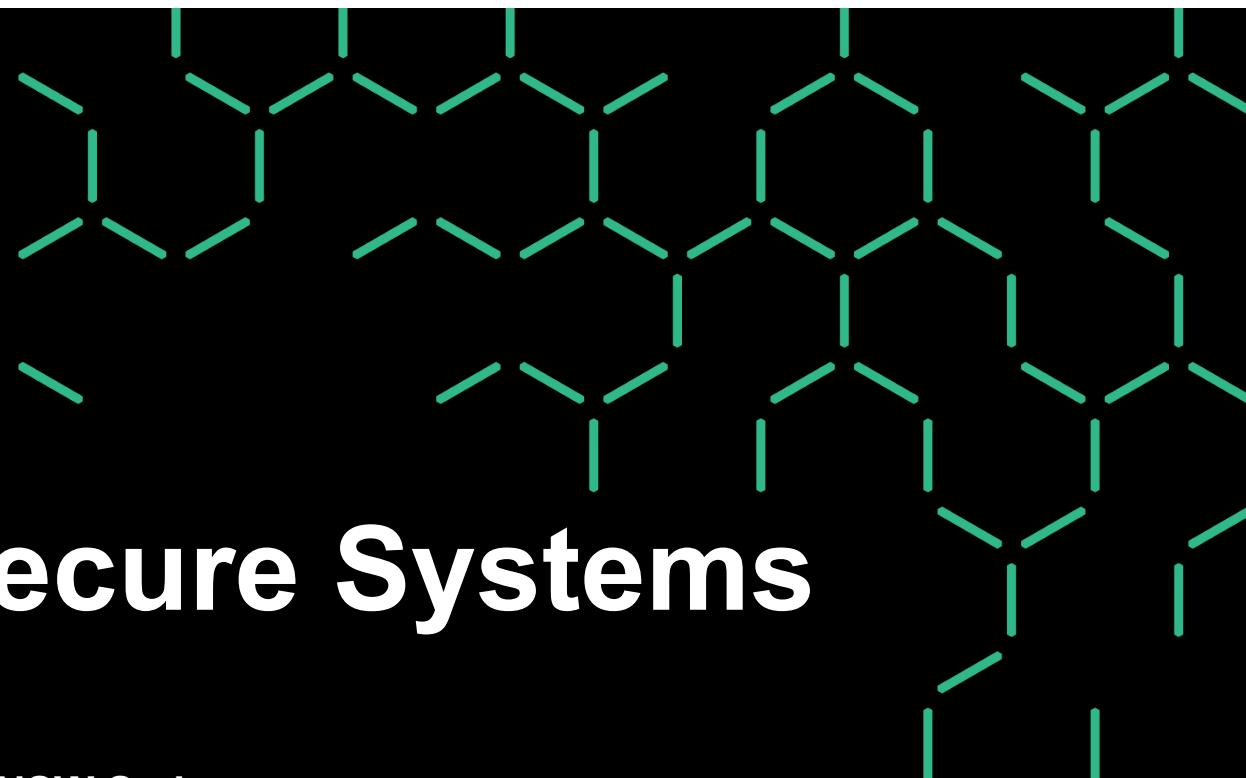




DATA  
61



# Building Secure Systems With seL4

Gernot Heiser | Data61 and UNSW Sydney

[Gernot.Heiser@data61.csiro.au](mailto:Gernot.Heiser@data61.csiro.au) | [@GernotHeiser](https://twitter.com/GernotHeiser)

FMATS, Cambridge, Sep'18

<https://sel4.systems>





# What is seL4?



seL4: The world's **only**  
operating-system kernel with  
**provable** security  
enforcement

seL4: The world's  
**only** protected-mode OS  
with complete, sound  
timeliness analysis

**Open Source**

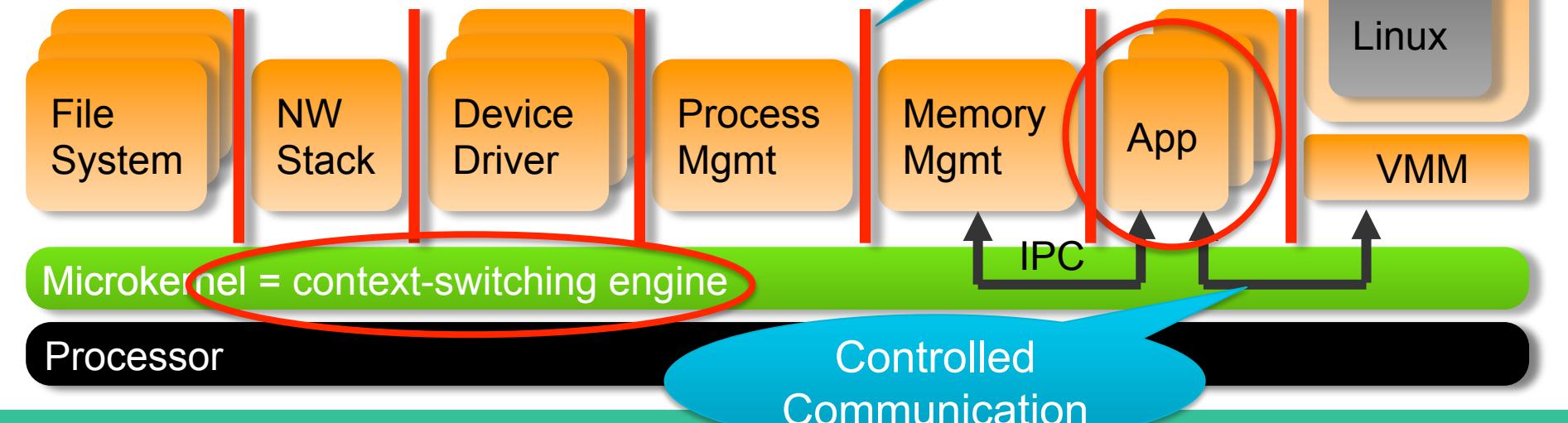
seL4: The world's  
**fastest** microkernel



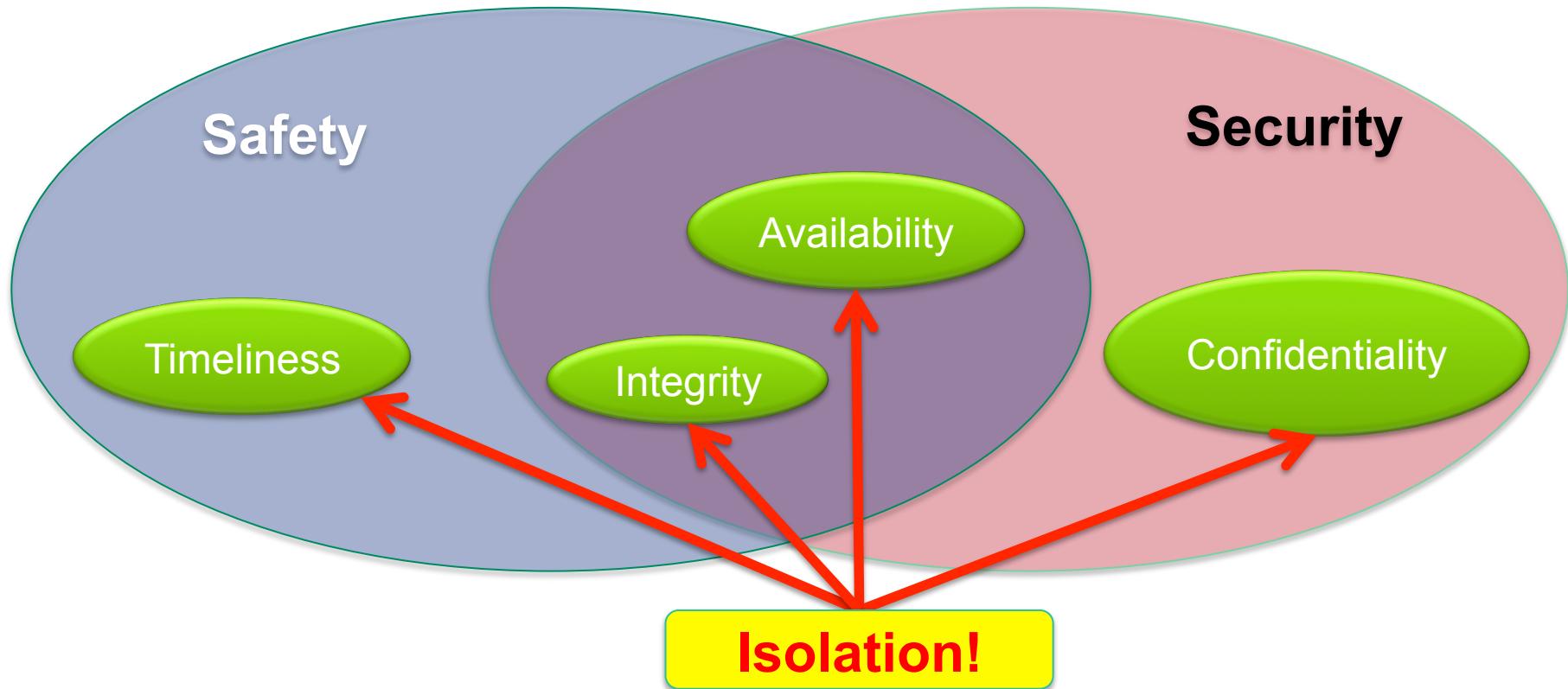
# A Microkernel is not an OS



Device drivers, file systems, crypto, power management, virtual-machine monitor are all usermode processes



# seL4 Security and Safety





# seL4 Provable Security Enforcement



Confidentiality

Integrity

Availability

World's fastest  
microkernel!

Abstract  
Model

Functional  
correctness

Translation  
correctness

Worst-case  
execution time

C Imple-  
mentation

Binary code

Proof

Proof

Proof

Proof

Exclusions (in progress):

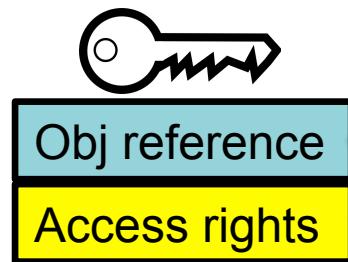
- Initialisation
- Privileged state & caches
- Security proofs on x86?
- Any proofs on RISC-V
- Multicore
- Covert *timing* channels



# Core Mechanism: Capability



**Capability = Access Token:**  
Prima-facie evidence of privilege



Eg. thread,  
address  
space

Eg. read,  
write, send,  
execute...

Any system call is invoking a capability:  
`err = method( cap, args );`

Capabilities provide:

- Fine-grained access control
- Reasoning about information flow

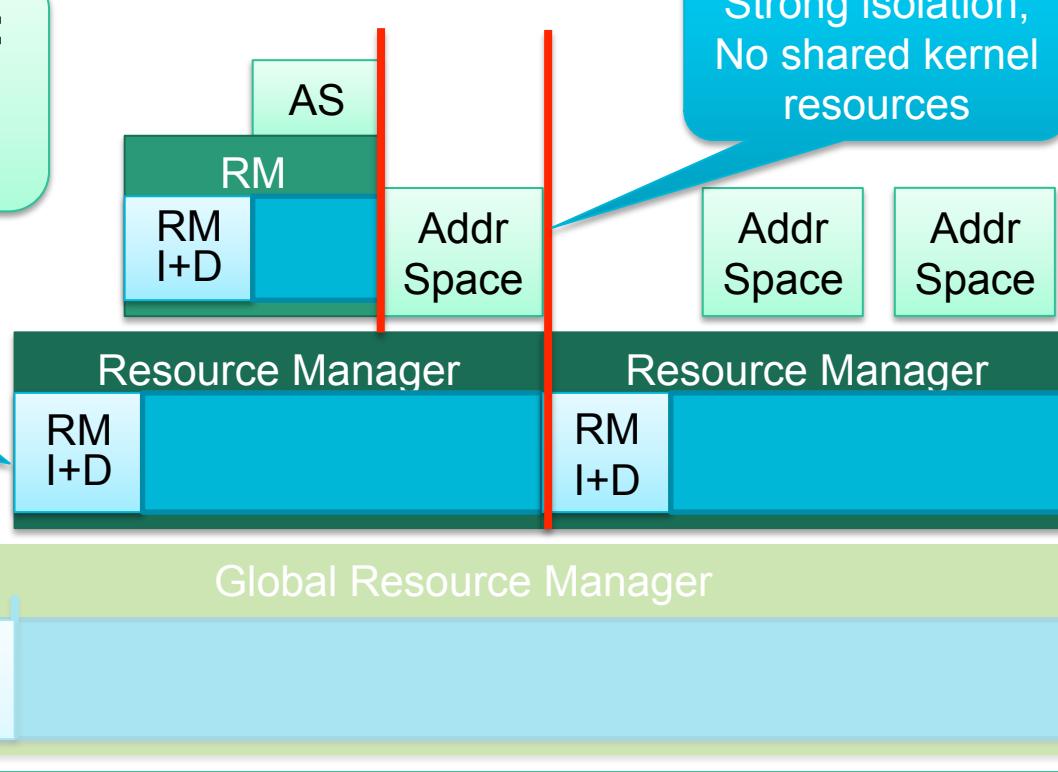


# seL4 Difference To Other Kernels

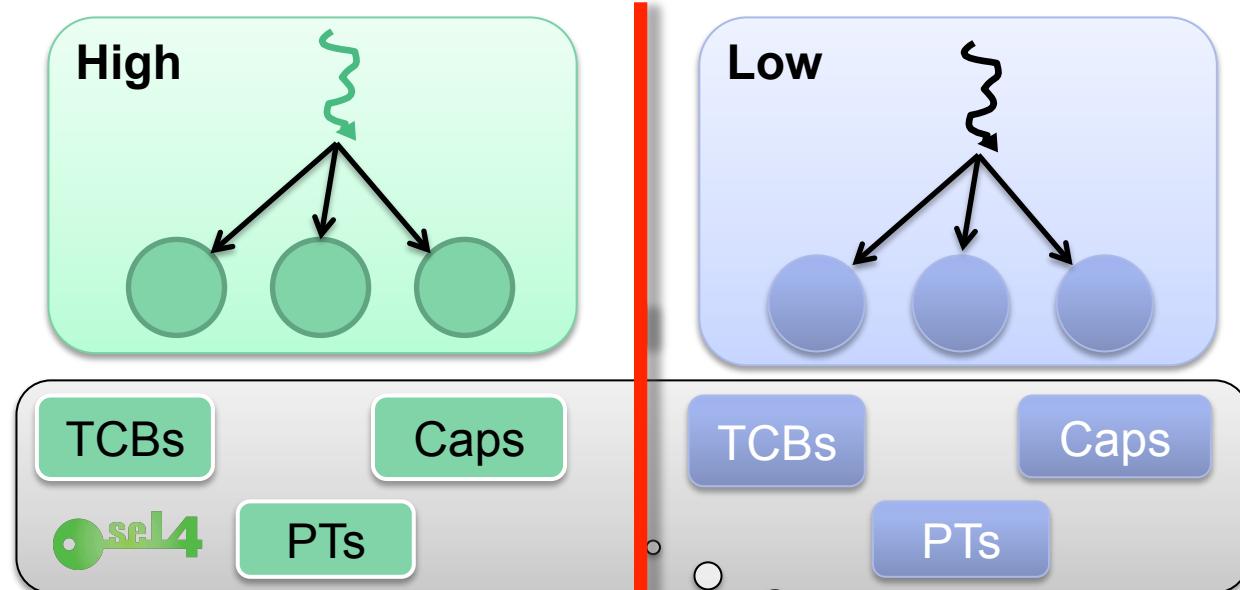


**Design for isolation:  
no memory  
allocation by kernel**

Resources fully  
delegated, allows  
autonomous  
operation



# •sel4 Isolation Goes Deep



Kernel data  
partitioned  
like user data

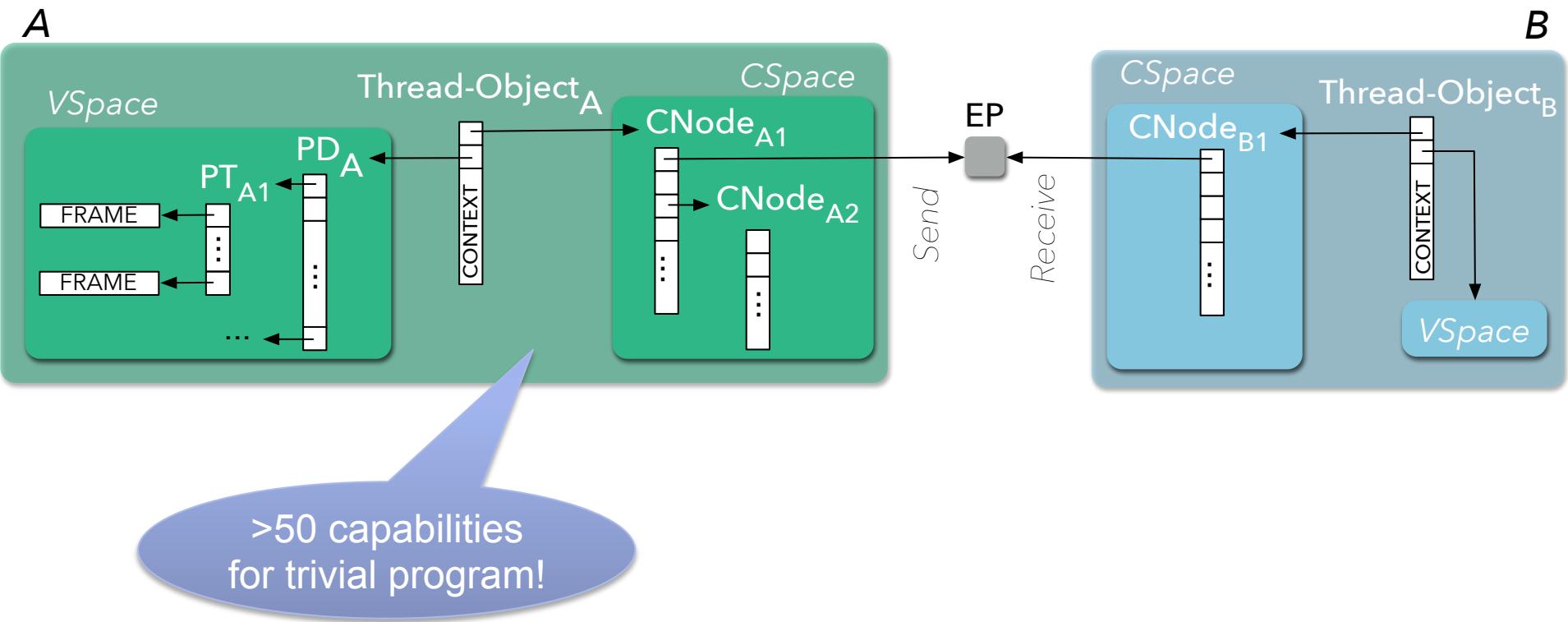


# CAmkES: Security by Architecture

sel4



# Issue: Capabilities are Low-Level

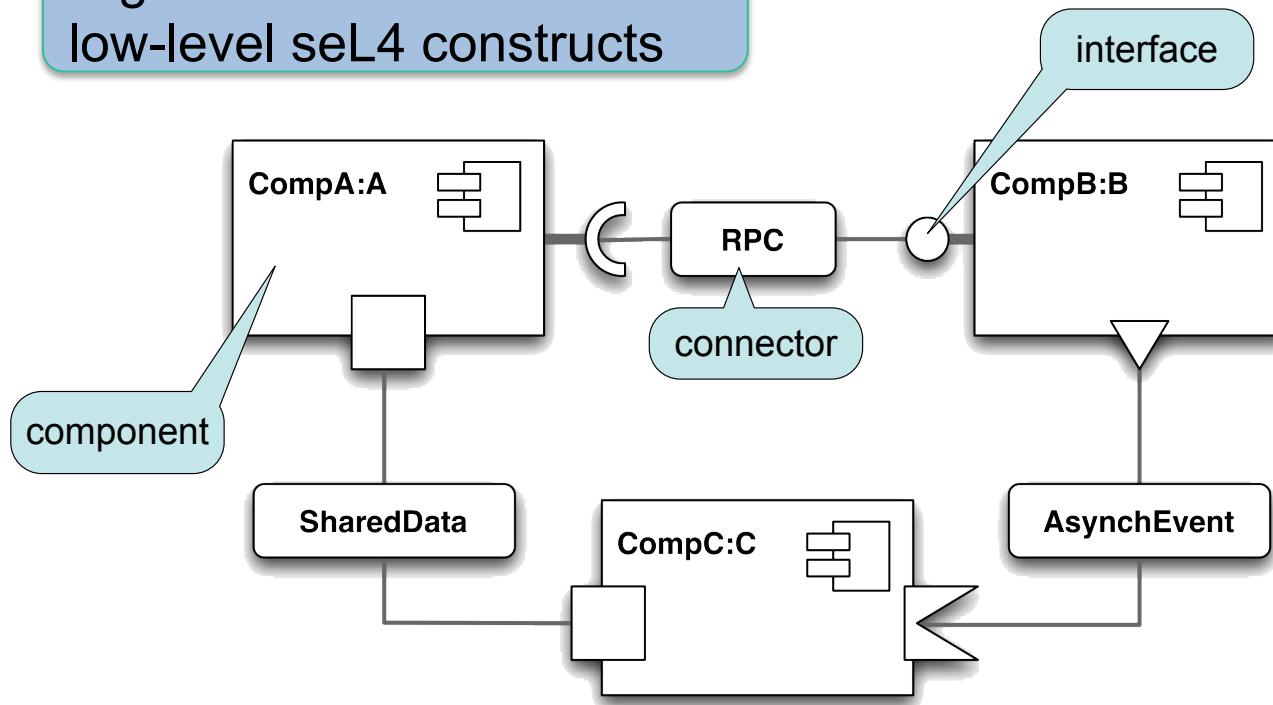




# Component Middleware: CAmkES



Higher-level abstractions of low-level seL4 constructs





# Example: DARPA HACMS



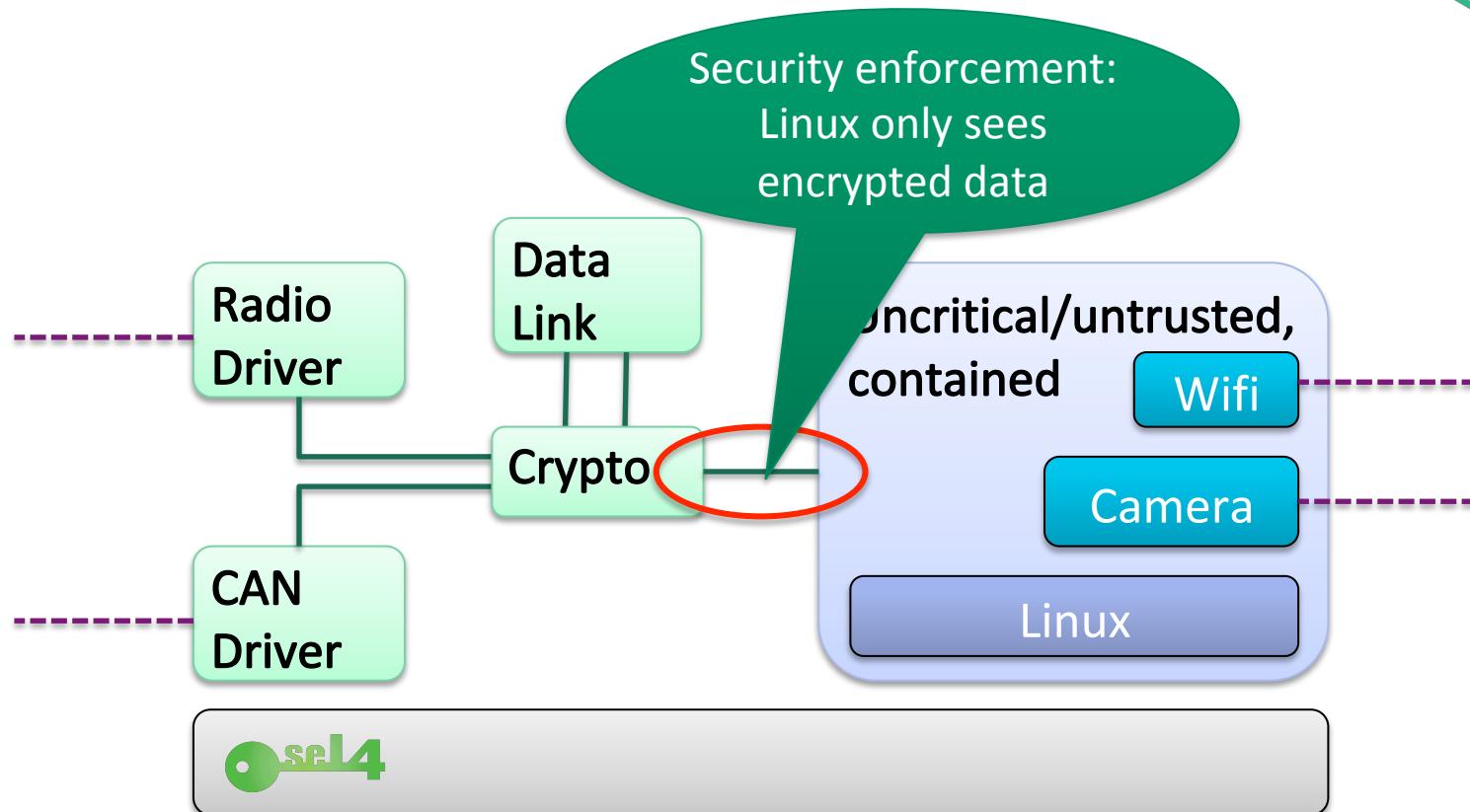
Retrofit  
existing  
system!



Develop  
technology

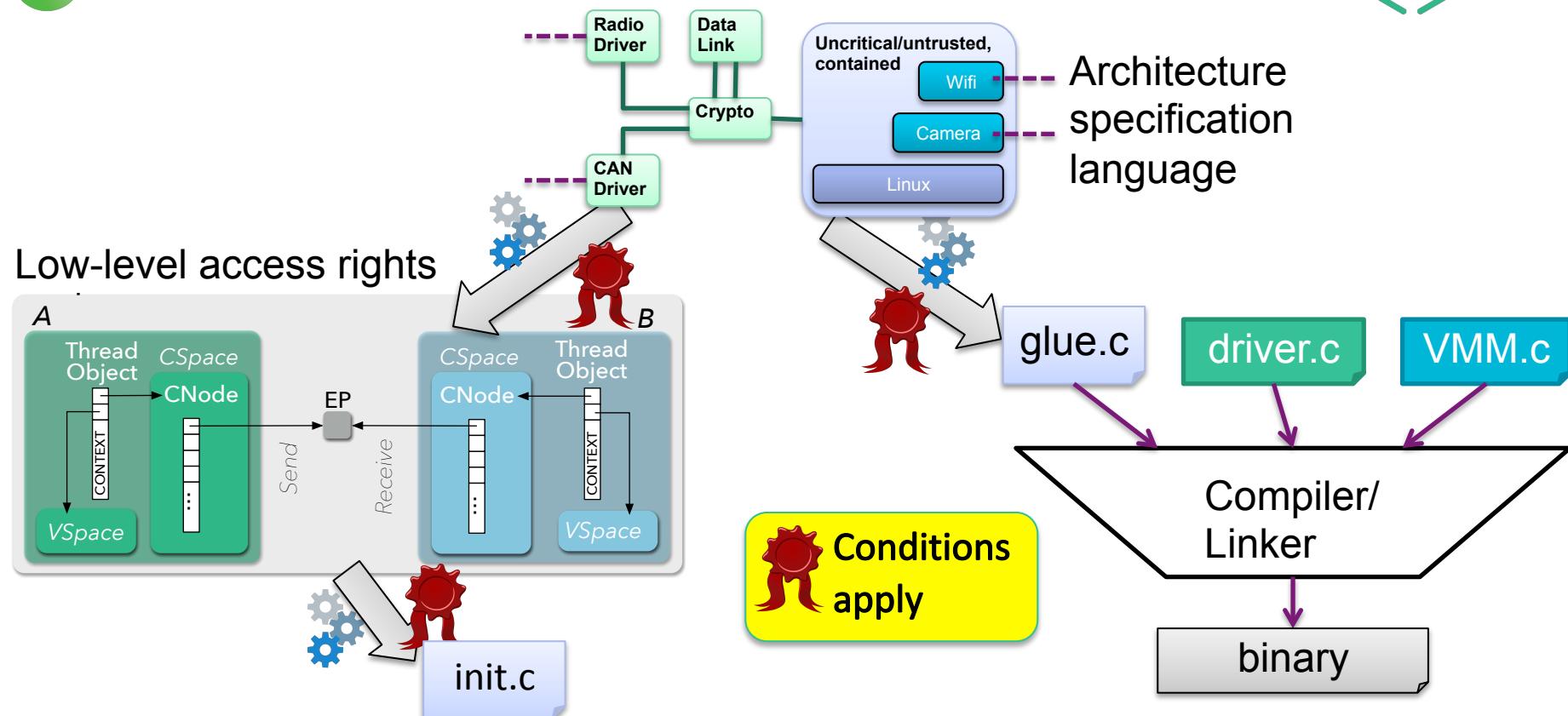


# seL4 HACMS UAV Architecture

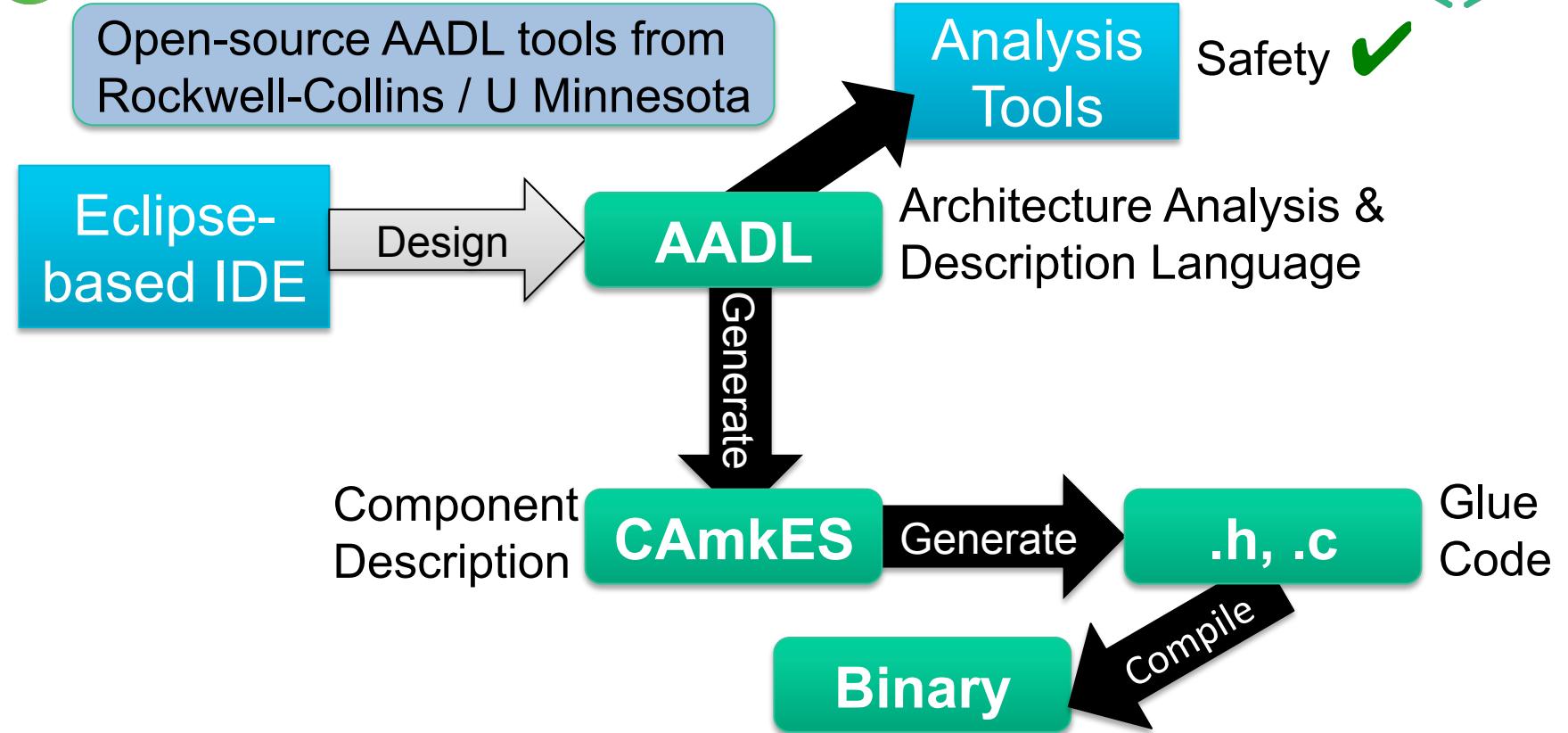




# sel4 Enforcing the Architecture



# seL4 Architecture Analysis



# seL4 Military-Grade Security



Cross-Domain Desktop Compositor

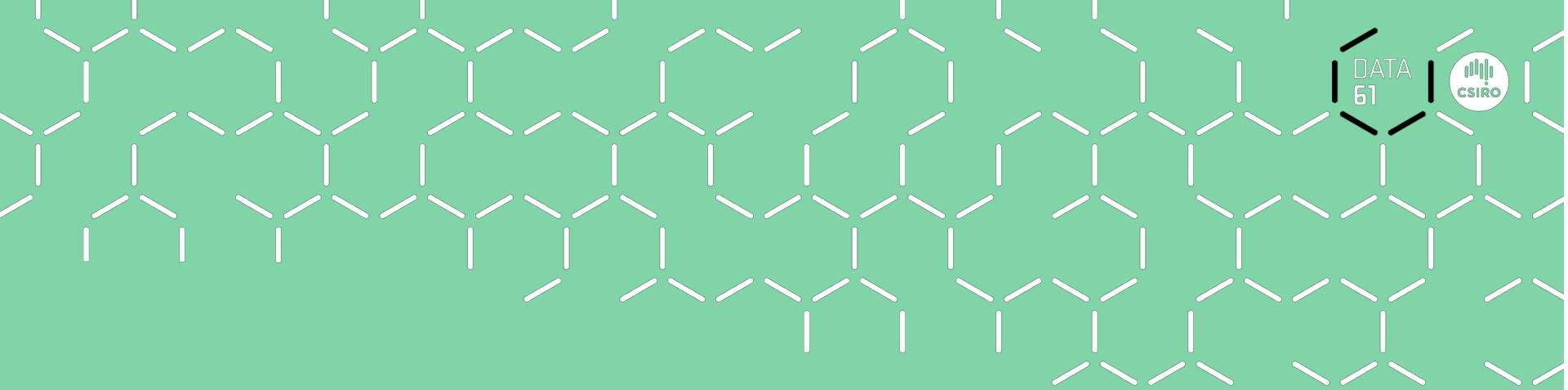


Multi-level secure terminal

- Successful defence trial in AU
- Evaluated in US, UK, CA
- Formal security evaluation soon

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





DATA  
61



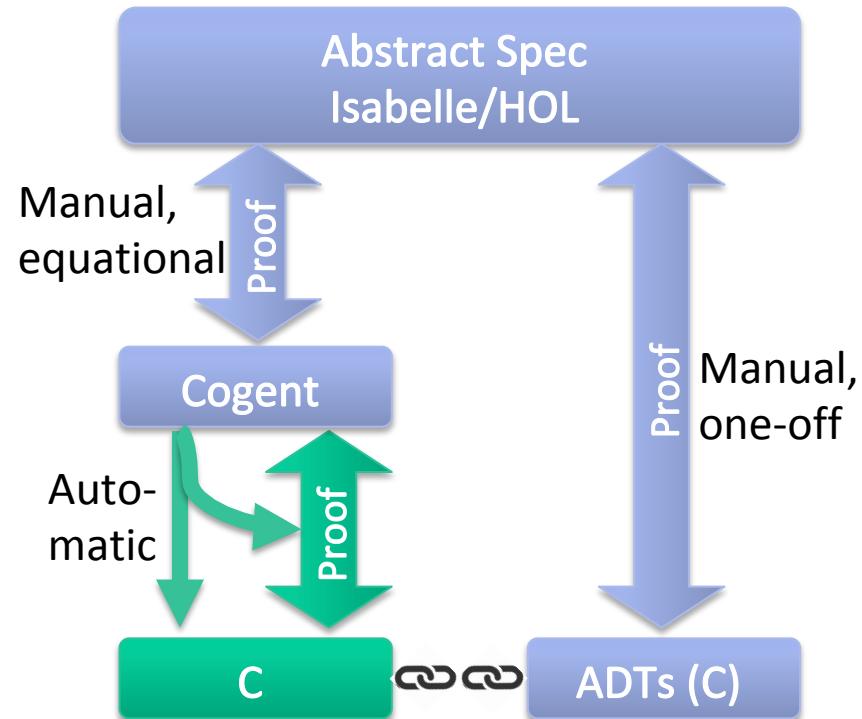
# Cogent: Reducing the Cost of Verification

# Cogent: Code & Proof Co-Generation



- Reduces the cost of formally verifying systems code
- Restricted, purely functional language
- Uniqueness type system
- Case-studies: BilbyFs, ext2, F2FS, VFAT

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



# Manual Proof Effort

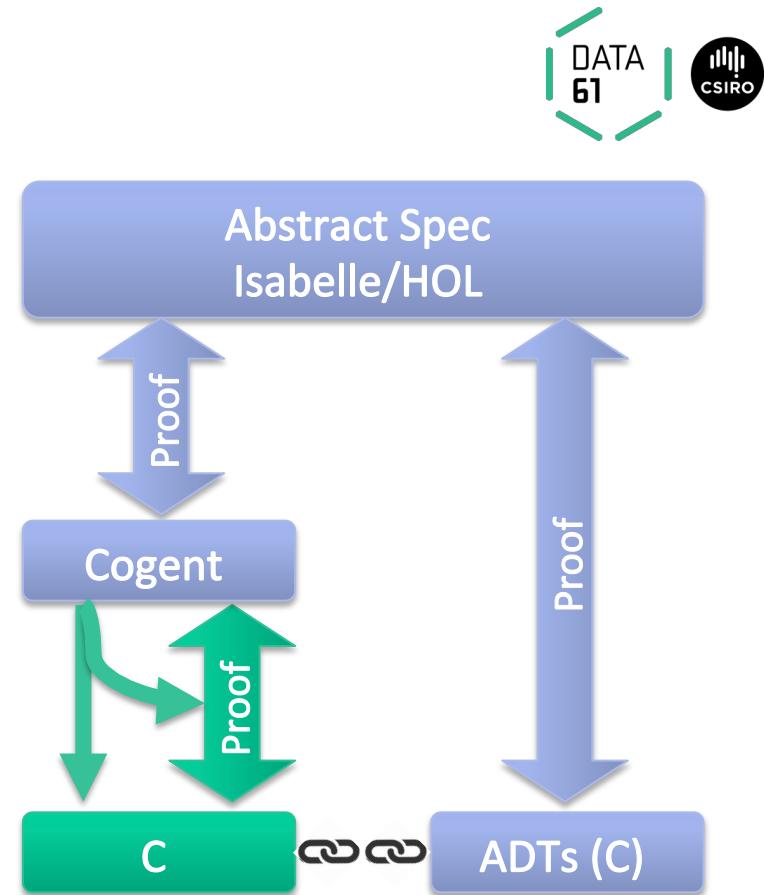


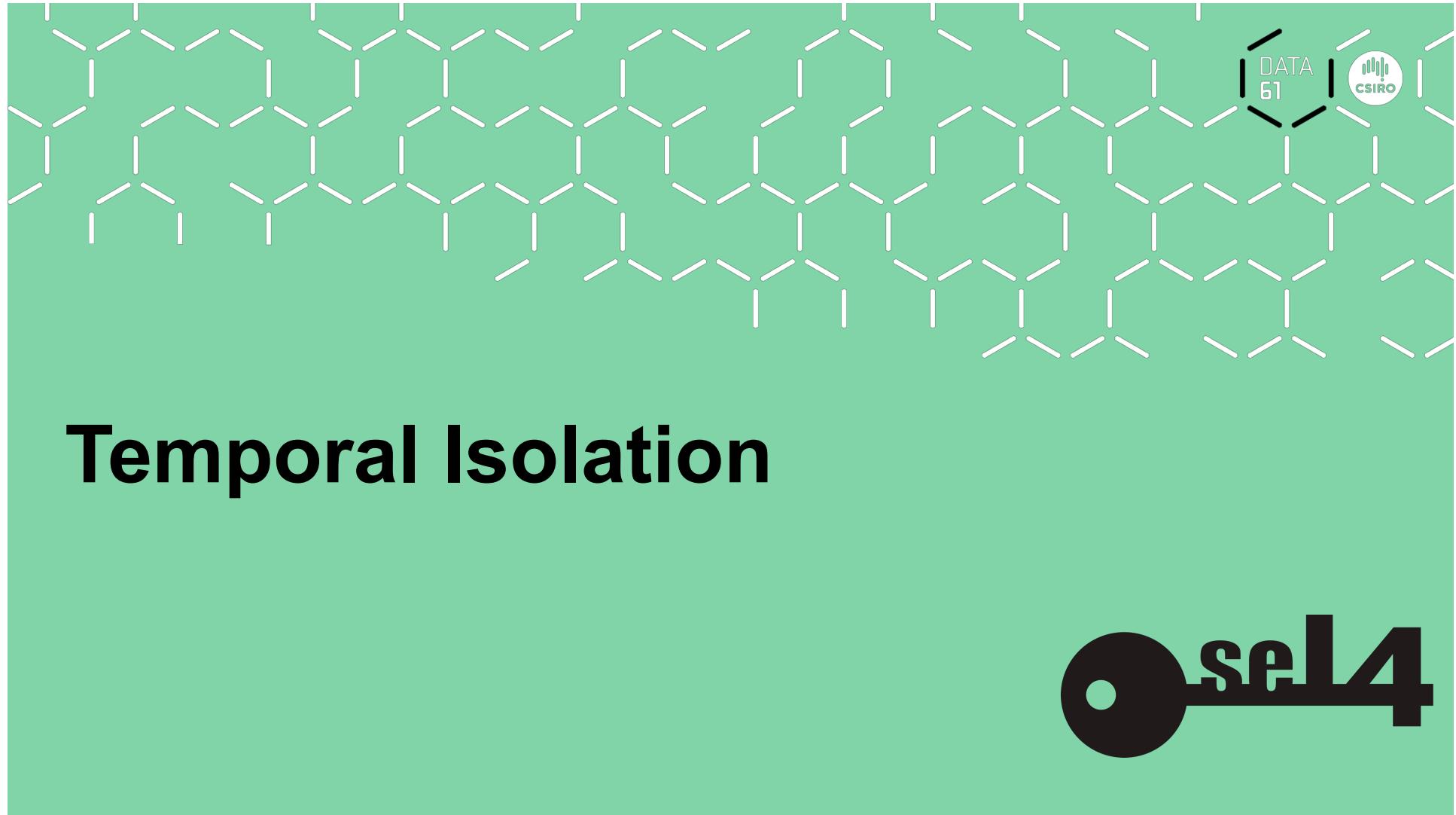
BilbyFS functions	Effort	Isabelle LoP	Cogent SLoC	Cost \$/SLoC	LoP/SLOC
isync() iget() library	9.25 pm	13,000	1,350	150	10
sync()-specific	3.75 pm	5,700	300	260	19
iget()-specific	1 pm	1,800	200	100	9
<b>seL4</b>	<b>12 py</b>	<b>180,000</b>	<b>8,700 C</b>	<b>350</b>	<b>20</b>

BilbyFS: 4,200 LoC Cogent

# Cogent: Present Work

- Relax type-system restrictions
  - purely linear types lead to excessive copying, 100% overhead
- Extend base language with domain-specific syntactic sugar
  - increased expressiveness → less code
  - automate boilerplate code
- Apply to other systems code
  - device drivers
  - network protocols stacks





# Temporal Isolation

# Temporal Interference

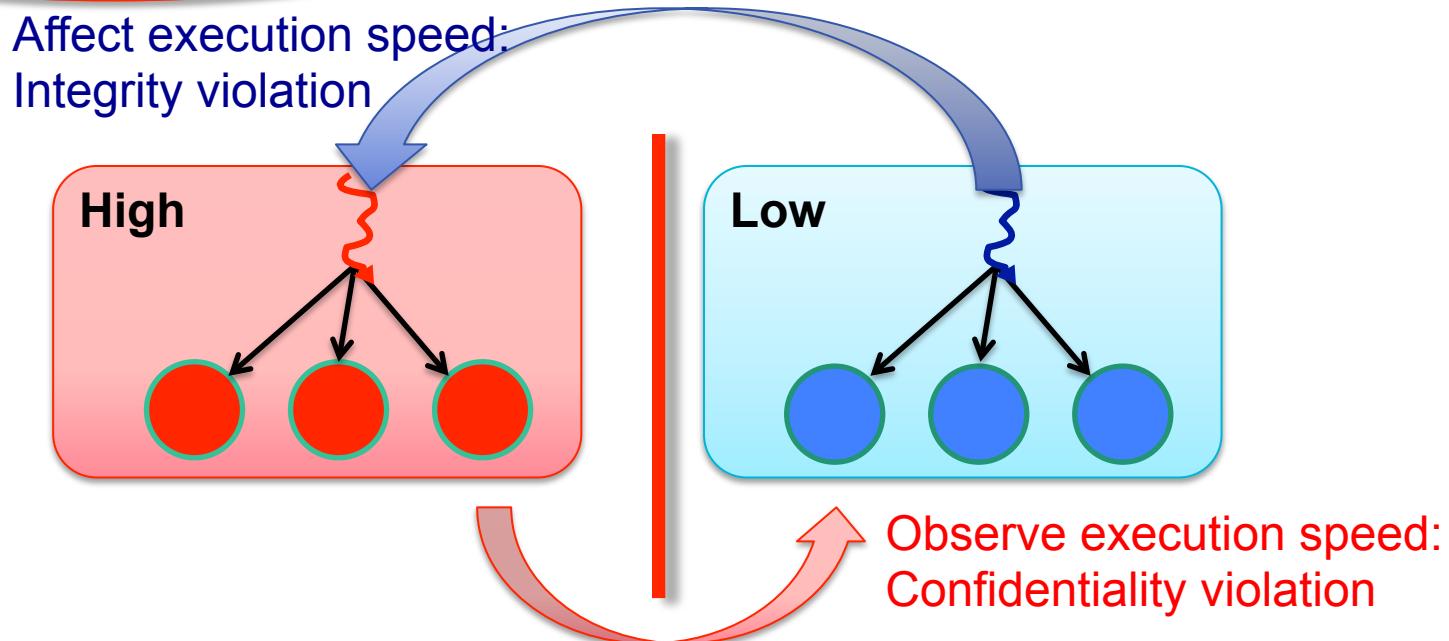


## Safety: Timeliness

- Execution interference

## Security: Confidentiality

- Leakage via timing channels





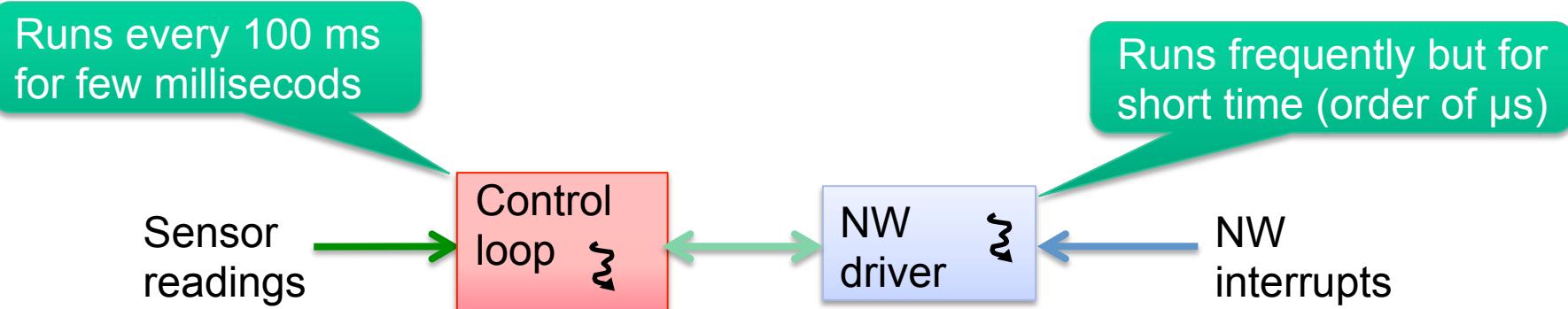
# Mixed-Criticality Scheduling: Enforcing Temporal Integrity

# Integration Challenge: Mixed Criticality

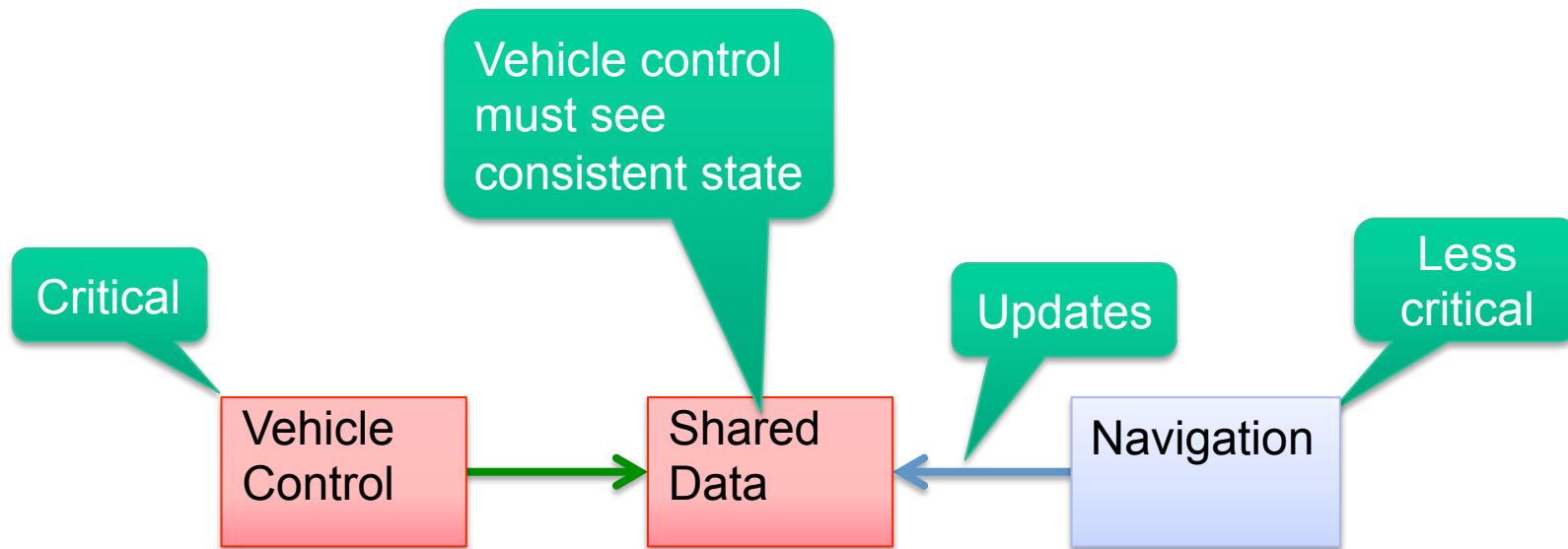


## NW driver must preempt control loop

- ... to avoid packet loss
- Driver must run at high prio
- Driver must be trusted not to monopolise CPU

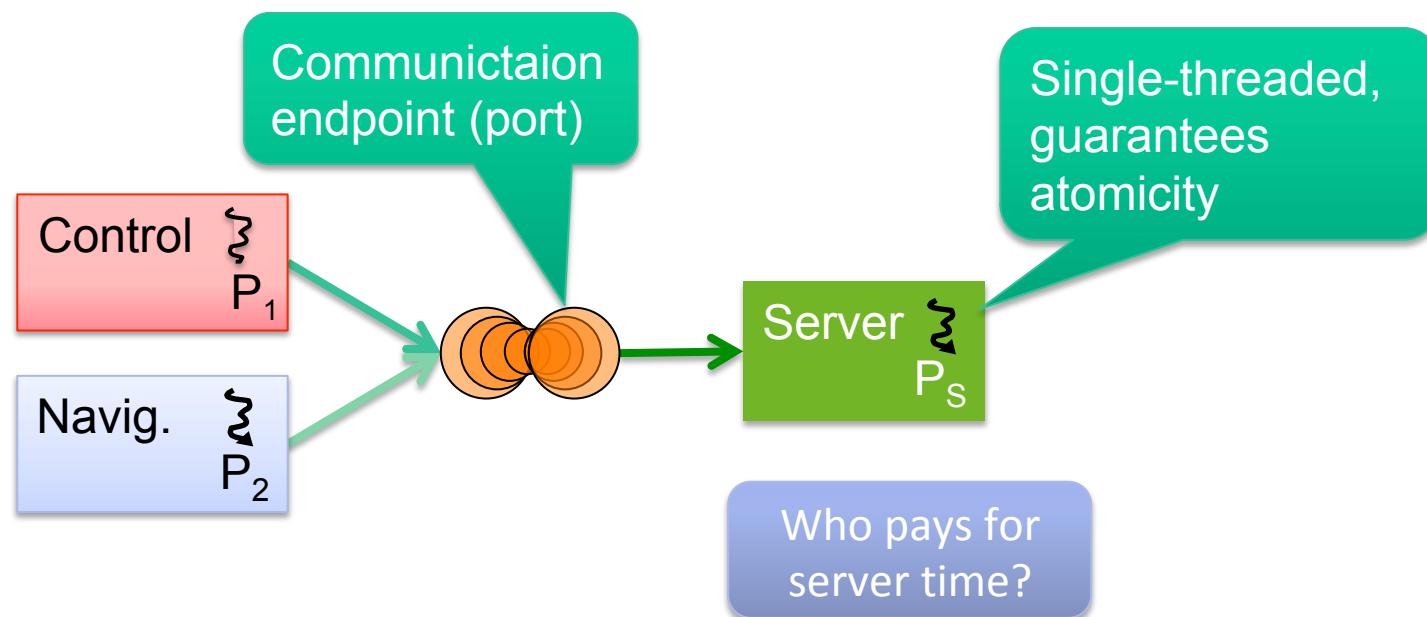


# Integration Challenge: Sharing





# Sharing Through *Resource Server*





# Scheduling Contexts: Time Caps



## Classical thread attributes

- Priority
- Time slice

Limits CPU access!

Not runnable if null

## New thread attributes

- Priority
- Scheduling context capability

Capability for time

### Scheduling context object

- T: period
- C: budget ( $\leq T$ )

$$C = 2 \\ T = 3$$

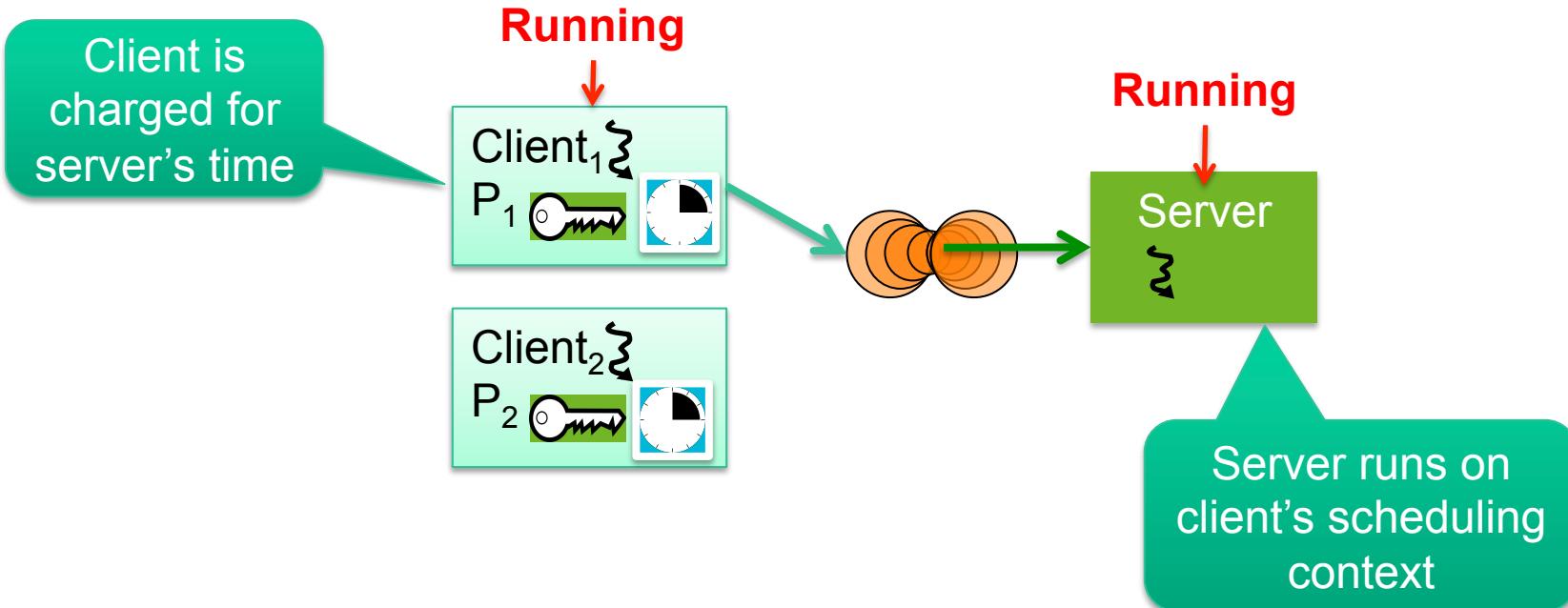


$$C = 250 \\ T = 1000$$



SchedControl capability conveys right to assign budgets (i.e. perform admission control)

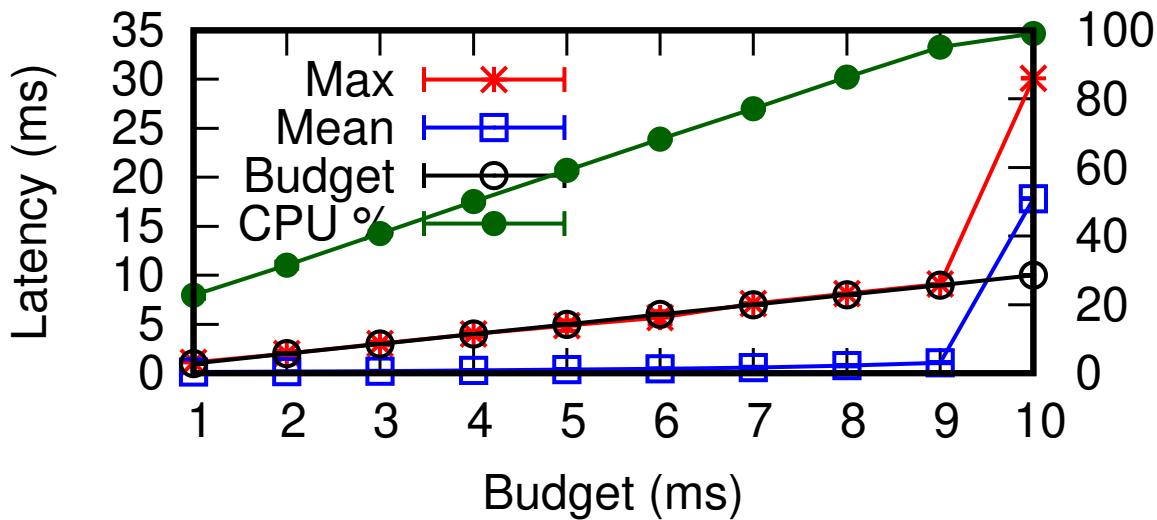
# seL4 Shared Server



# seL4 Isolation in Action



- High-prio CPU hog, budget limited, 10ms period
- Lower-prio UDP echo server, 10ms period



CPU utilisation (%)

MCS kernel will replace current mainline once verified

Scheduling-context capabilities: a principled, light-weight OS mechanism for managing time [Lyons et al, EuroSys'18]



# Time Protection: Removing Timing Channels

# Temporal Interference

## Safety: Timeliness

- Execution interference

Affect execution speed:  
Integrity violation

High



## Security: Confidentiality

- Leakage via timing channels

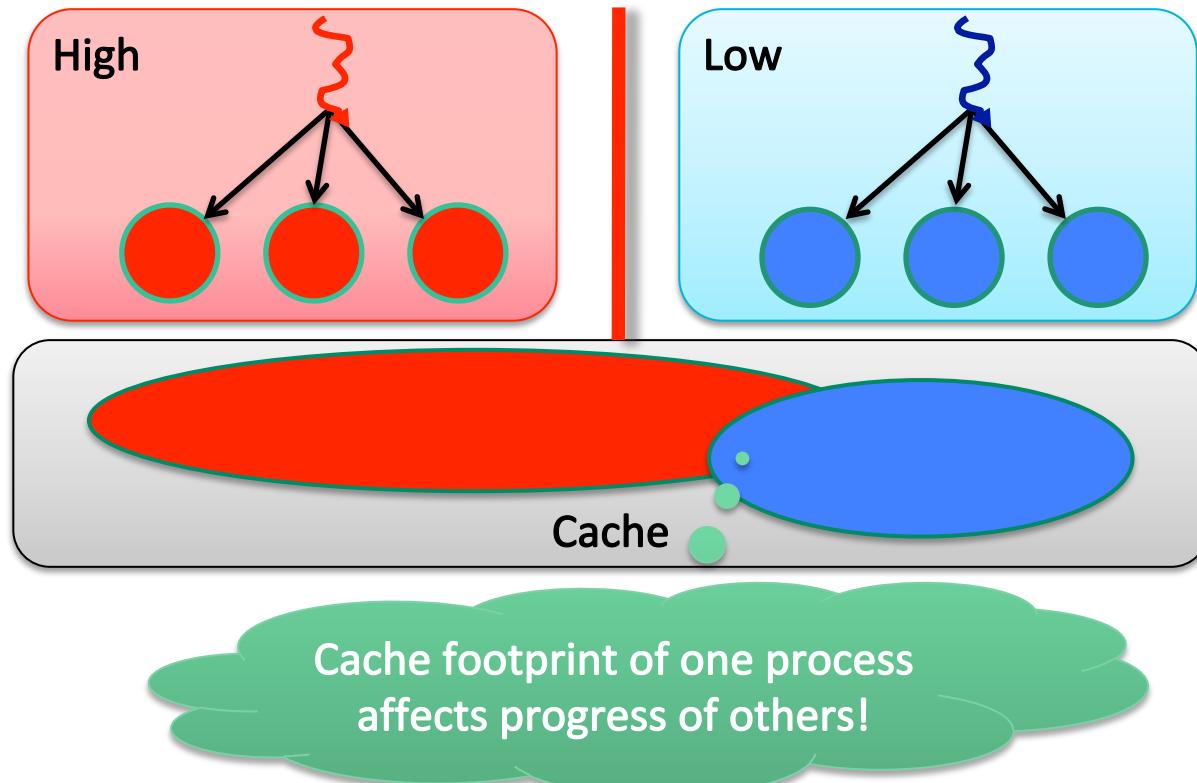
Low



Observe execution speed:  
Confidentiality violation



# Cause: Conflicts on Shared HW



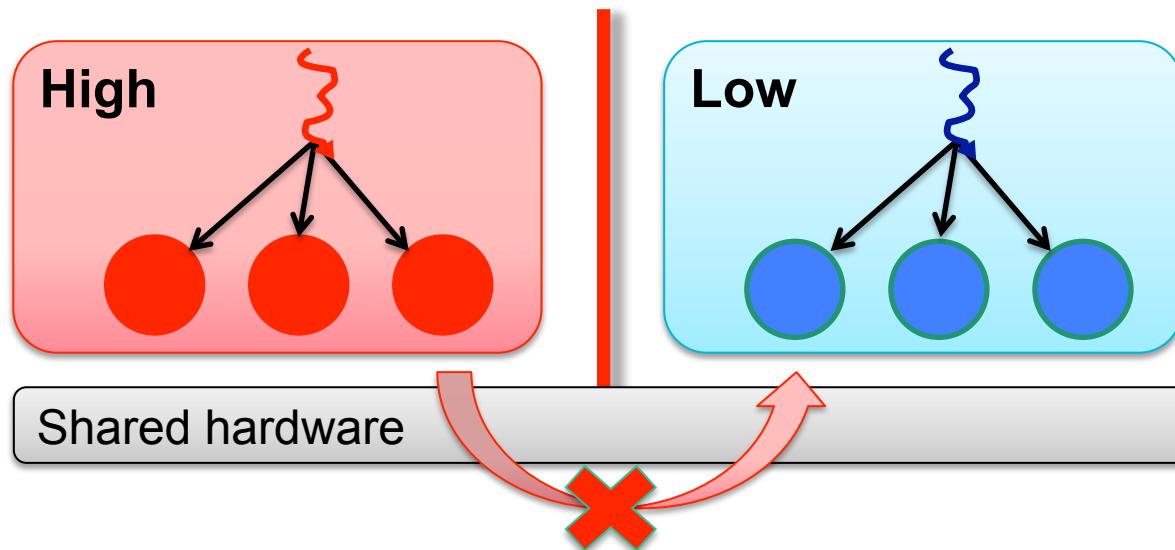
**Sharing can be:**

- Concurrent multicore, HW thread
- Time-shared

**“Caches” include:**

- L1 I-, D-cache
- TLB
- Branch predictor
- Instr. prefetcher
- Data prefetcher
- off-core caches & busses

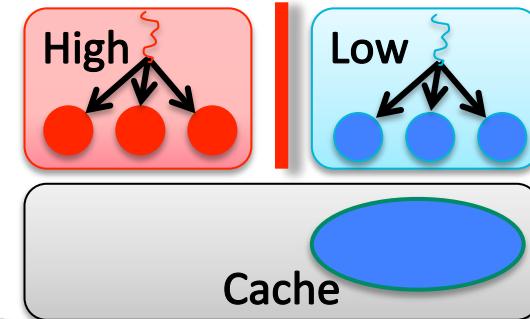
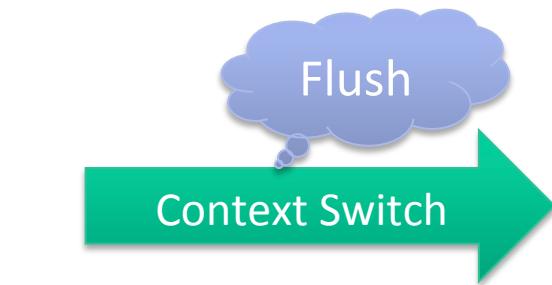
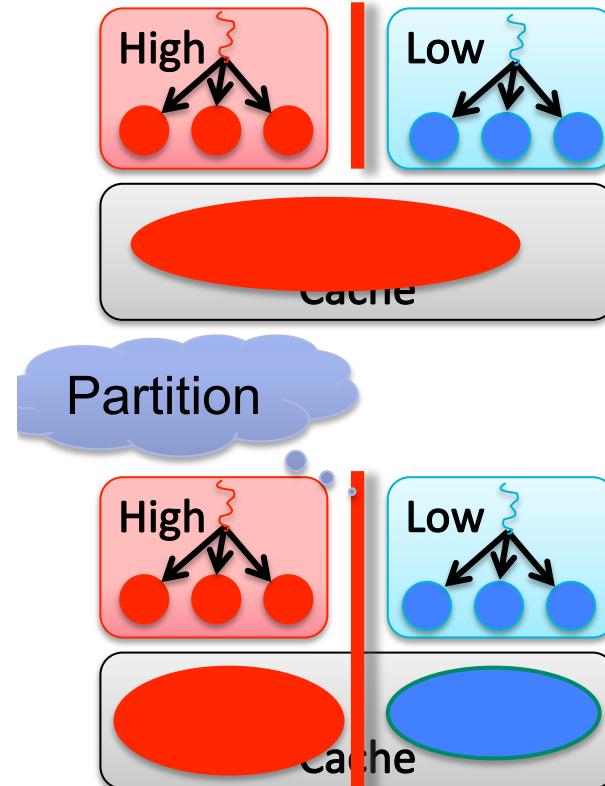
# OS Must Enforce *Time Protection*



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

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

# Time Protection: No Sharing of State



Need both!

Cannot partition on-core caches (L1, TLB, branch predictor, prefetchers)

- virtually-indexed
- OS cannot control

Flushing useless for concurrent access

- between HW threads, cores
- for stateless HW



# seL4 Flushing on Domain Switch



1.  $T_0 = \text{current\_time}()$
2. Switch context
3. Flush caches
4. Touch all code/data needed for return
5. Reprogram timer
6.  $\text{while } (T_0 + \text{WCET} < \text{current\_time}()) ;$
7. return

Latency depends  
on prior execution!

Ensure  
deterministic  
execution

Remove  
dependency



# Can Time Protection Be Verified?



1. Correct treatment of partitionable state:
  - Need hardware model that identifies all such state (augmented ISA)
  - Enables *functional correctness* argument:  
**No two domains can access the same physical state**
2. Correct flushing of non-partitionable 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!

Transforms timing  
channels into  
storage channels!



# How Can We Prove Time Padding?



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

To prove: padding loop terminates  
as soon as timer value  $\geq T_0 + WCET$

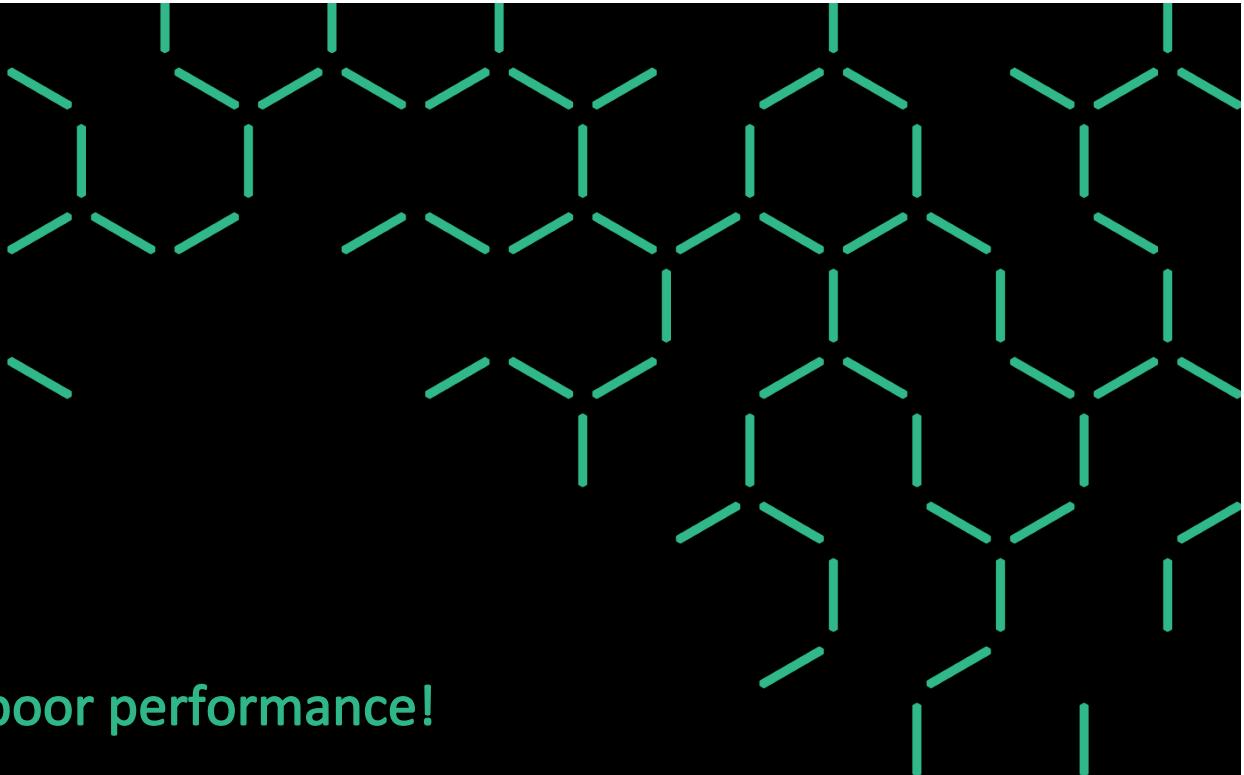
Functional  
property

 **seL4** Summary

- **seL4** has most comprehensive verified spatial isolation
- **Cogent** project aims to reduce cost of verified system components
  - file systems, device drivers, network protocol stacks
- **Scheduling contexts** introduce capability-controlled time management
  - support *mixed-criticality systems* where priority ≠ criticality
  - temporal integrity enforcement
- **Time protection** aims to eliminate microarchitectural timing channels
  - extend isolation to time
  - ideas on how to verify



# Thank you!



**Security is no excuse for poor performance!**

Gernot Heiser | Microkernel Dude

Gernot.Heiser@data61.csiro.au | @GernotHeiser

<https://sel4.systems>

