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

---

# 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
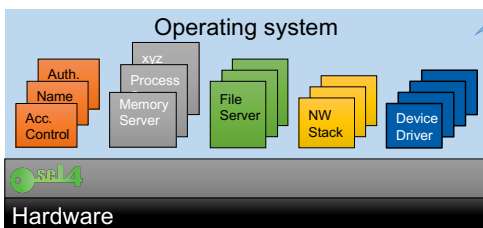> • no defence available
> • confirmed

> 115 critical
> Linux CVEs
> to Nov'17

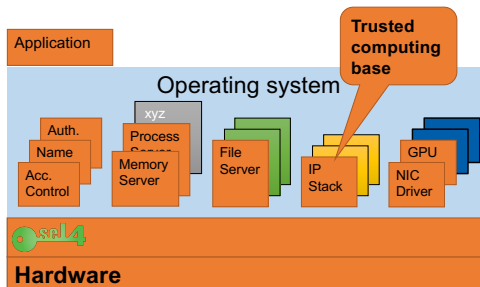• For each establish how microkernel-based design would change impact

---

## seL4 Hypothetical seL4-based OS

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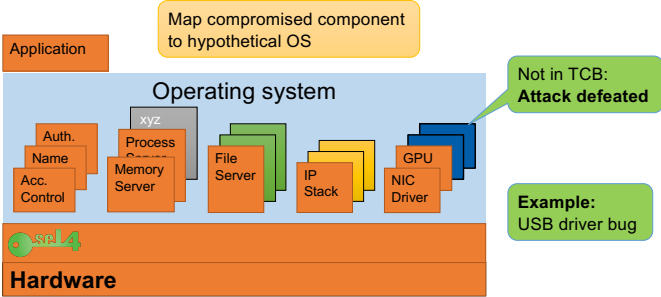
> Functionality comparable to Linux

**Operating system**

Auth.
Name
Acc. Control

xyz
Process Server
Memory Server

File Server

NW Stack

Device Driver

seL4

**Hardware**

---

## seL4 Hypothetical Security-Critical App

Application

**Trusted computing base**

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

**Operating system**

Auth.
Name
Acc. Control

xyz
Process Server
Memory Server

File Server

IP Stack

GPU

NIC Driver

seL4

**Hardware**

## Slide 6

seL4 **Analysing CVEs**

Application

Operating system

Map compromised component to hypothetical OS

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

seL4 **Hardware**

## Slide 7

seL4 **Analysing CVEs**

Application

Operating system

Map compromised component to hypothetical OS

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

In microkernel: **Attack defeated by verifiation**

Auth. Name Acc. Control | xyz Process Server Memory Server | File Server | IP Stack | GPU NIC Driver

seL4 **Hardware**

## Slide 8

seL4 **Analysing CVEs**

Application

Operating system

Map compromised component to hypothetical OS

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

Auth. Name Acc. Control | xyz Process Server Memory Server | File Server | IP Stack | GPU NIC Driver

Example: File system compromised

seL4 **Hardware**

## Slide 9

seL4 **Analysing CVEs**

Application

Operating system

Map compromised component to hypothetical OS

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

Auth. Name Acc. Control | xyz Process Server Memory Server | File Server | IP Stack | GPU NIC Driver

**Example:** GPU compromised

seL4 **Hardware**

## Slide 10

seL4 **Analysing CVEs**

Application

Operating system

Map compromised component to hypothetical OS

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

Auth. Name Acc. Control | xyz Process Server Memory Server | File Server | IP Stack | GPU NIC Driver

Full system compromise: **No effect**

seL4 **Hardware**

## Slide 11

seL4 **All Critical Linux CVEs to 2017**

Still full system compromise: **No effect** — 4%

Not in TCB: **Attack defeated** — 30%

No full compromise, but violates integrity or confidentiality: **Weakly mitigated** — 38%

Only *crash* essential service (availability): **Strongly mitigated** — 17%

In microkernel: **Attack defeated by verification** — 11%

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

## Summary

**OS structure matters!**

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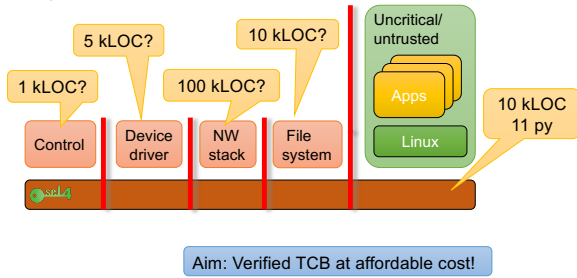- Monolithic OS design is *fundamentally flawed from security point of view*

[Biggs et al., APSys'18]

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

---

# Cogent

---

## Beyond the Kernel



- 1 kLOC? — Control
- 5 kLOC? — Device driver
- 100 kLOC? — NW stack
- 10 kLOC? — File system
- Uncritical/untrusted: Apps, Linux
- 10 kLOC 11 py
- seL4

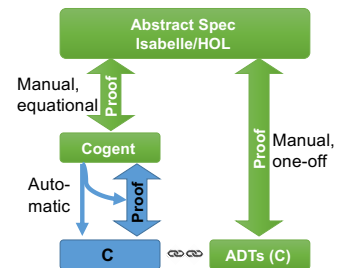Aim: Verified TCB at affordable cost!

---

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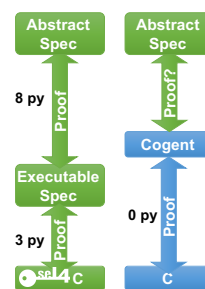
[O'Connor et al, ICFP'16; Amani et al, ASPLOS'16]



Abstract Spec Isabelle/HOL

Manual, equational — Proof

Cogent

Auto-matic — Proof

C ⊂⊃⊂⊃ ADTs (C)

Manual, one-off — Proof

---

## 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

---

## Addressing Verification Cost



Abstract Spec — Proof — Executable Spec — Proof — seL4 C
8 py
3 py

Abstract Spec — Proof? — Cogent — Proof — C
0 py

**Dependability-cost tradeoff:**
- Reduced faults through safe language
- Property-based testing (QuickCheck)
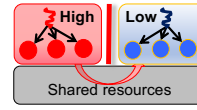- Model checking
- Full functional correctness proof

**Spec reuse!**

**Work in progress:**
- Language expressiveness
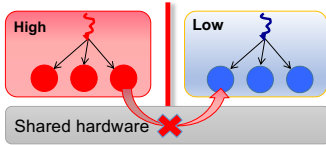- Reduce boiler-plate code
- Network stacks
- Device drivers

# Time Protection

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

---

## Refresh: Microarchitectural Timing Channels

**High** **Low**

Shared resources

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

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

---

## OS Must Enforce *Time Protection*

**High** **Low**

Shared hardware
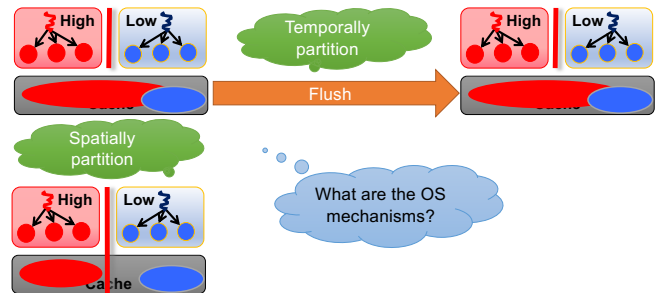
**Preventing interference is core duty of the OS!**
- *Memory protection* is well established
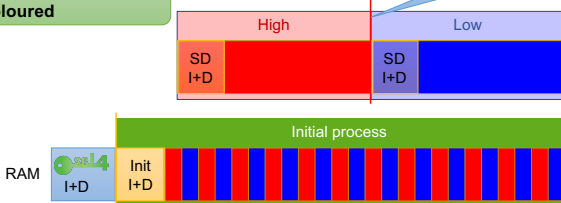- *Time protection* is completely absent

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

---

## Time Protection: No Sharing of HW State

**High** **Low** — Temporally partition

Cache — Flush — Cache

Spatially partition

**High** **Low**

Cache

What are the OS mechanisms?

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

---

## seL4 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

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

---

## seL4 Spatial Partitioning: Cache Colouring

**High** **Low**

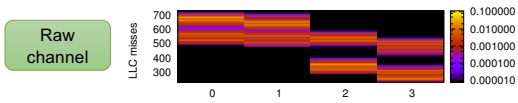TCB PT — TCB PT

seL4

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

Shared kernel image

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

## Channel Through Kernel Code

Raw channel



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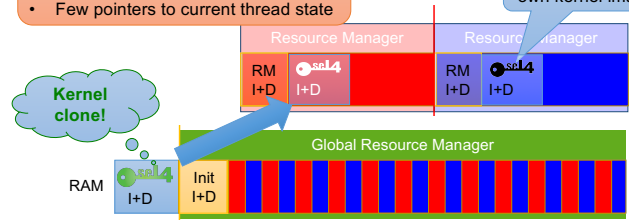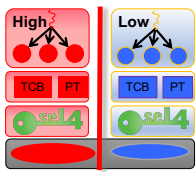
---

## Colouring the Kernel

Ensure deterministic access!

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

Each partition has own kernel image



Kernel clone!

RAM

---

## 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
- Per-partition kernel image to colour kernel

Ensure deterministic access!

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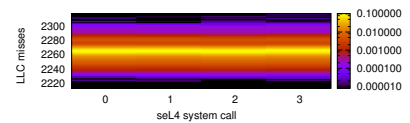- Scheduler queue array & bitmap
- Few pointers to current thread state

---

## Channel Through Kernel Code

Raw channel

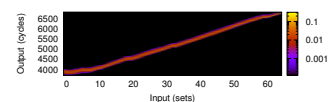Channel with cloned kernel

---

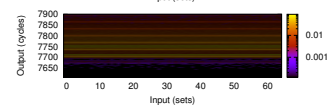## Temporal Partitioning: Flush on Switch

Must remove any history dependence!

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

6. Reprogram timer
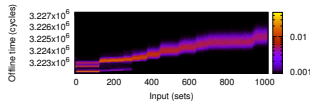7. return

---

## D-Cache Channel

Raw channel

Channel with flushing

## Flush-Time Channel

Raw channel



Offline time (cycles)
3.227x10^6
3.226x10^6
3.225x10^6
3.224x10^6
3.224x10^6

Input (sets)

---

## 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()) ;
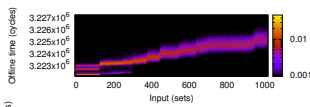6. Reprogram timer
7. return

Latency depends on prior execution!

Time padding to remove dependency

Ensure deterministic execution

---

## Flush-Time Channel

Raw channel

Channel with deterministic flushing

---

## Performance Impact of Colouring

Splash-2 benchmarks on Arm A9



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 |

---

## 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
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

Cannot share HW threads across security domains!

[Ge et al., APSys'18]

---

## 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!

## 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 $\geq T_0 + WCET$

[Heiser et al., HotOS'19]

Functional property

---

# Making COTS Hardware Dependable

---

## Satellites: SWaP vs Dependability

Space is becoming commodisized:
- many, small (micro-) satellites
- increasing cost pressure

Harsh evironment for electronics:
- temperature fluctuations
- ionising radiation

Radiation-hardened processors are slow, bulky and expensive

Use redundancy of cheap COTS multicores

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

---

## Traditional Redundancy Approaches

Sphere of replication

Master | Slave | Slave
App / Lib / OS

Watch-dog | Syscall Emulation Layer | OS | CPU(s) | Devices

HW lockstepping/voting infrastructure
CPU | Dev | CPU | Dev | CPU | Dev

Fault-tolerant HW:
- expensive

SW replication:
- cheap
- incomplete

---

## Redundant Co-Execution (RCoE)

Userland transparently replicated

Sphere of replication

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

Device interface

Primary | Secondary | Secondary
App / Lib / Driver

Device | Core | Core | Core

Vote & sync | Vote & sync

No master-slave, but peer-to-peer

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

---

## 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

## Slide 42

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



Vote → Leading Replica? → Wait on Barrier → Synced

Set Breakpoint on IP and catch up → Breakpoint Exception → Compare <branches, IP> → Clear Breakpoint

Next / No / Yes

---

## Slide 43

|  | 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

---

## Slide 44

| 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-core VM

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×

---

## Slide 45

Yahoo! Cloud Service Benchmark

In-memory key-value store



Load generator: YCSB / OS / NIC

System under test:
Redis / lwIP / NIC driver — seL4 — Core (×3)

---

## Slide 46

1000 transactions/s

LC: loosely-coupled
CC: closely-coupled
D: DMR
T: TMR
A: vote on interrupt
S: also vote on syscall

Base / LC-D-A / CC-D-A / LC-D-S / CC-D-S / LC-T-A / CC-T-A / LC-T-S / CC-T-S

Overhead is 1.2–3 depending on configuration

---

## Slide 47

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 |

## Comparison to Rad-Hardened Processor

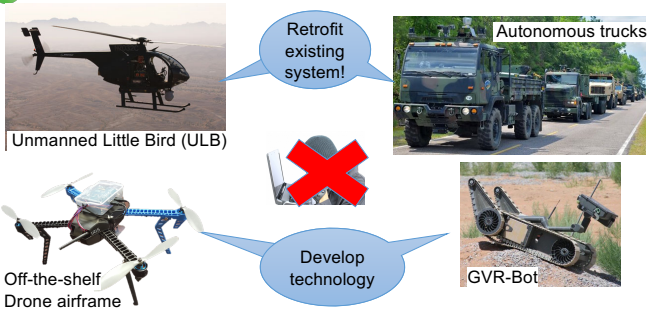| | Sabre Lite | RAD750 |
|---|---|---|
| Cores @ clock | 4 @ 800 MHz | 1 @ 133 MHz |
| Performance | 4 × 2,000 DMIPS | 240 DMIPS |
| Power | < 5 W | < 6 W |
| Energy Efficiency | 200 DMIPS/W | 40 DMIPS/W |
| Cost | $200 | $200,000 |
| Perf/Cost | 5 DMIPS/$ | 0.0002 DMIPS/$ |

2002 price

Assuming 2× overhead, TMR

[Shen et al., DSN'19]

---

# Real-World Use

---

## DARPA HACMS

Retrofit existing system!

Autonomous trucks

Unmanned Little Bird (ULB)

Off-the-shelf Drone airframe

Develop technology

GVR-Bot

---

## ULB Architecture

Ground Station Link — Mission Computer — GPS, Camera

Network

Sensors — Flight Computer — Motors

---

## 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
Linux
Local NW | VMM
seL4
Camera
Linux
VMM

---

## Incremental Cyber Retrofit

Original Mission Computer

**Trusted**
Mission Manager
Crypto | Cam
Local NW | G
Ground Stn Lir
Linux

**Trusted**
GS Lk
Miss Mgr
Crypto
GPS
Linux
Local NW | VMM
seL4
Camera
Linux
VMM

**Trusted**
Crypto | Mission Mngr
Local NW | Comms
seL4
Camera
Linux
GPS | VMM

## Incremental Cyber Retrofit

Original Mission Computer

[Klein et al, CACM, Oct'18]

Cyber-secure Mission Computer

**Trusted**
- Mission Manager
- Crypto | Camera
- Local NW | GPS
- Ground Stn Link
- Linux

**Trusted**
- Crypto | Mission Mngr
- Local NW | Comms

Camera | Linux | GPS | VMM

seL4

## Issue: Capabilities are Low-Level



>50 capabilities for trivial program!

## Simple But Non-Trivial System

## Component Middleware: CAmkES

Higher-level abstractions of low-level seL4 constructs

Interface

Comp A — RPC — Comp B

Component | Connector

Shared memory | Comp C | Semaphore

## HACMS UAV Architecture

Security enforcement: Linux only sees encrypted data

Radio Driver | Data Link

Crypto

Uncritical/untrusted, contained
- Wifi
- Camera
- Linux

CAN Driver

seL4

## Enforcing the Architecture

Radio Driver | Data Link | Uncritical/untrusted, contained | Wifi | Camera | Linux

CAN Driver | Crypto

glue.c | driver.c | VMM.c

**Conditions apply**

Compiler/Linker

init.c

binary

## Architecture Analysis

```
Eclipse-based IDE  --Design-->  AADL  -->  Analysis Tools   Safety ✔
                                   |
                                   | Generate      Architecture analysis
                                   v                and design language
                                 CAmkES  --Generate-->  .h, .c
                                                          |
                                                          | Compile
                                                          v
                                                        Binary
```

## Military-Grade Security

Cross-Domain Desktop Compositor

Multi-level secure terminal
- Successful defence trial in AU
- Evaluated in US, UK, CA
- Formal security evaluation soon

Pen10.com.au crypto communication device in use in AU, UK defence

## Real-World Use
### Courtesy Boeing, DARPA