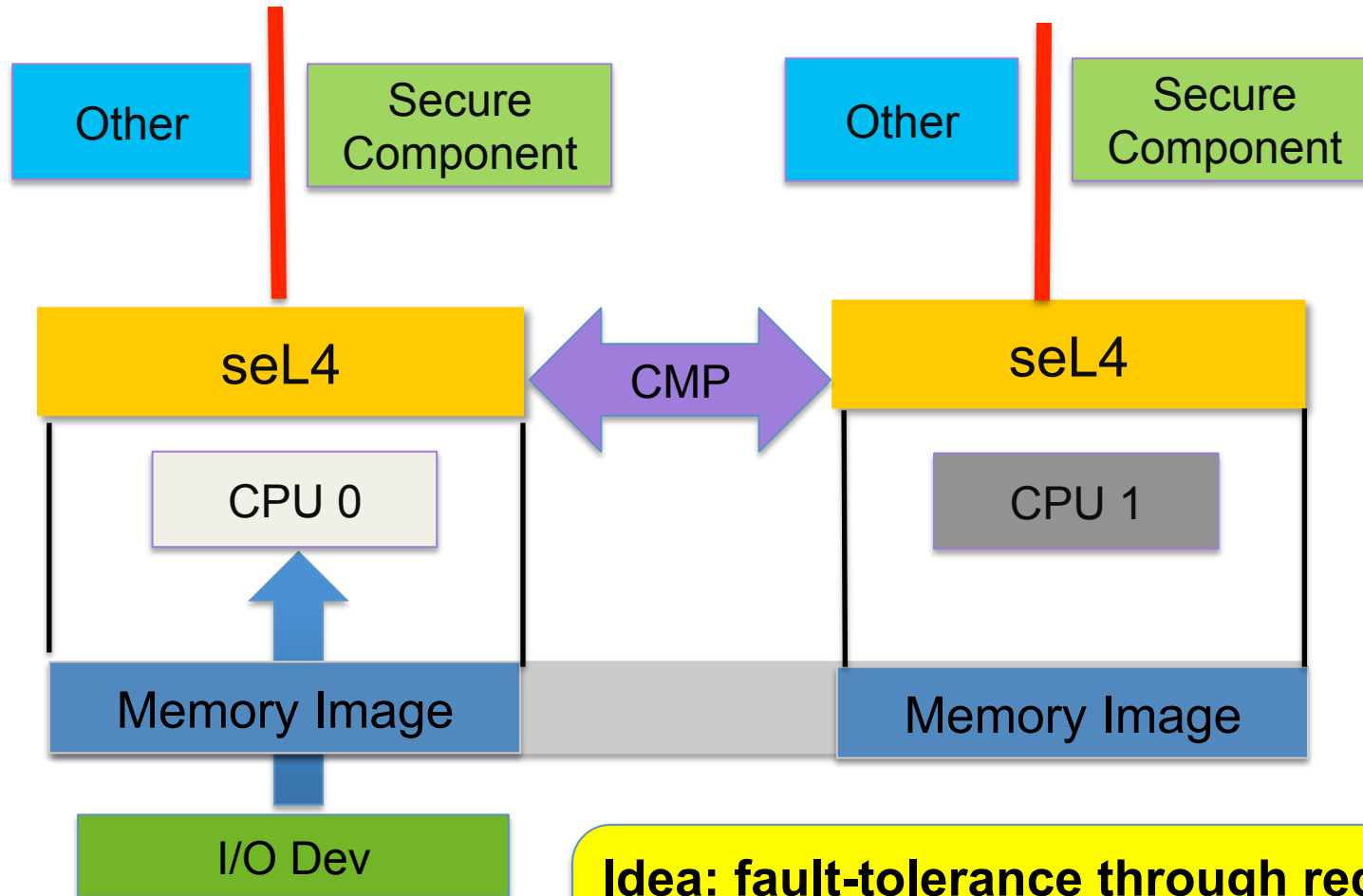AUSTRALIA

Big-Lock Scalability

# Hardware Faults

# How About Hardware Faults?



- Single-event upset: Random (transient) bit-flips due to cosmic rays, natural radioactivity

- May break "proved" isolation

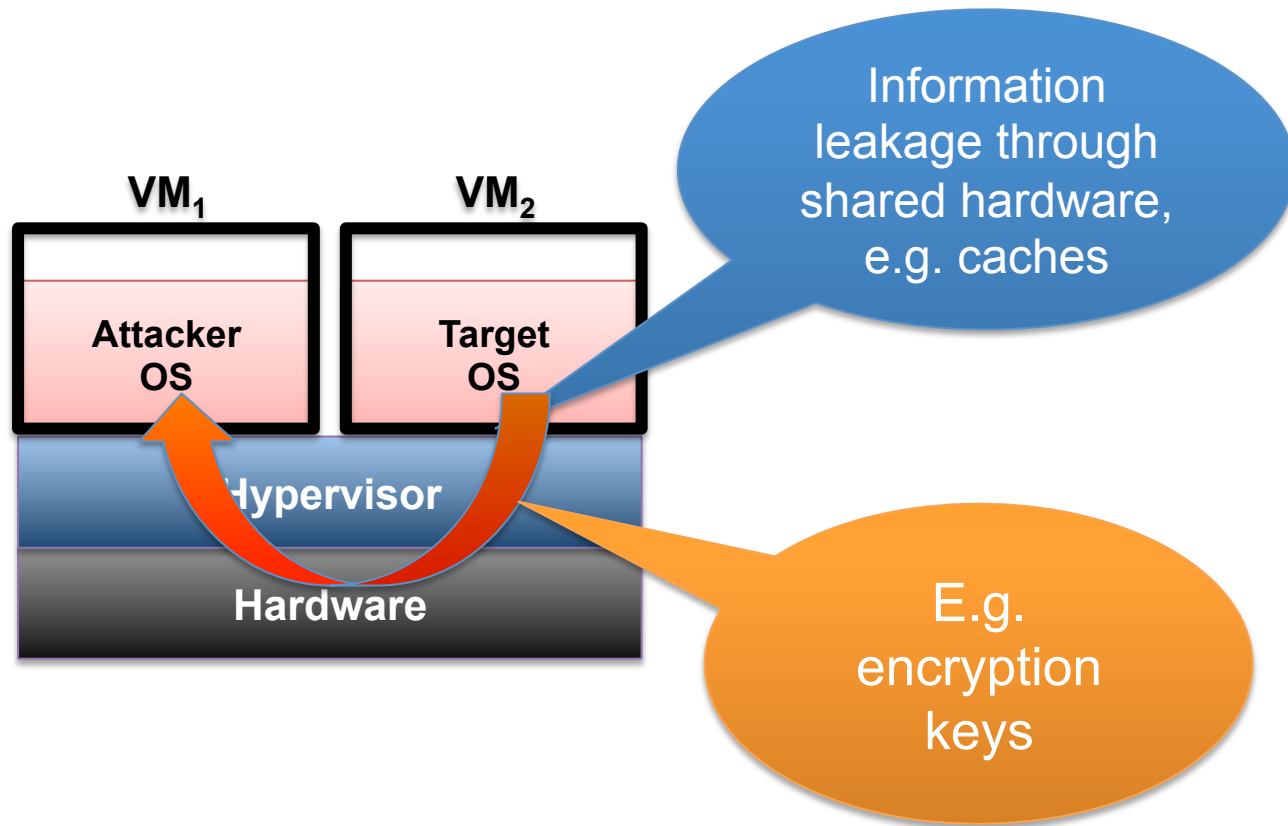UNSW
AUSTRALIA

# Redundant Execution



**Idea: fault-tolerance through redundancy**
- Compare & vote at kernel entry/exit
- Work in progress (Yanyan's PhD)

# Side Channels

# Side Channel Attacks



VM₁ — Attacker OS

VM₂ — Target OS

Hypervisor

Hardware

Information leakage through shared hardware, e.g. caches

E.g. encryption keys

UNSW
AUSTRALIA

# Types of Side Channels

## Storage Channels

- Use some shared state

- Could be inside the OS/ hypervisor
  - Eg existence of a file
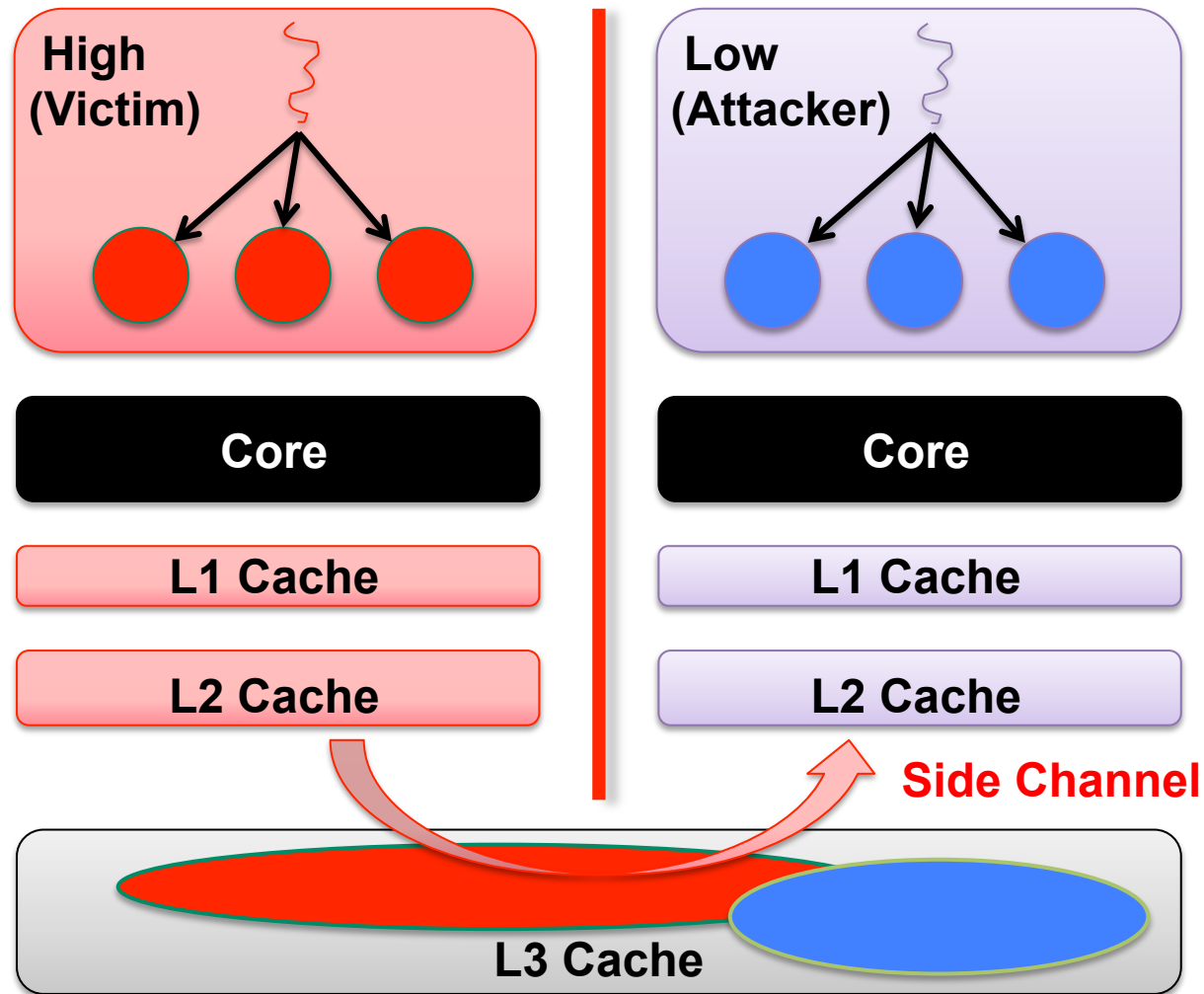  - Eg accessibility of an object

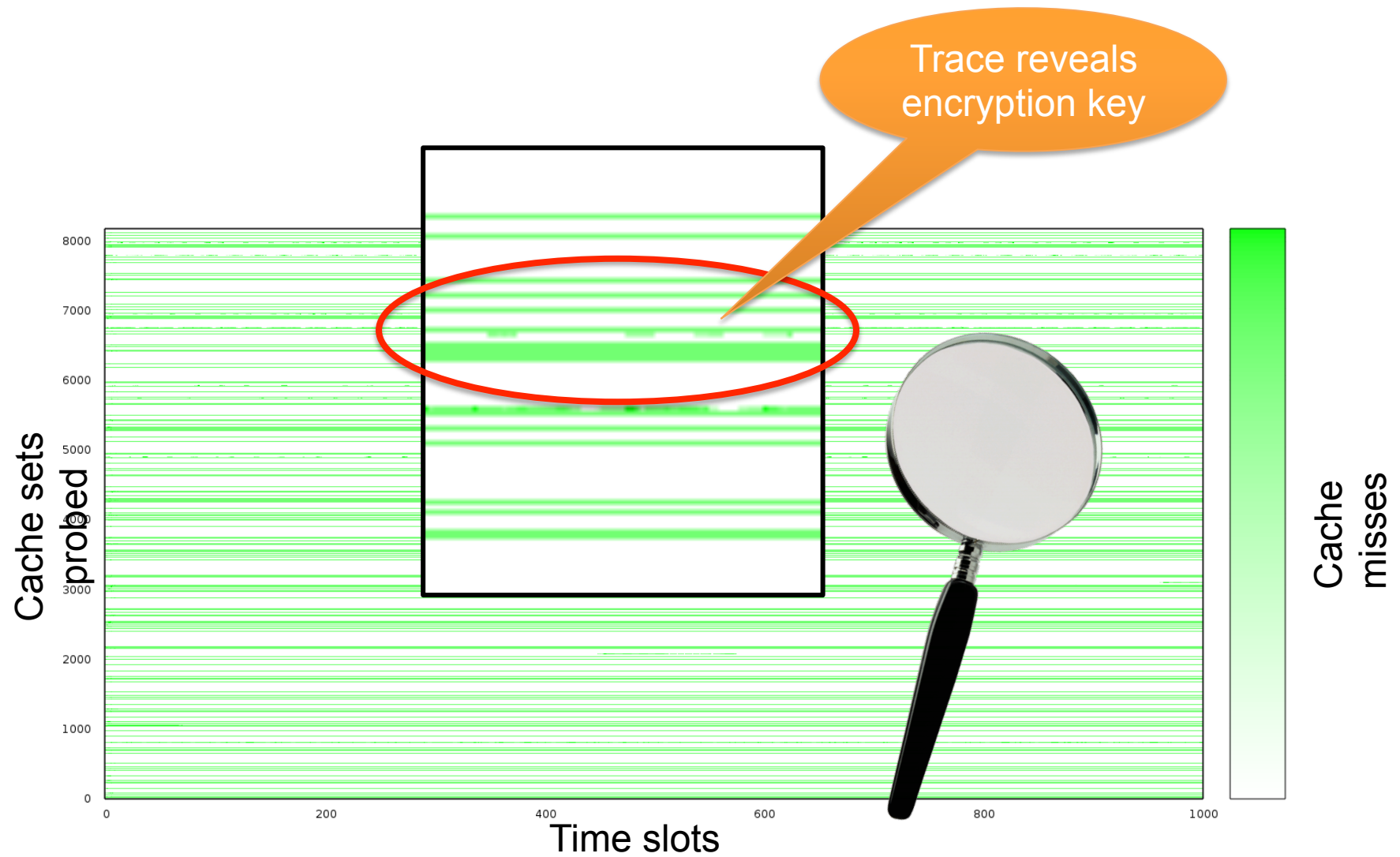**seL4: The world's only OS proved free of storage channels!**

## Timing Channels

- Observe timing of events

- Eg memory access latency
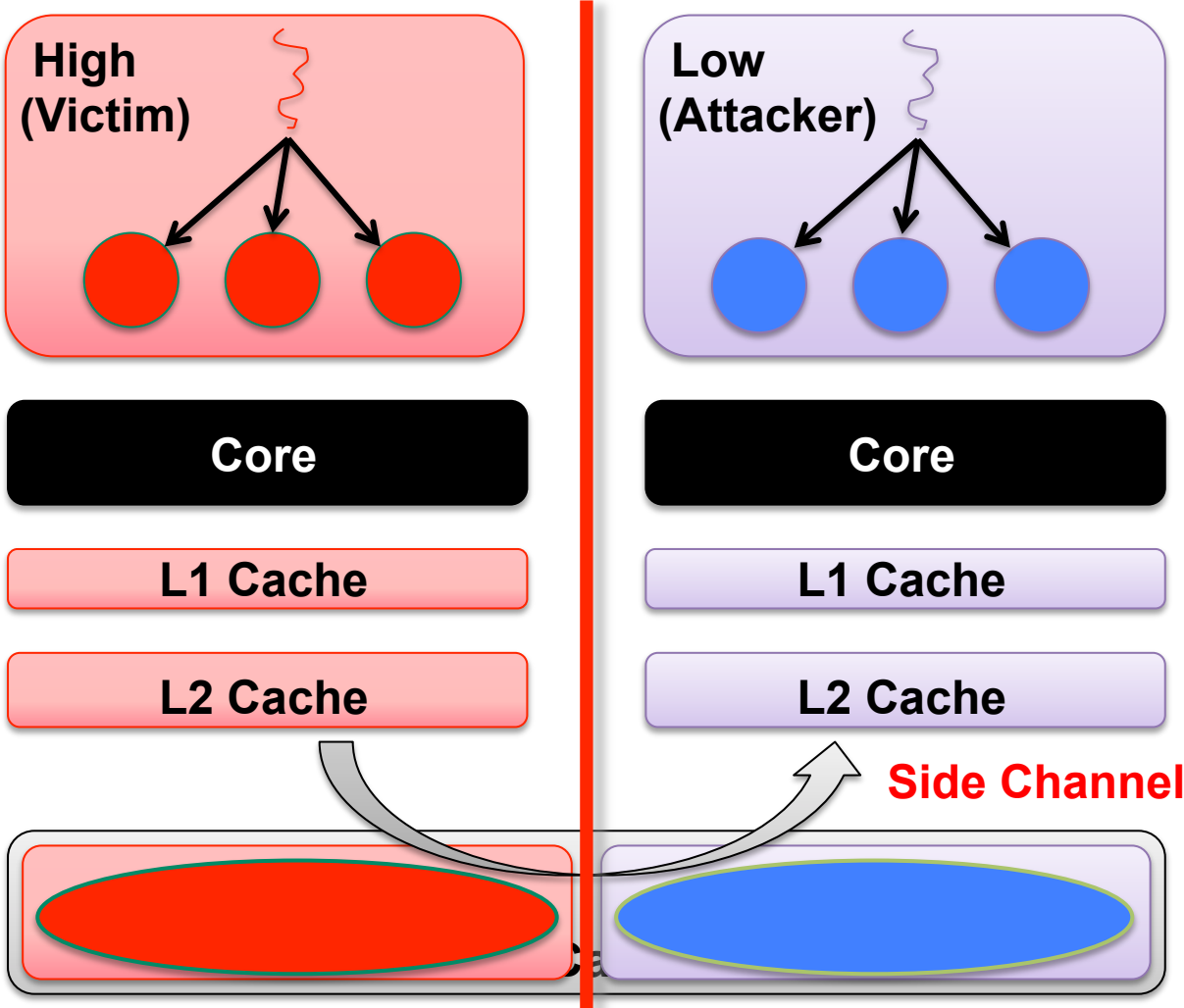  - Senses victim's cache footprint

How about timing channels?

© 2016 Gernot Heiser. Distributed under CC Attribution License

UNSW
AUSTRALIA

# Timing Side-Channel Attack in Public Cloud



© 2016 Gernot Heiser. Distributed under CC Attribution License

# Analysing Memory Access Latency



Trace reveals encryption key

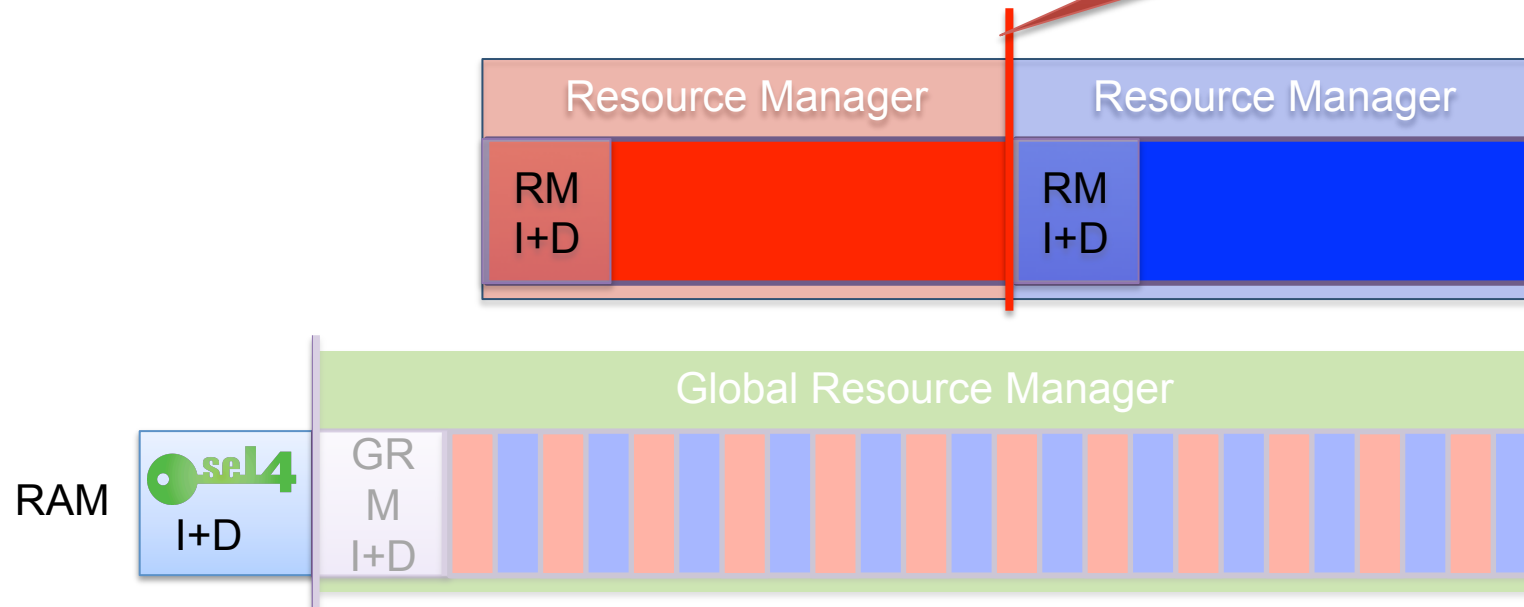© 2016 Gernot Heiser. Distributed under CC Attribution License

# Mitigation: Partition Cache (Colouring)

# Colouring the System is Easy

System permanently coloured

Partitions restricted to coloured memory
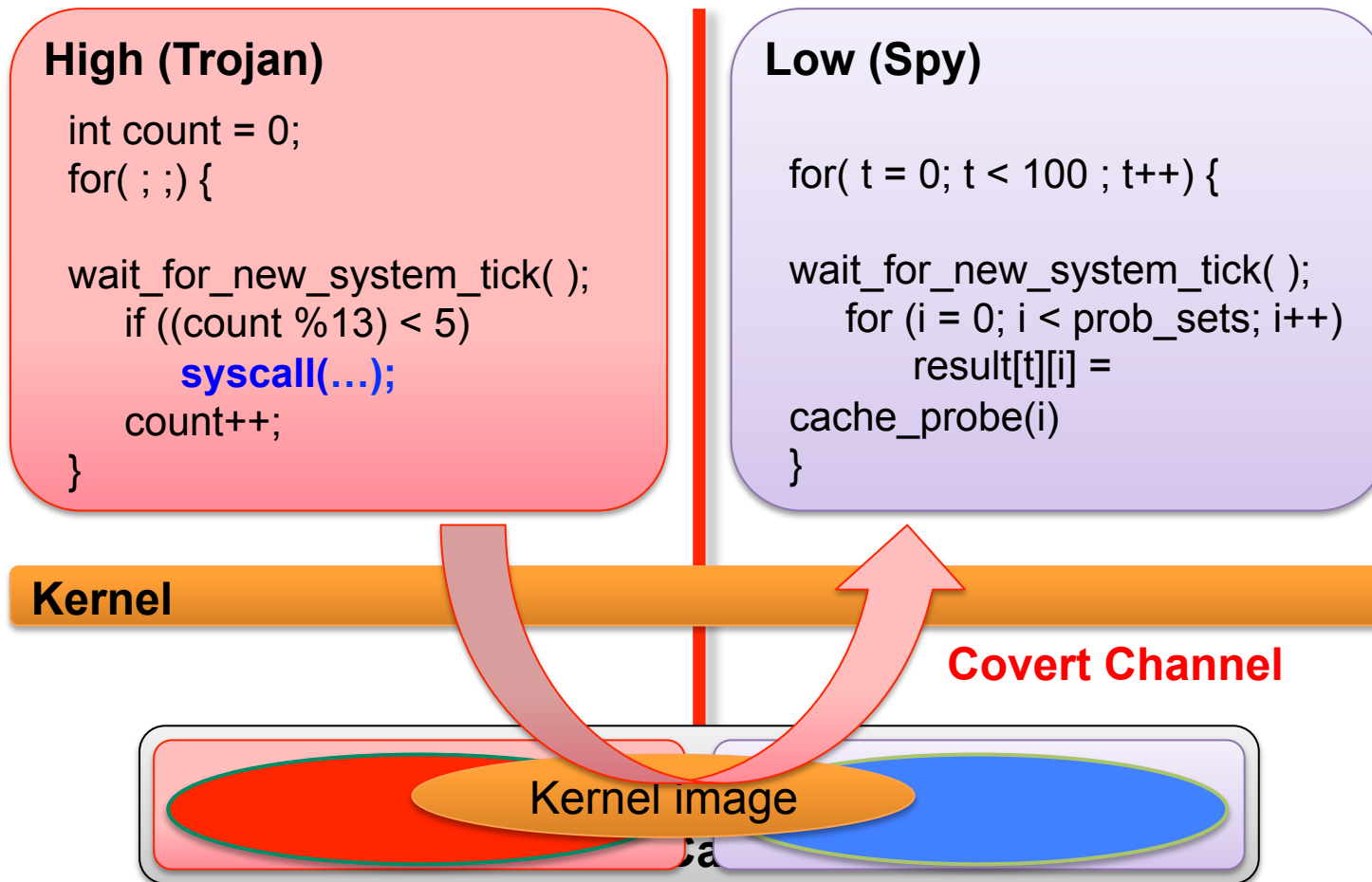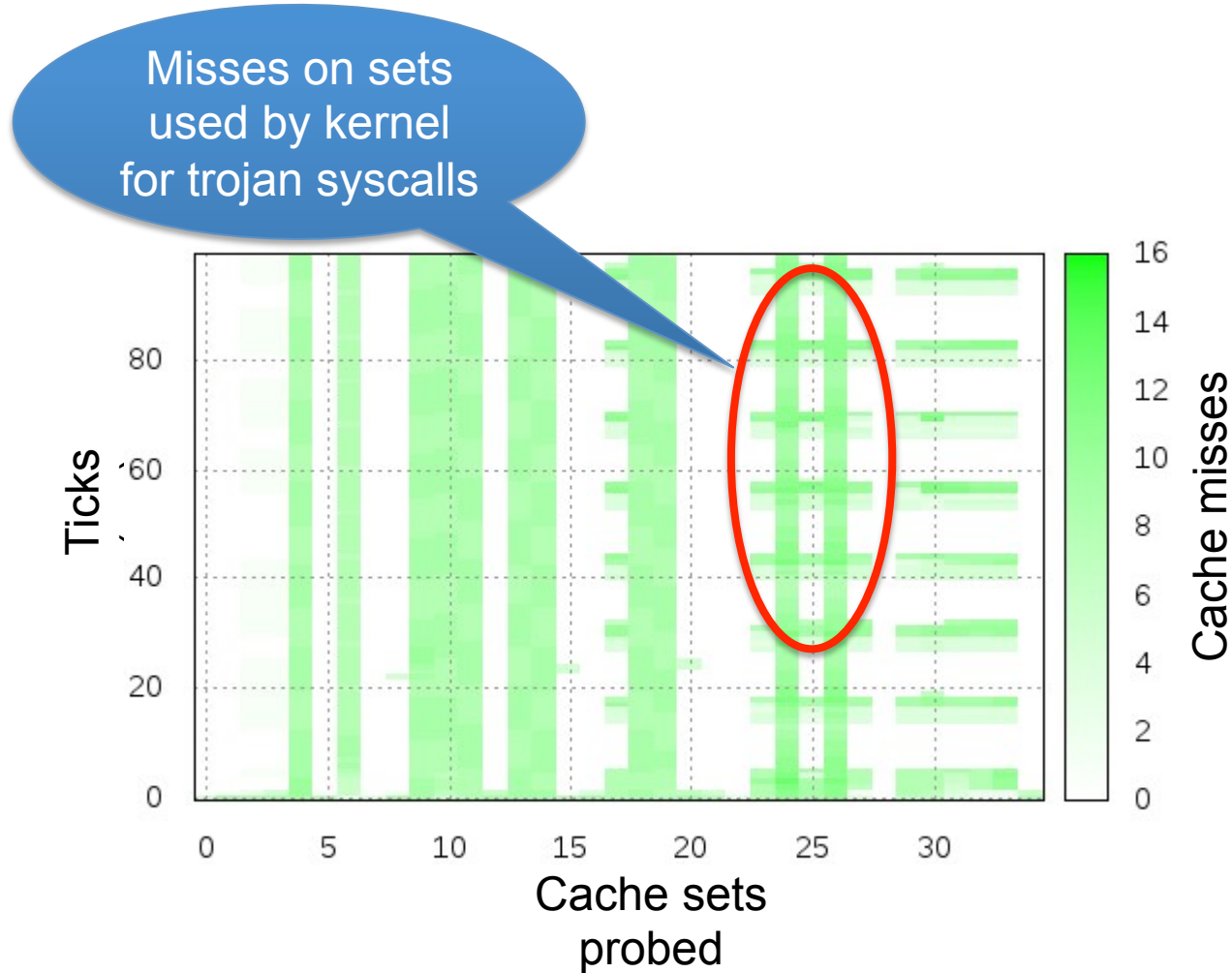
Resource Manager

Resource Manager

RM I+D

RM I+D

RAM

sel4 I+D

GRM I+D

Global Resource Manager

UNSW
AUSTRALIA

# Timing Channel Through Kernel

**High (Trojan)**

```
int count = 0;
for( ; ;) {

wait_for_new_system_tick( );
    if ((count %13) < 5)
        syscall(…);
    count++;
}
```

**Low (Spy)**

```
for( t = 0; t < 100 ; t++) {

wait_for_new_system_tick( );
    for (i = 0; i < prob_sets; i++)
        result[t][i] =
cache_probe(i)
}
```

**Kernel**

**Covert Channel**

Kernel image

UNSW
AUSTRALIA

# Cache Covert Channel Through Kernel
## Spy observations



Misses on sets used by kernel for trojan syscalls

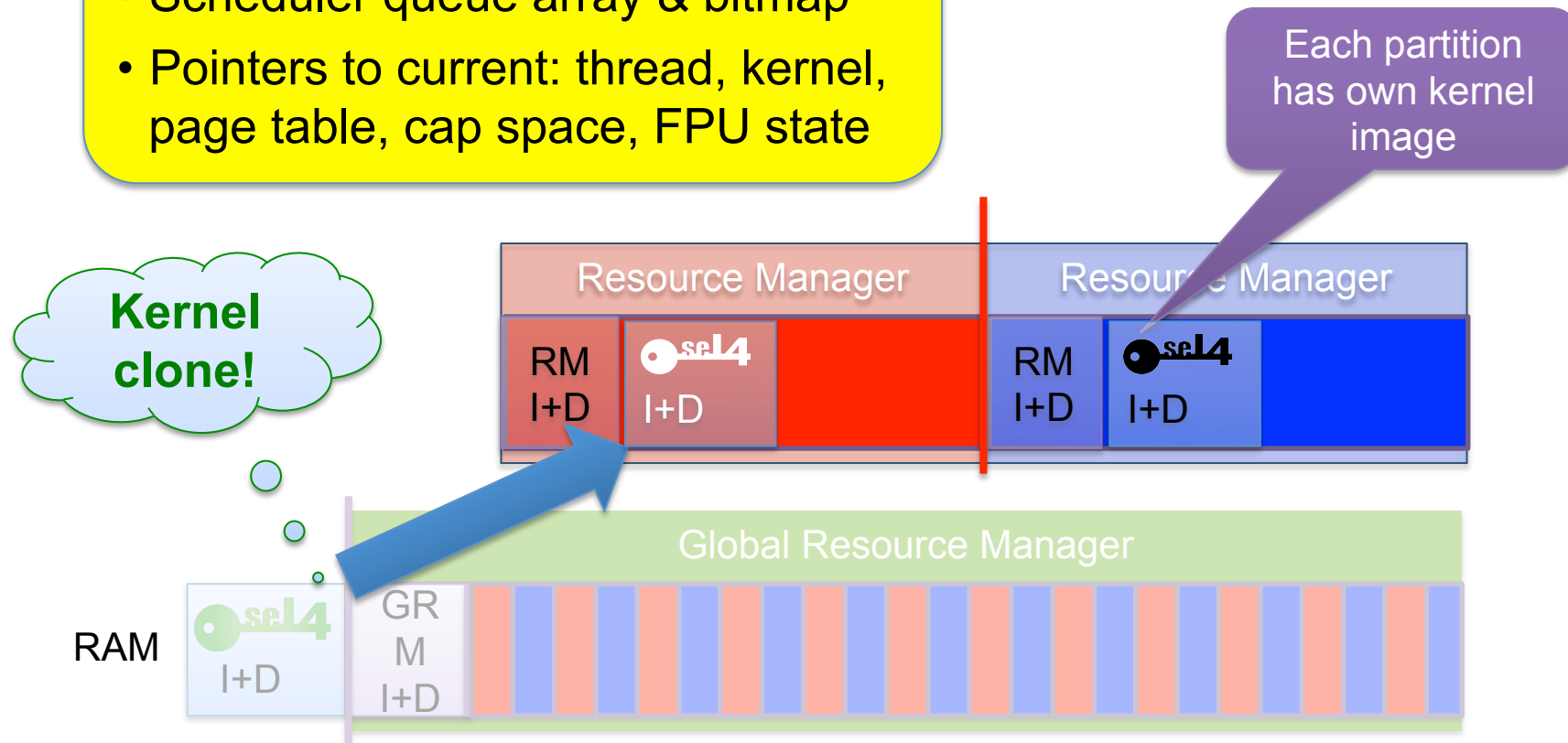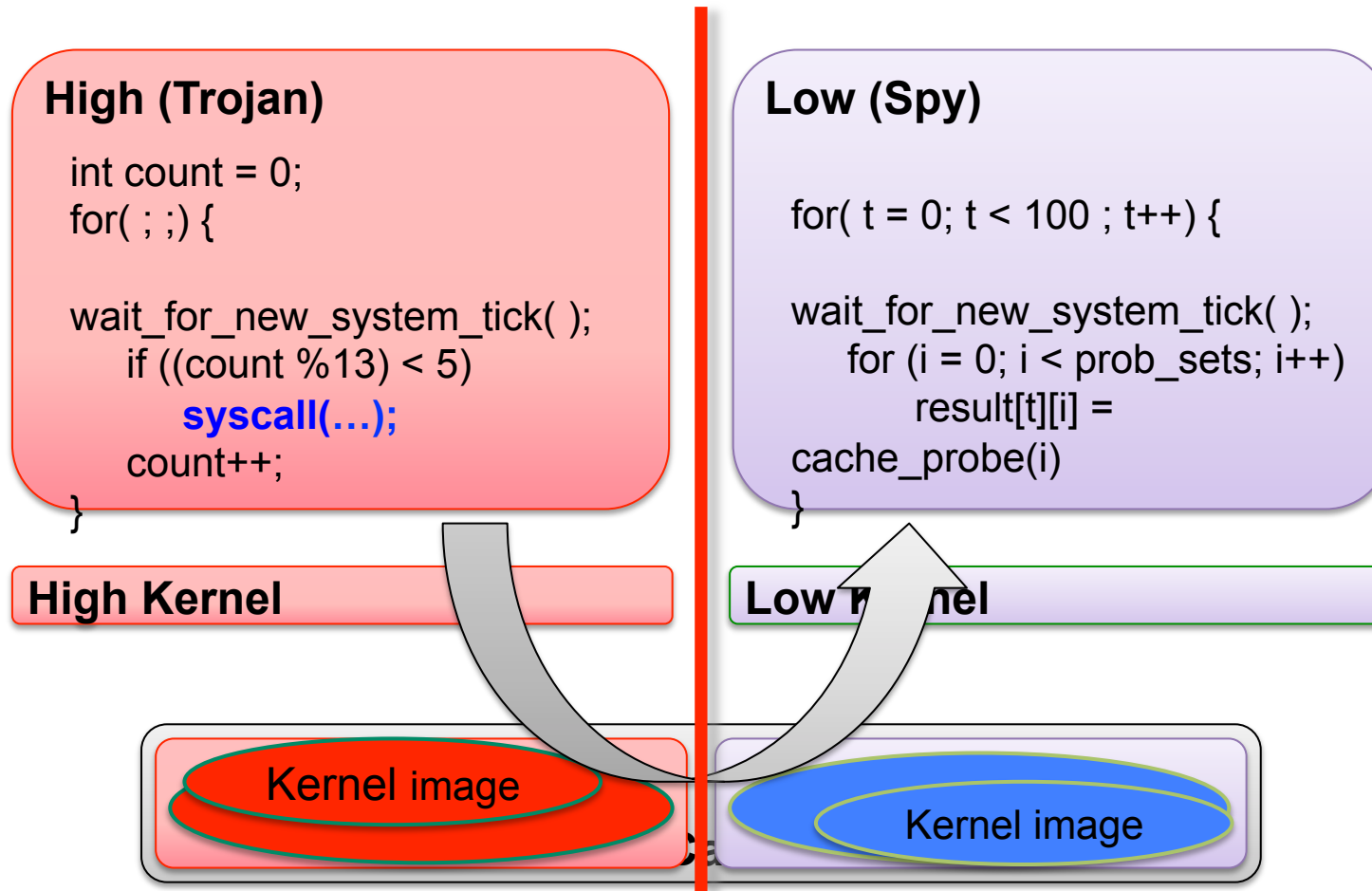© 2016 Gernot Heiser. Distributed under CC Attribution License

UNSW
AUSTRALIA

# Colouring the Kernel

**Only shared kernel data:**

- Scheduler queue array & bitmap
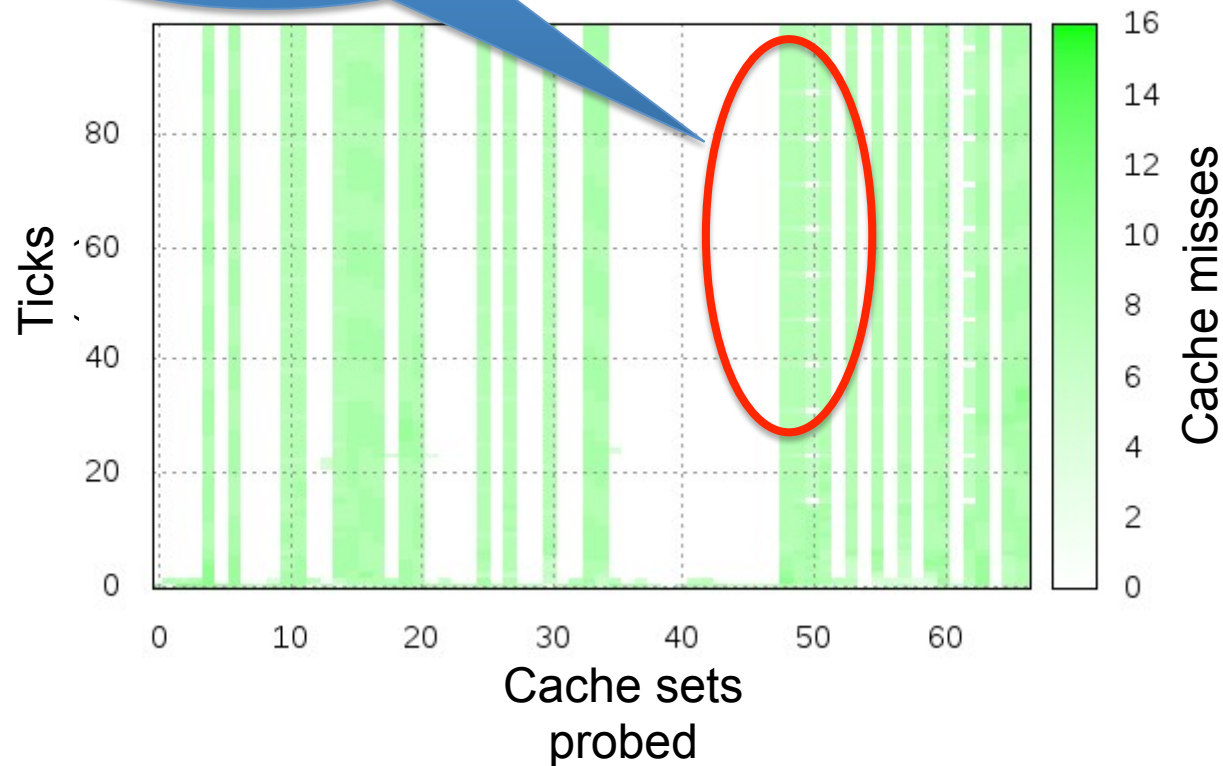- Pointers to current: thread, kernel, page table, cap space, FPU state

Each partition has own kernel image

Kernel clone!

Resource Manager

Resource Manager

RM I+D

I+D

RM I+D

I+D

Global Resource Manager

RAM

I+D

GR M I+D

UNSW
AUSTRALIA

# Timing Channel Through Kernel

**High (Trojan)**

```
int count = 0;
for( ; ;) {

wait_for_new_system_tick( );
    if ((count %13) < 5)
        syscall(…);
    count++;
}
```

**Low (Spy)**

```
for( t = 0; t < 100 ; t++) {

wait_for_new_system_tick( );
    for (i = 0; i < prob_sets; i++)
        result[t][i] =
cache_probe(i)
}
```

**High Kernel**

**Low Kernel**

Kernel image

Kernel image

© 2016 Gernot Heiser. Distributed under CC Attribution License

UNSW
AUSTRALIA

# Cache Covert Channel Through Kernel
## Spy observations with coloured kernel
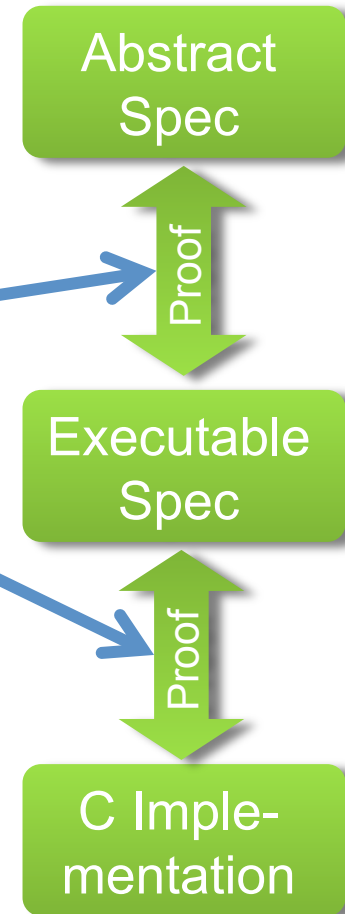


Only self-conflict missess,
no time signal!

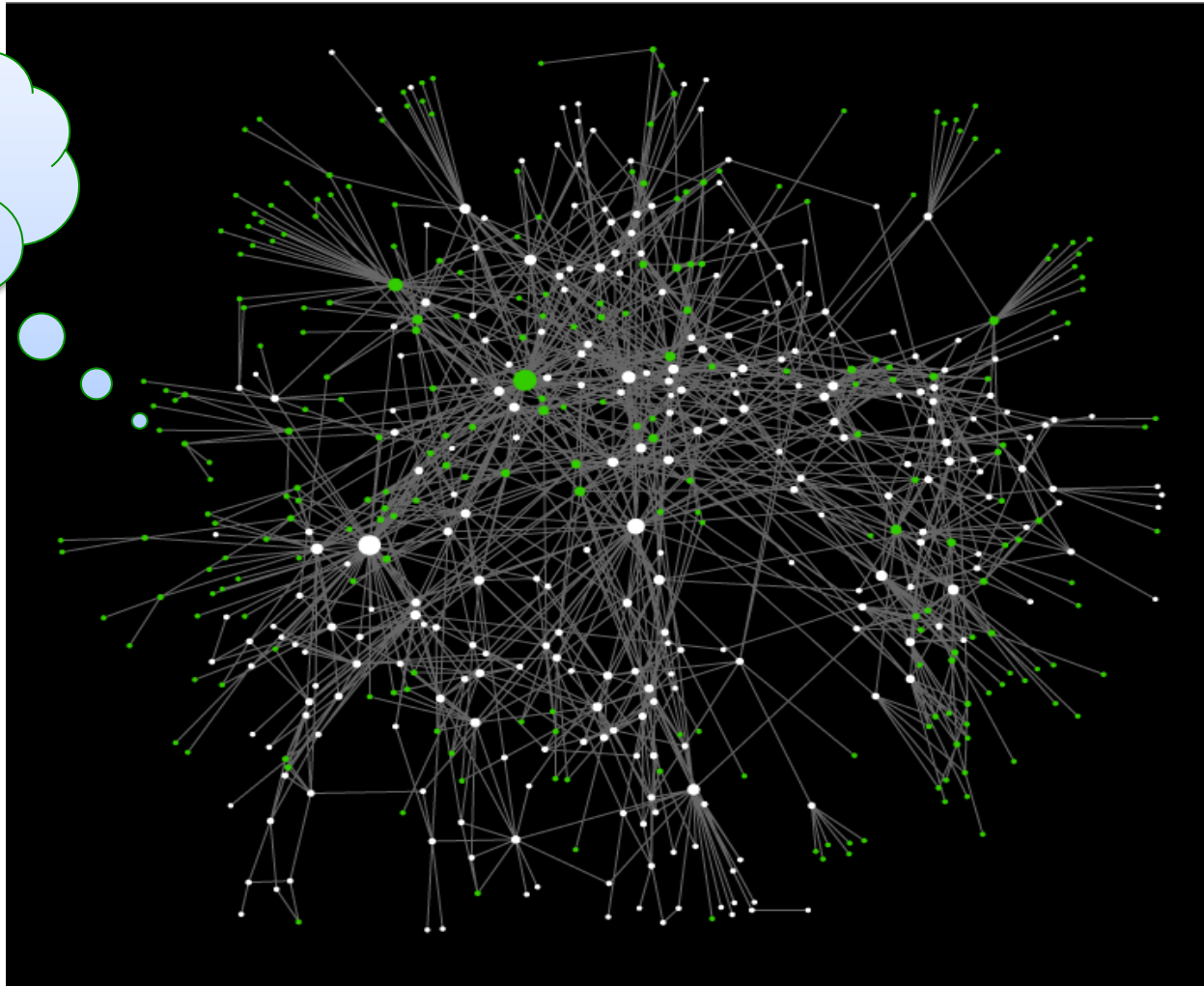# Tackling Verification Cost

# Verification Cost Breakdown

| | |
|---|---|
| Haskell design | 2 py |
| C implementation | 2 months |
| Debugging/Testing | 2 months |
| Abstract spec refinement | 8 py |
| Executable spec refinement | 3 py |
| Fastpath verification | 5 months |
| Formal frameworks | 9 py |
| **Total** | **24 py** |
| Repeat (estimated) | 6 py |
| Traditional engineering | 4–6 py |

Reusable!

Abstract Spec

Proof

Executable Spec

Proof

C Imple-mentation

UNSW
AUSTRALIA

# Why So Hard for 9,000 LOC?

seL4 call graph

# Cost of Assurance

Confiden-tiality

Availability

Integrity

Proof

Proof

Proof

4.5 py

Abstract Model

1 py
4 months

21 py
4.5 years

0 py
By construction

2 py, 1.5 years
Mostly for tools

C Imple-mentation

Proof

$400 per line of code!

2 py, 1 year
Mostly for tools

Binary code

Proof

Estimate repeat cost: $200/LOC

UNSW
AUSTRALIA

# Microkernel Life-Cycle Cost in Context



© 2016 Gernot Heiser. Distributed under CC Attribution License

# Cost of Assurance

**Industry Best Practice:**

- "High assurance": $1,000/LOC, no guarantees, *unoptimised*
- Low assurance: $100–200/LOC, 1–5 faults/kLOC, *optimised*

**State of the Art – seL4:**

- – $400/LOC, 0 faults/kLOC, *optimised*
- Estimate repeat would cost half
  - – that's about twice the development cost of the predecessor Pistachio!
- Aggressive optimisation [APSys'12]
  - – much faster than traditional high-assurance kernels
  - – as fast as best-performing low-assurance kernels

# What Have We Learnt?

**Formal verification *probably* didn't produce a more *secure* kernel**

- In reality, traditional separation kernels are *probably* secure

**But:**

- We now have certainty
- We did it *probably* at less cost

**Real achievement:**

- Cost-competitive at a scale where traditional approaches still work
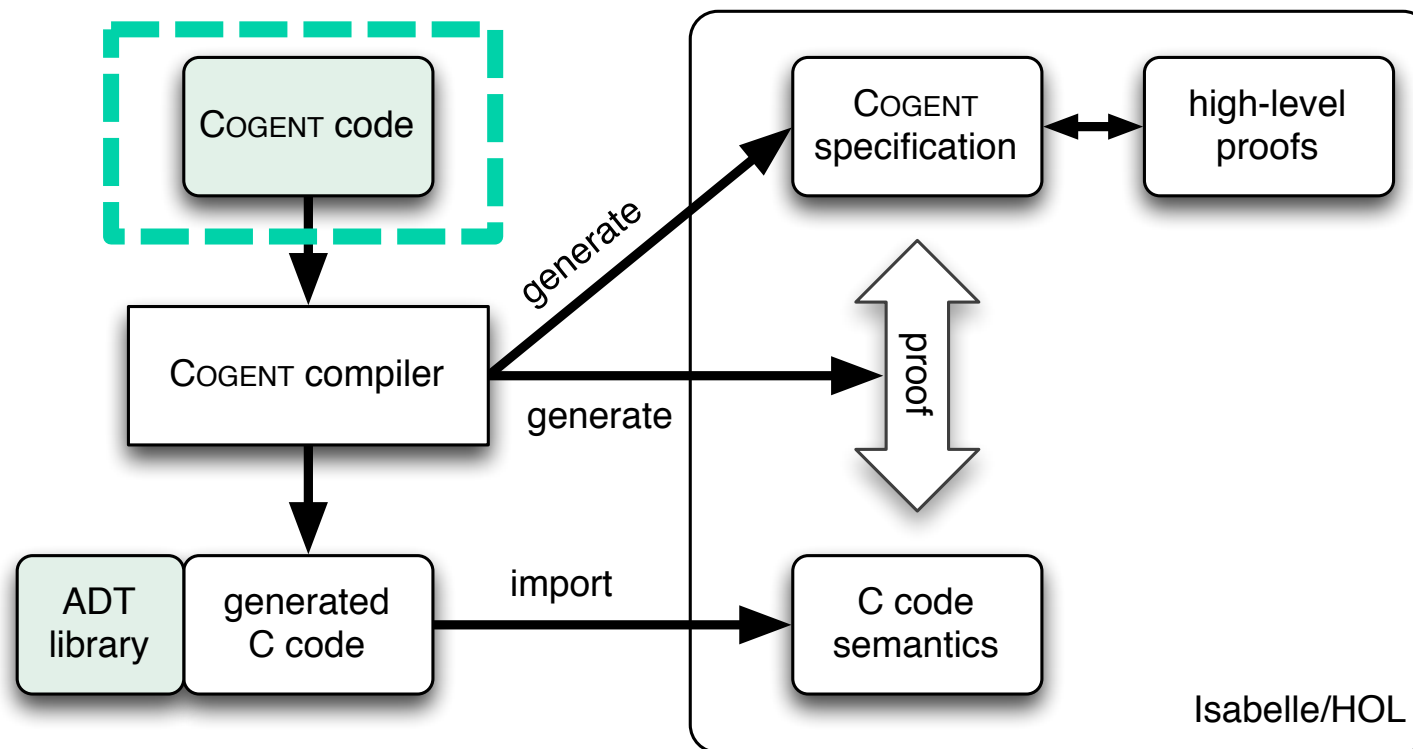- Foundation for scaling beyond: **2 × cheaper, 10 × bigger!**

**How?**

- Combine theorem proving with
  - synthesis
  - domain–specific languages (DSLs)

UNSW
AUSTRALIA

# Our approach

- Cogent: code and proof co-generation
  - Implement FS in high-level functional language (and reason about it)
  - Generate efficient low-level code in C
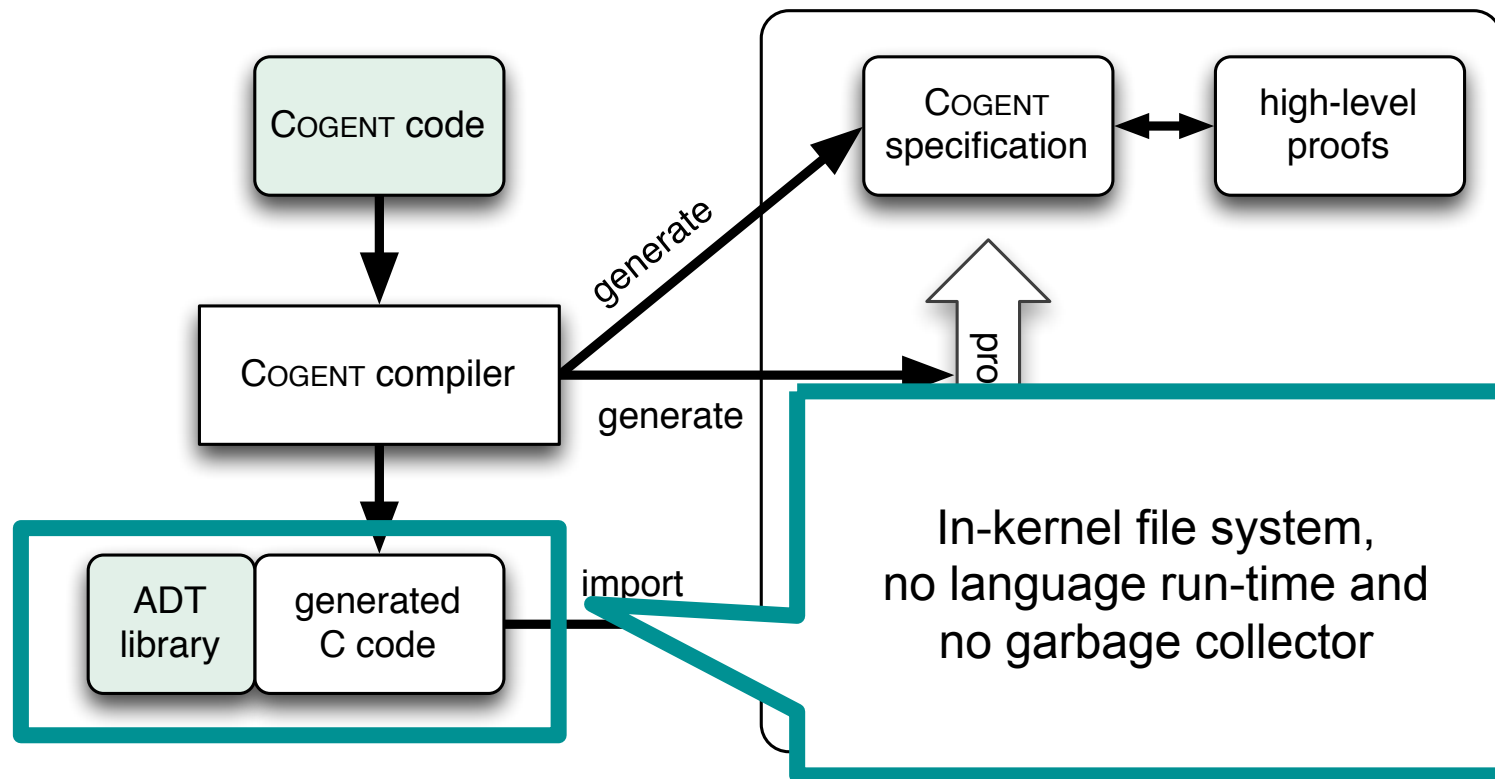  - Automatically prove correspondence between the two

# Cogent Workflow

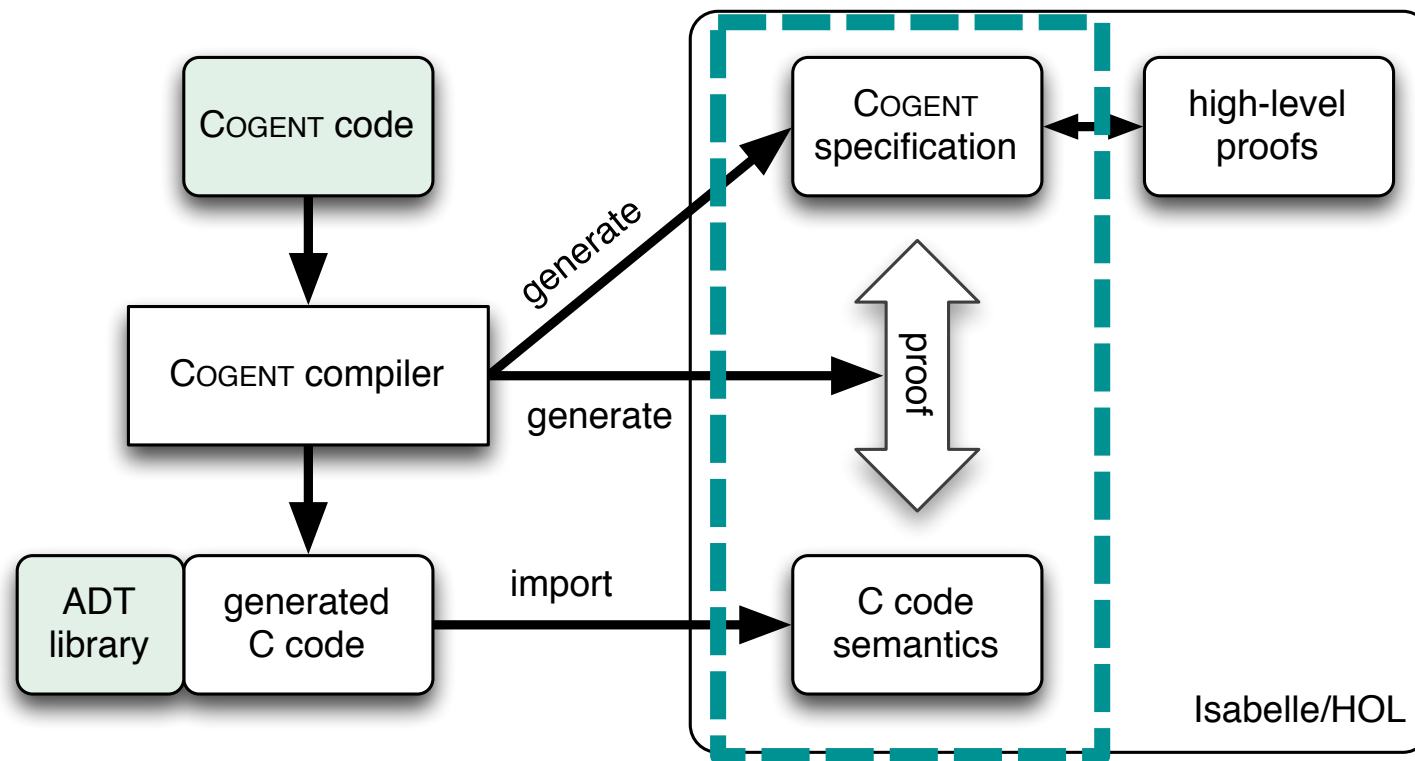- Cogent: purely functional memory-safe language



© 2016 Gernot Heiser. Distributed under CC Attribution License

# Cogent workflow

- Cogent's certifying compiler generates an C implementation



COGENT code

COGENT compiler

COGENT specification ↔ high-level proofs

generate

generate

proof

ADT library | generated C code

import

In-kernel file system,
no language run-time and
no garbage collector

UNSW
AUSTRALIA

# Cogent workflow

- Cogent generates a specification and a proof that links it to the C code

© 2016 Gernot Heiser. Distributed under CC Attribution License
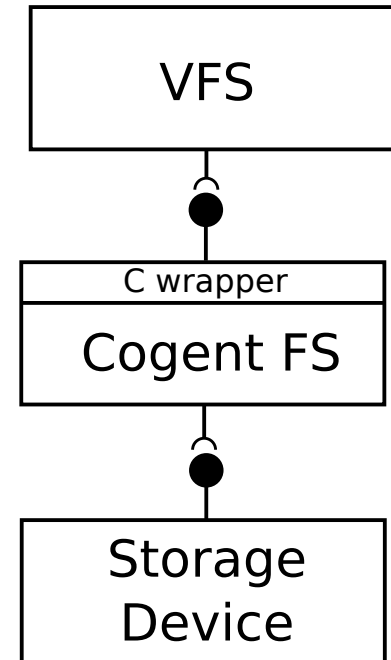
UNSW
AUSTRALIA

# Cogent workflow

- Prove high-level properties about Cogent-generated specifications using a proof assistant

# Cogent File Systems

- We implemented two Linux FSs:
  - Ext2: functionally complete original spec
    - No ACLs, symlinks
  - BilbyFs: custom flash file system

- Invoked from VFS via a small C wrapper, which:
  - Uses a global lock to prevent concurrent execution of FS operations
  - Handles VFS caches
  - Calls Cogent FS entry points

- FSs interface with the storage device via external ADT functions

VFS

C wrapper

Cogent FS

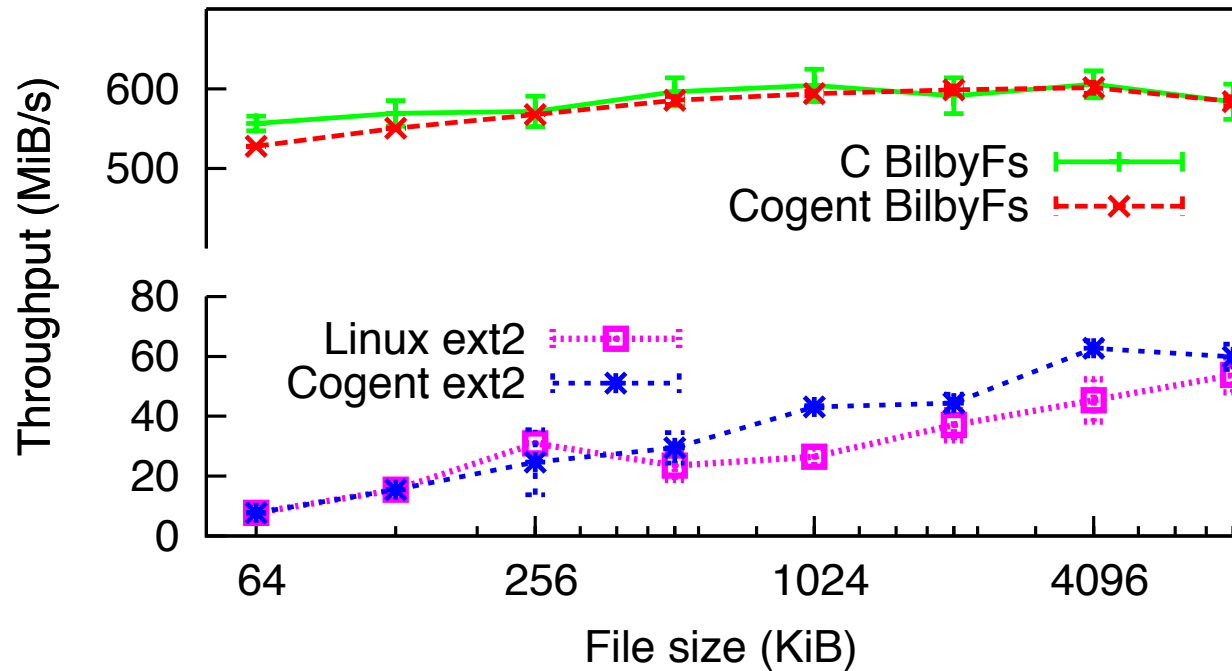Storage
Device

UNSW
A U S T R A L I A

# Evaluation

- Compare ext2 with Linux's native implementation
  - Hardware:
    - 4 core i7-6700 running at 3.1 GHz,
    - Samsung HD501JL 7200RPM 500G SATA disk

- Compare BilbyFs with handwritten C implementation
  - Hardware:
    - Mirabox development board
    - Marvell Armada 370 single-core 1.2 GHz ARMv7 processor
    - 1 GiB of NAND flash

UNSW
AUSTRALIA

# IOZone random 4k writes



- 20% CPU load for Cogent BilbyFs vs 15% for C
- Both ext2 implementations have the same CPU load

UNSW
AUSTRALIA

# Postmark on RAM-disk

| System | Total time sec | creation files/sec | read rate kB/sec |
|---|---|---|---|
| C ext2 | 10 | 5025 | 248 |
| COGENT ext2 | 21 | 2393 | 118 |
| C BilbyFs | 6 | 33375 | 431 |
| COGENT BilbyFs | 10 | 20025 | 259 |

- Degradation of a factor 2 for Cogent FSs

UNSW
AUSTRALIA

# Postmark on RAM-disk

| System | Total time<br>sec | creation<br>files/sec | read rate<br>kB/sec |
|---|---|---|---|
| C ext2 | 10 | 5025 | 248 |
| COGENT ext2 | 21 | 2393 | 118 |
| C BilbyFs | 6 | 33375 | 431 |
| COGENT BilbyFs | 10 | 20025 | 259 |

- Degradation of a factor 2 for Cogent FSs

- Overhead is due to two reasons:
  - extra copying involved when converting in-buffer directory entries into Cogent's internal data type
  - Cogent compiler is overly reliant on C compiler's optimiser to convert automatically C structs passed by copy to pointers

UNSW
AUSTRALIA

# ⊙seL4 Remember: Verification Cost Breakdown



**Abstract Spec**

**8 py** — Proof

**Executable Spec**

Cogent spec higher level than seL4 exec spec

**3 py** — Proof

Fully automated in Cogent

**C Implementation**

UNSW AUSTRALIA