



# How to Build Truly Trustworthy Systems

Gernot Heiser

NICTA and University of New South Wales  
Sydney, Australia



Australian Government

Department of Broadband, Communications  
and the Digital Economy

Australian Research Council

## NICTA Funding and Supporting Members and Partners



## Windows

An exception 06 has occurred at 0028:C11B3ADC in VxD DiskTSD(03) + 00001660. This was called from 0028:C11B40C8 in VxD voltrack(04) + 00000000. It may be possible to continue normally.

- \* Press any key to attempt to continue.
- \* Press CTRL+ALT+RESET to restart your computer. You will lose any unsaved information in all applications.

Press any key to continue

# Present Systems are *NOT* Trustworthy!



# What's Next?

---



So, why don't  
we prove  
trustworthiness  
?

**Claim:**

**A system must be considered *not trustworthy*  
unless *proved* otherwise!**

*Corollary [with apologies to Dijkstra]:*

Testing, code inspection, etc. can only show  
*lack of trustworthiness!*

# Core Issue: Complexity

- Massive functionality of C devices  
⇒ huge software stacks

- How secure are your payments?



- Increasing usability requirements

- Wearable or implanted

- Patient-operated

- GUIs next to life-critical

**Systems far too complex to prove their trustworthiness!**

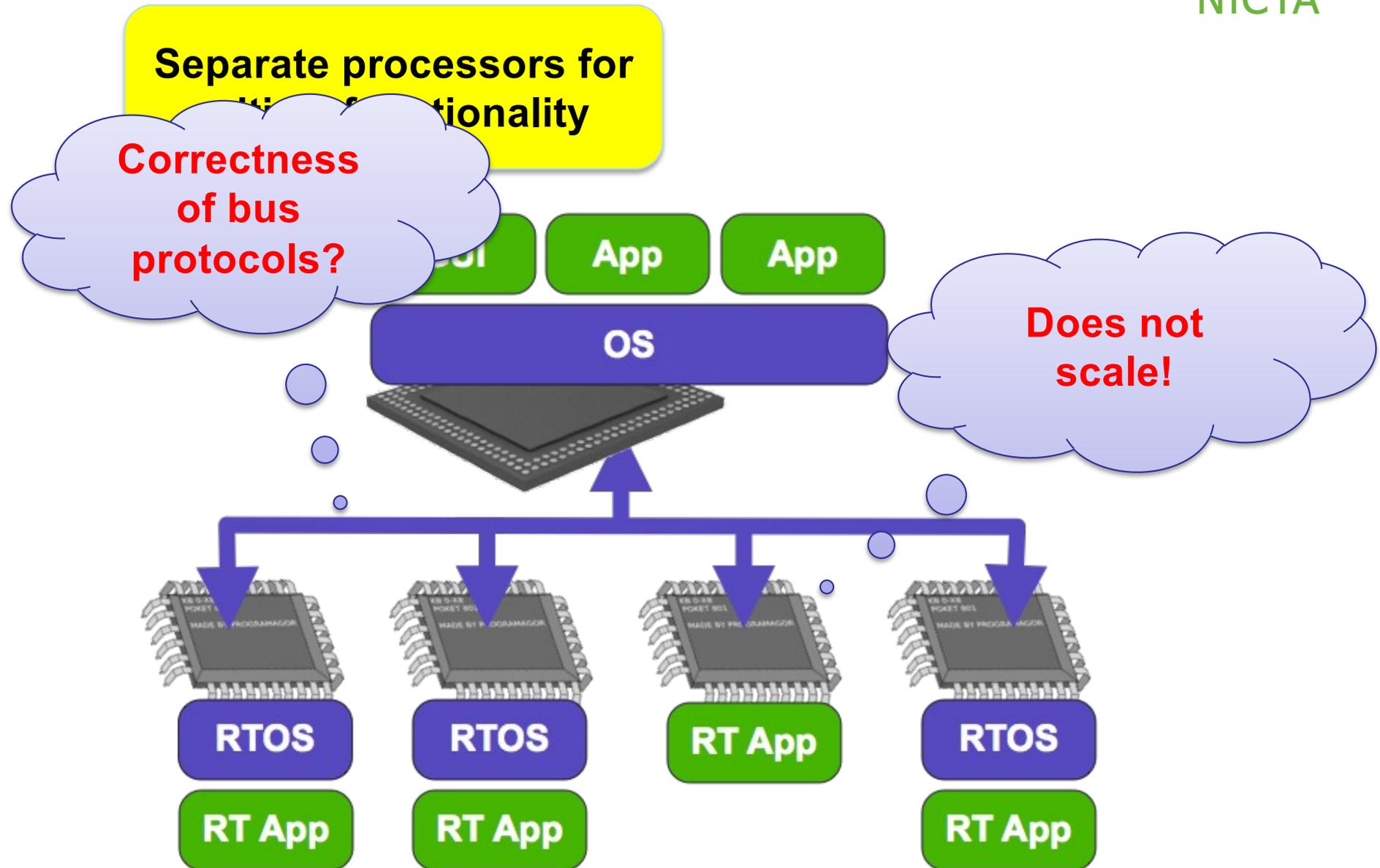
- On-going integration of new services

- Automotive infotainment and navigation

- Gigabytes of software on 100 CPUs...



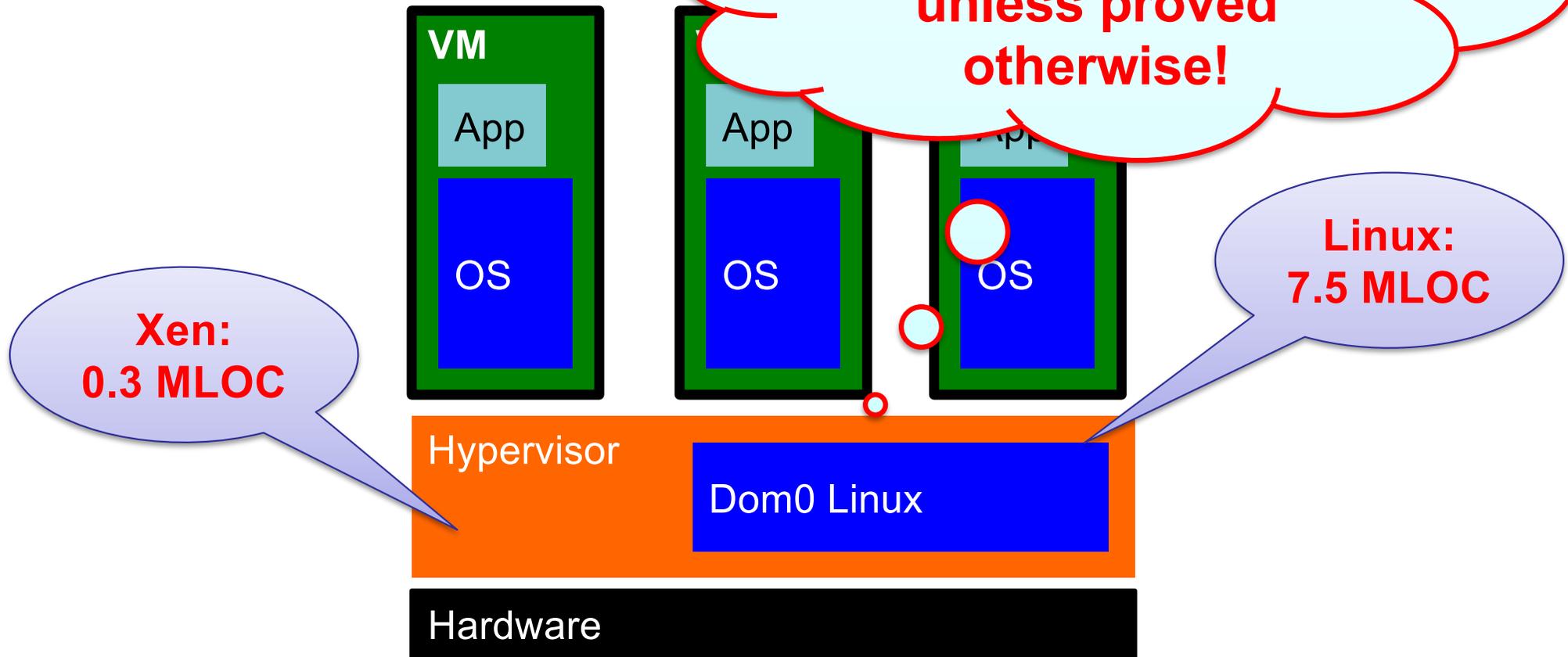
# Dealing with Complexity: Physical Isolation



# How About Logical Isolation?

Shared processor with software isolation

Remember: A system is *not trustworthy* unless proved otherwise!



# Our Vision: Trustworthy Systems



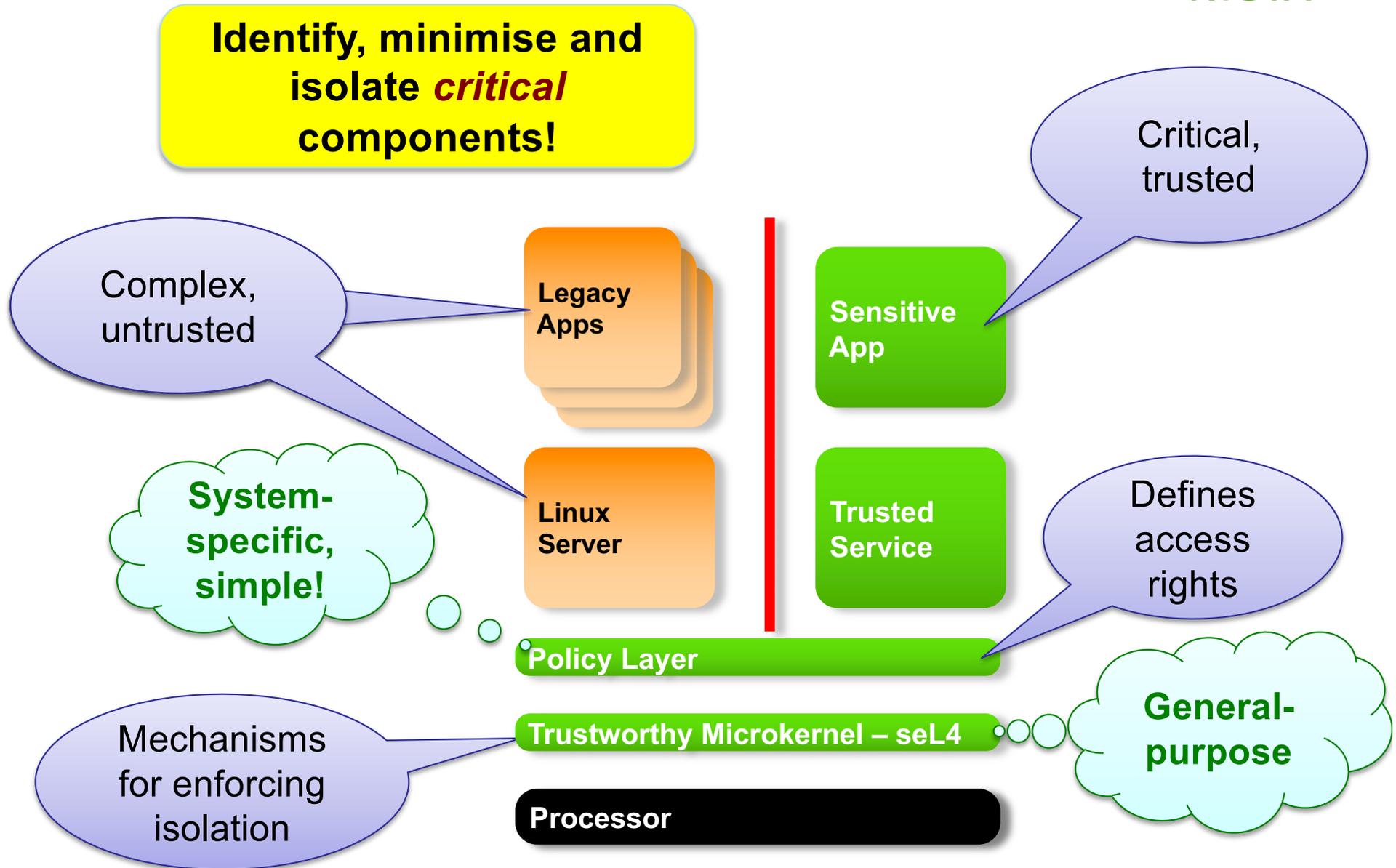
Suitable for  
real-world  
systems

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

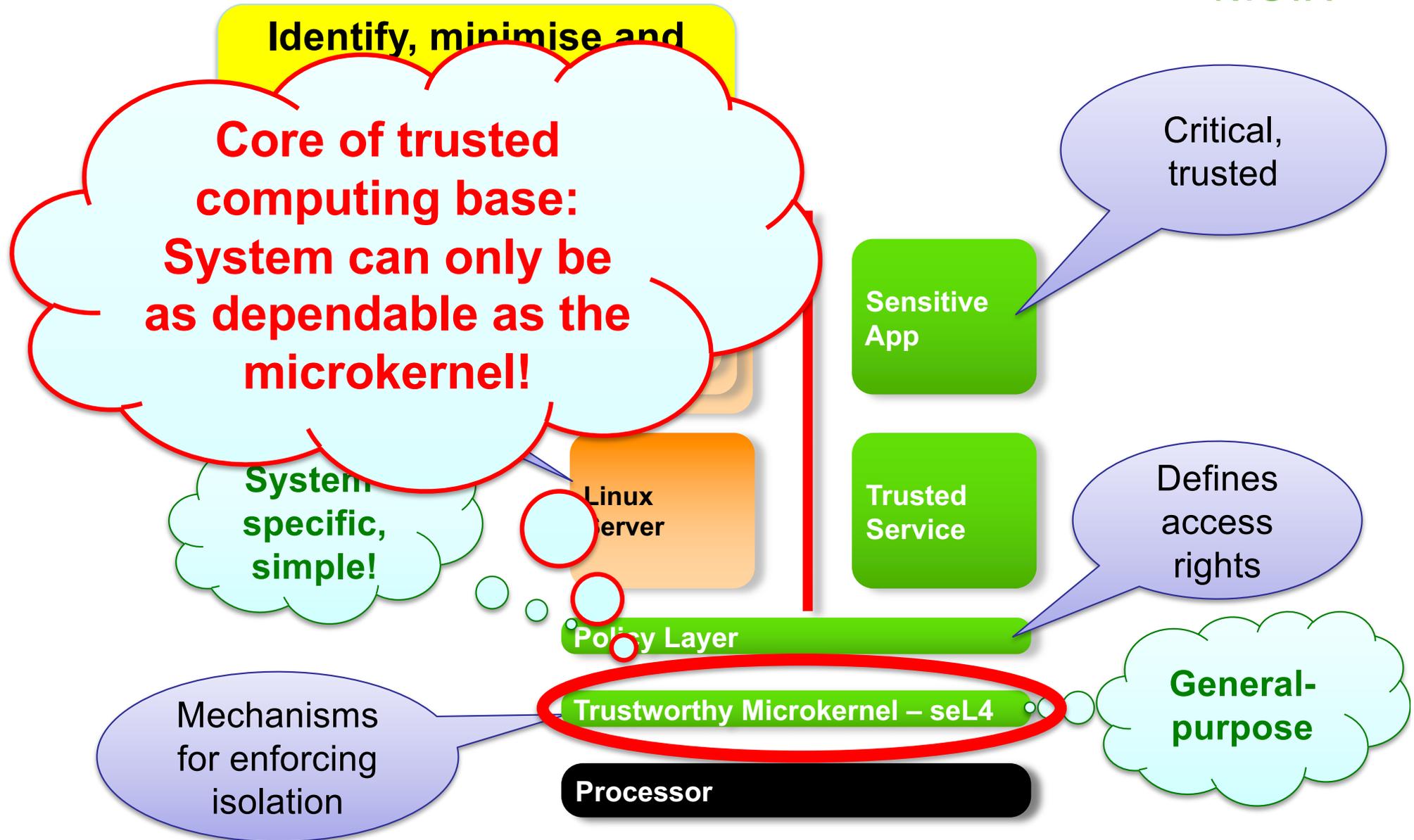
Hard  
*guarantees* on  
safety/security/  
reliability



# Isolation is Key!



# Isolation is Key!



# NICTA 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



## 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!



# Agenda

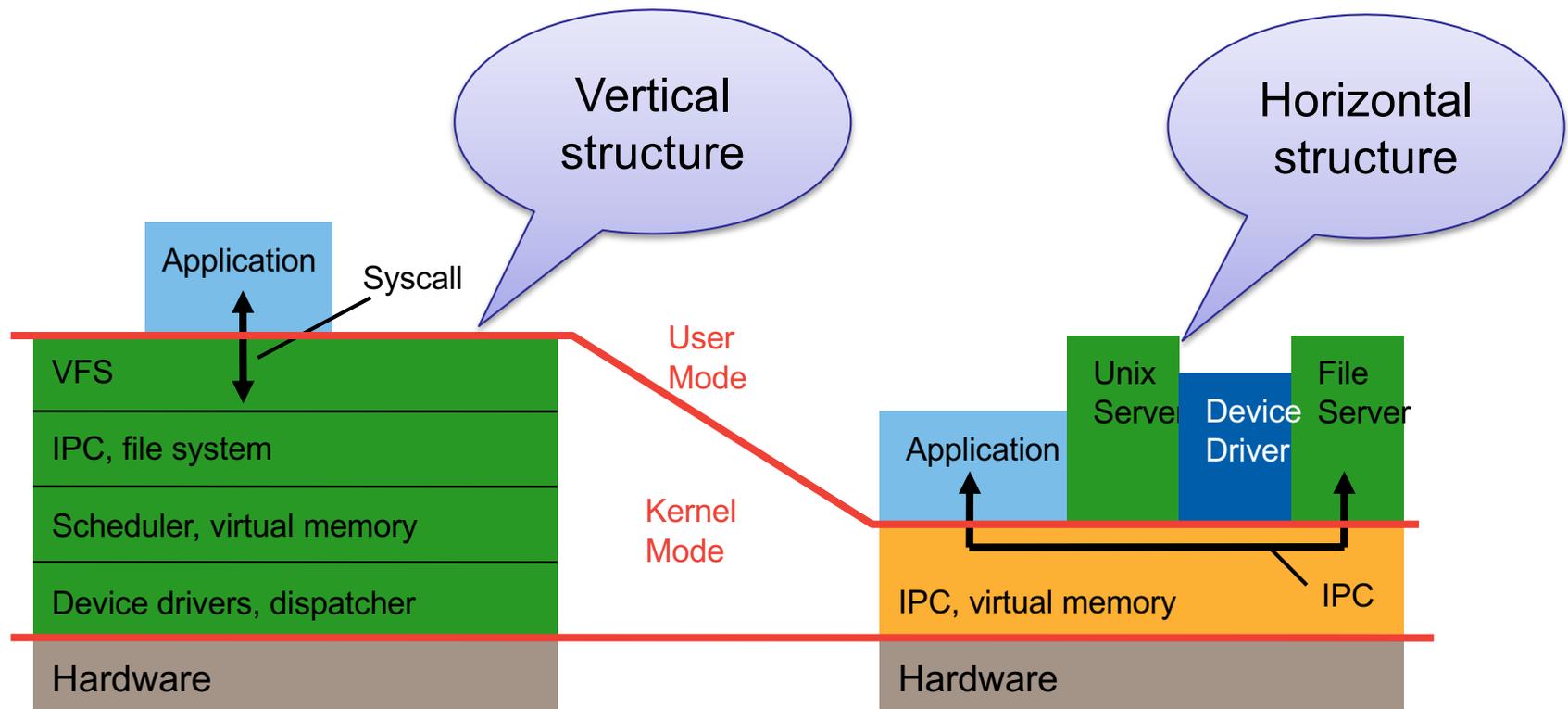
---



- Motivation
- **What is a microkernel, and what is L4?**
- seL4 – designed for trustworthiness
- Establishing trustworthiness
- From kernel to system
- Sample system 1: Secure access controller
- Sample system 2: RapiLog

# Monolithic Kernels vs Microkernels

- Idea of microkernel:
  - Flexible, minimal platform, extensible
  - Mechanisms, not policies
  - Goes back to Nucleus [Brinch Hansen, CACM'70]

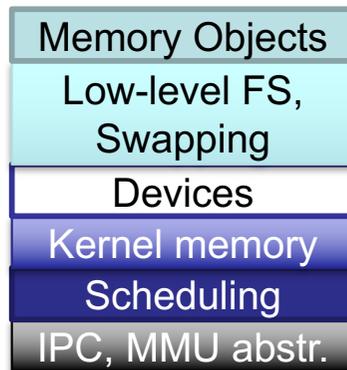


# Microkernel Evolution



## First generation

- Eg Mach ('87)



- 180 syscalls
- 100 kLOC
- 100  $\mu$ s IPC

## Second generation

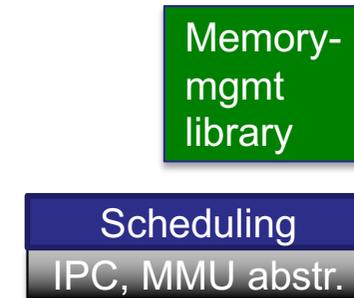
- Eg L4 ('95)



- ~7 syscalls
- ~10 kLOC
- ~ 1  $\mu$ s IPC

## Third generation

- seL4 ('09)



- ~3 syscalls
- 9 kLOC
- < 1  $\mu$ s IPC

## 2<sup>nd</sup>-Generation Microkernels

- 1<sup>st</sup>-generation kernels (Mach, Chorus) were a failure
  - Complex, inflexible, slow
- L4 was first 2<sup>nd</sup>-G microkernel [Liedtke, SOSP'93, SOSP'95]
  - Radical simplification & manual micro-optimisation, fast IPC

A concept is tolerated inside the microkernel only if moving it outside the kernel, i.e. permitting competing implementations, would prevent the implementation of the system's required functionality

- Family of L4 kernels:
  - Original GMD assembler kernel ('95)
  - Fiasco (Dresden '98), Hazelnut (Karlsruhe '99), Pistachio (Karlsruhe/UNSW '02), L4-embedded (NICTA '04)
    - L4-embedded commercialised as OKL4 by Open Kernel Labs
    - Deployed in >1.5 billion phones
  - Commercial clones (PikeOS, P4, CodeZero, ...)
  - Approach adopted e.g. in QNX ('82) and Green Hills Integrity ('90s)

# Microkernel Principles: Minimality

---



Strict adherence to minimality leads to a very small kernel

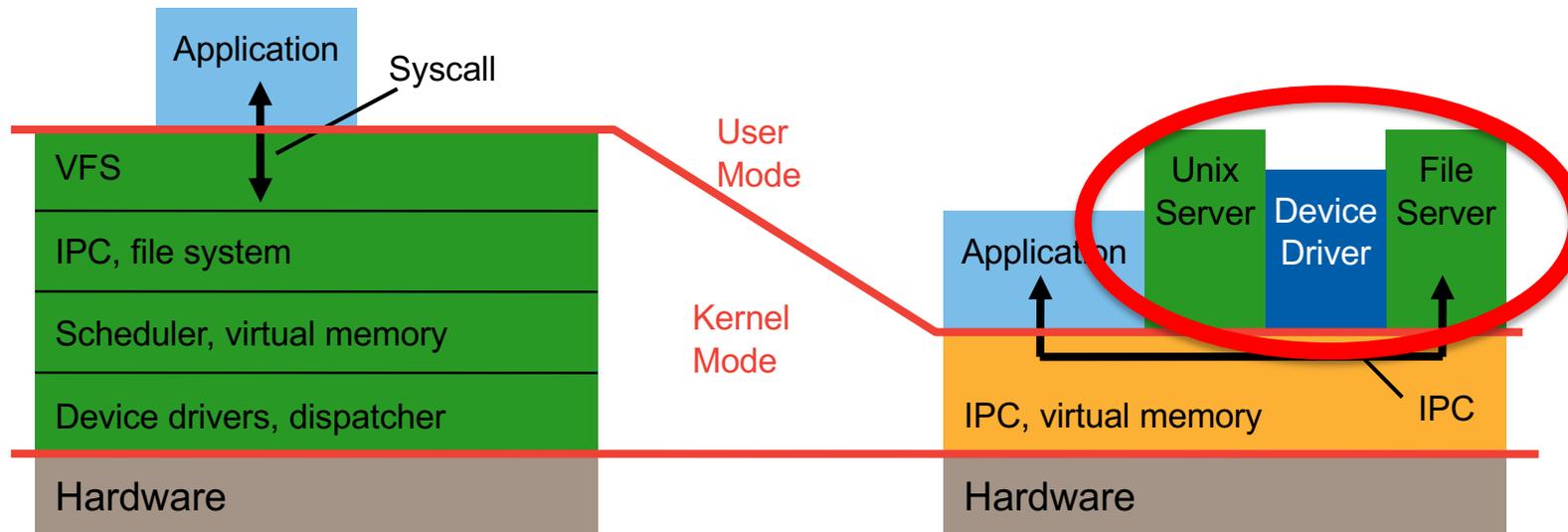
## Advantages:

- Easy to implement, port?
  - in practice limited architecture-specific micro-optimization
- Less code to optimise
- Hopefully enables a minimal *trusted computing base* (TCB)
  - small attack surface, fewer failure modes
- Easier debug, maybe even *prove* correct?

## Challenges:

- API design: generality with small code base
- Kernel design and implementation for high performance
  - ... and correctness!

# Consequence of Minimality: User-level Services



- Kernel provides no services, only mechanisms
- Strongly dependent on fast IPC and exception handling

# Microkernel Principles: Policy Freedom

---



Consequence  
of generality  
& minimality

**A true microkernel must be free of policy!**

## Policies limit

- May be good for many cases, but always bad for some
- Example: disk pre-fetching

## “General” policies lead to bloat

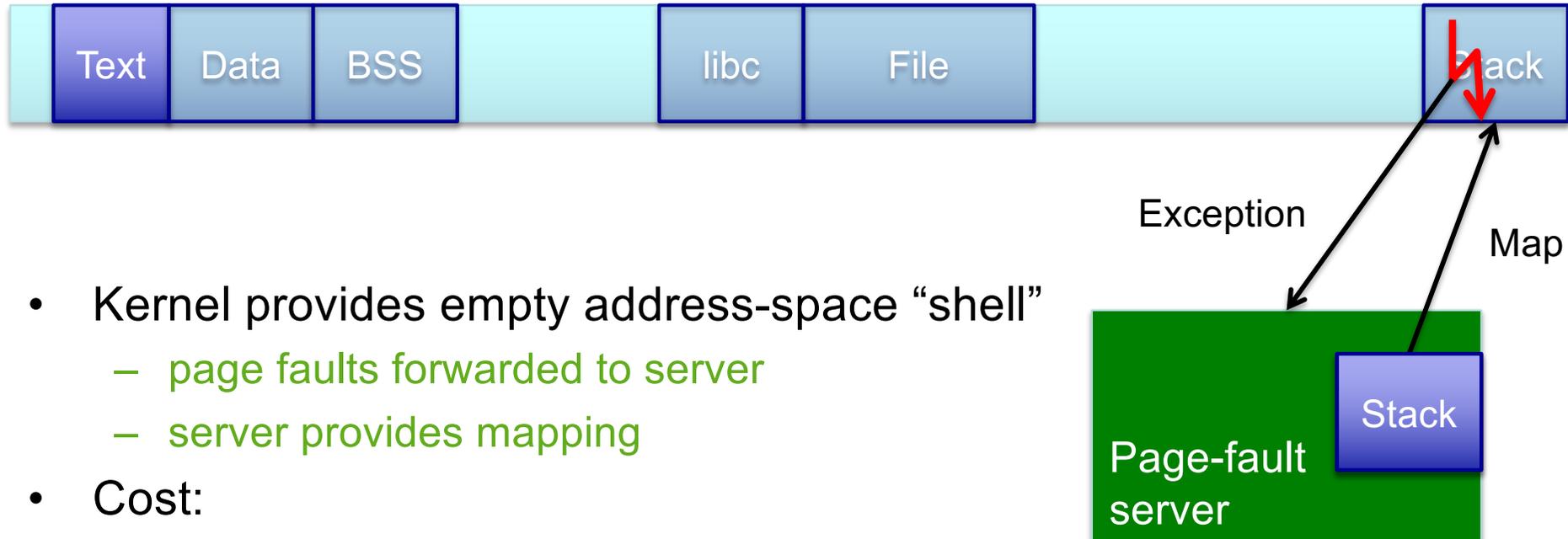
- Implementing combination of policies
- Try to determine most appropriate one at run-time

# Policy Example: Address-Space Layout



- Kernel determines layout, knows executable format, allocates stack
  - limits ability to import from other OSes
  - cannot change layout
    - small non-overlapping address spaces beneficial on some archs
  - kernel loads apps, sets up mappings, allocates stack
    - requires file system in kernel or interfaced to kernel
    - bookkeeping for revocation & resource management
    - heavyweight processes
  - memory-mapped file API

# Policy-Free Address-Space Management



- Kernel provides empty address-space “shell”
  - page faults forwarded to server
  - server provides mapping
- Cost:
  - 1 round-trip IPC, plus mapping operation
    - mapping may be side effect of IPC
    - kernel may expose data structure
  - kernel mechanism for forwarding page-fault exception
- “External pagers” first appeared in Mach [Rashid et al, '88]
  - ... but were optional

# What Mechanisms?

---

- Fundamentally, the microkernel must abstract
  - *Physical memory*
  - *CPU*
  - *Interrupts/Exceptions*
- Unfettered access to any of these bypasses security
  - No further abstraction needed for devices
    - memory-mapping device registers and interrupt abstraction suffices
    - ...but some generalised memory abstraction needed for I/O space
- Above isolates execution units, hence microkernel must also provide
  - *Communication* (traditionally referred to as *IPC*)
  - *Synchronization*

# What Mechanisms?

---



## Traditional hypervisor vs microkernel abstractions

Resource	Hypervisor	Microkernel
Memory	Virtual MMU (vMMU)	Address space
CPU	Virtual CPU (vCPU)	Thread or scheduler activation
Interrupt	Virtual IRQ (vIRQ)	IPC message or signal
Communication	Virtual NIC	Message-passing IPC
Synchronization	Virtual IRQ	IPC message

# Issues of 2G L4 Kernels

---

- L4 solved performance issue [Härtig et al, SOSPP'97]  
... but left a number of security issues unsolved
- Problem: ad-hoc approach to protection and resource management
  - Global thread name space  $\Rightarrow$  covert channels
  - Threads as IPC targets  $\Rightarrow$  insufficient encapsulation
  - Single kernel memory pool  $\Rightarrow$  DoS attacks
  - Insufficient delegation of authority  $\Rightarrow$  limited flexibility, performance
- Addressed by seL4
  - Designed to support safety- and security-critical systems

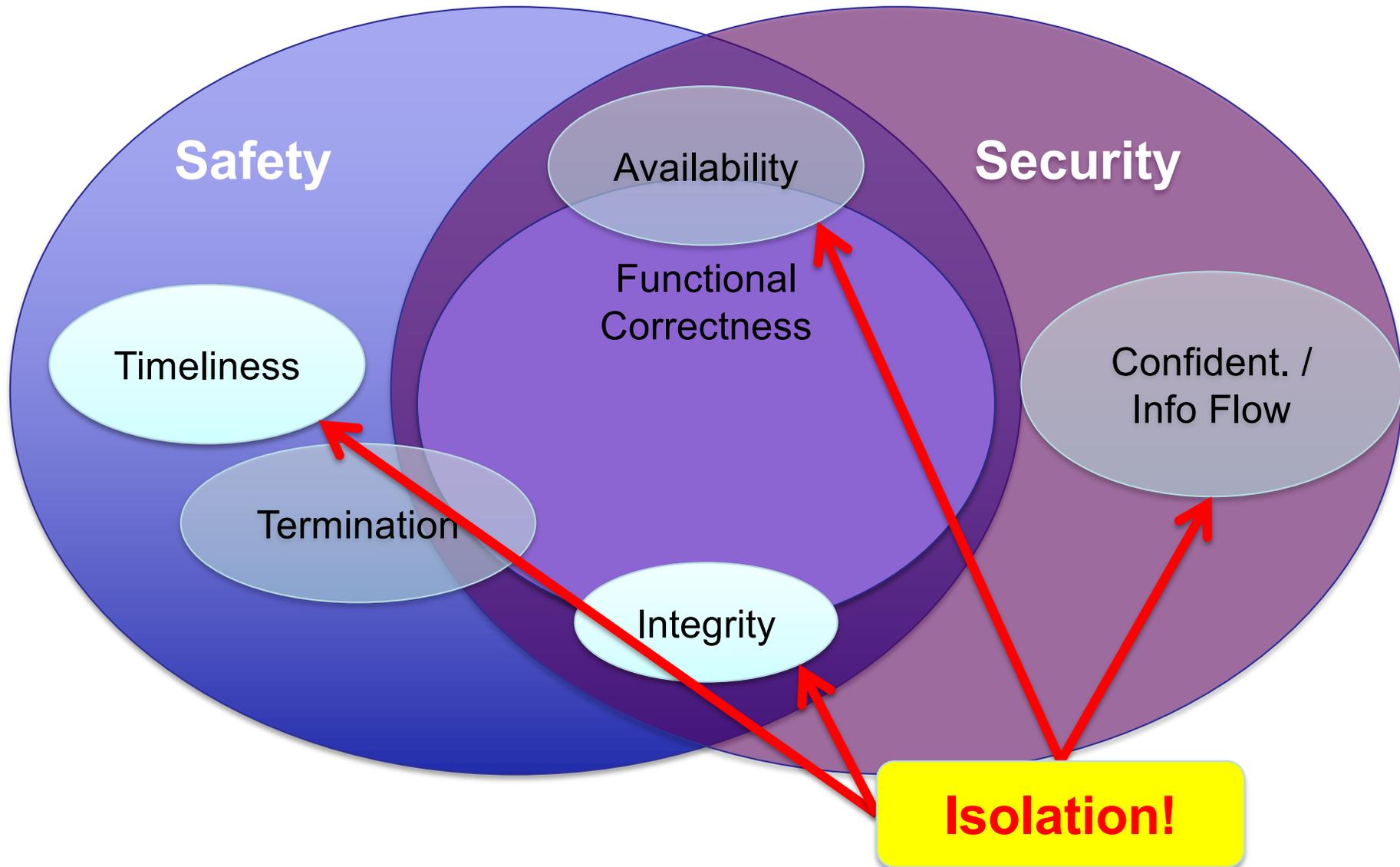
# Agenda

---

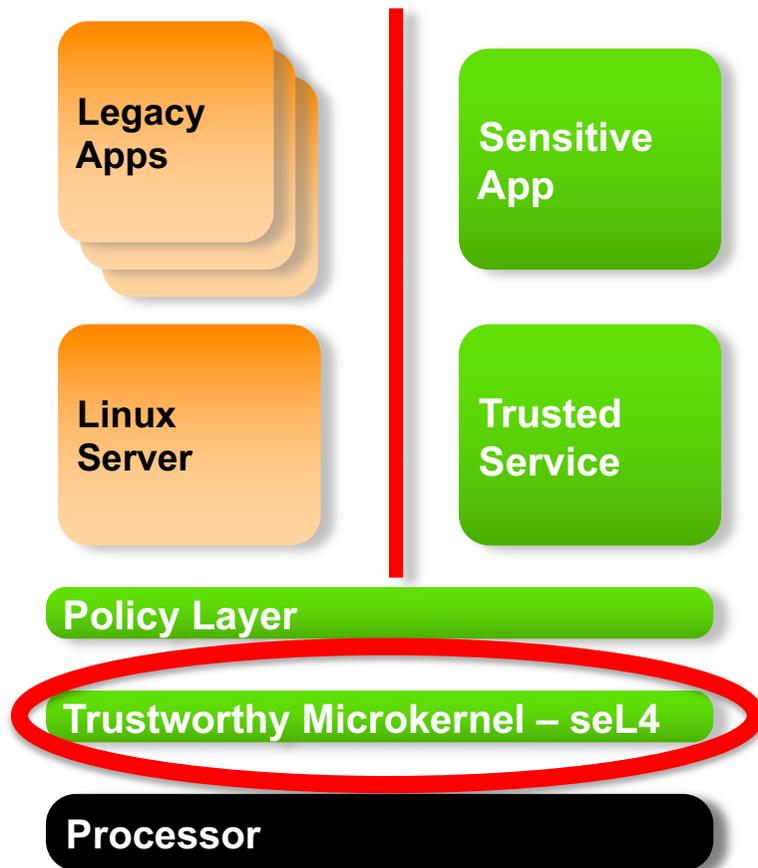


- Motivation
- What is a microkernel, and what is L4?
- **seL4 – designed for trustworthiness**
- Establishing trustworthiness
- From kernel to system
- Sample system 1: Secure access controller
- Sample system 2: RapiLog

# Requirements for Trustworthy Systems



# seL4 Design Goals



1. **Isolation**
  - **Strong partitioning!**
2. **Formal verification**
  - **Provably trustworthy!**
3. **Performance**
  - **Suitable for real world!**

# Fundamental Design Decisions for seL4



1. Memory management is user-level responsibility

- Kernel never allocates memory (post-boot)
- Kernel objects controlled by user-mode servers



2. Memory management is fully delegatable

- Supports hierarchical system design
- Enabled by *capability-based access control*



3. “Incremental consistency” design pattern

- Fast transitions between consistent states
- Restartable operations with progress guarantee



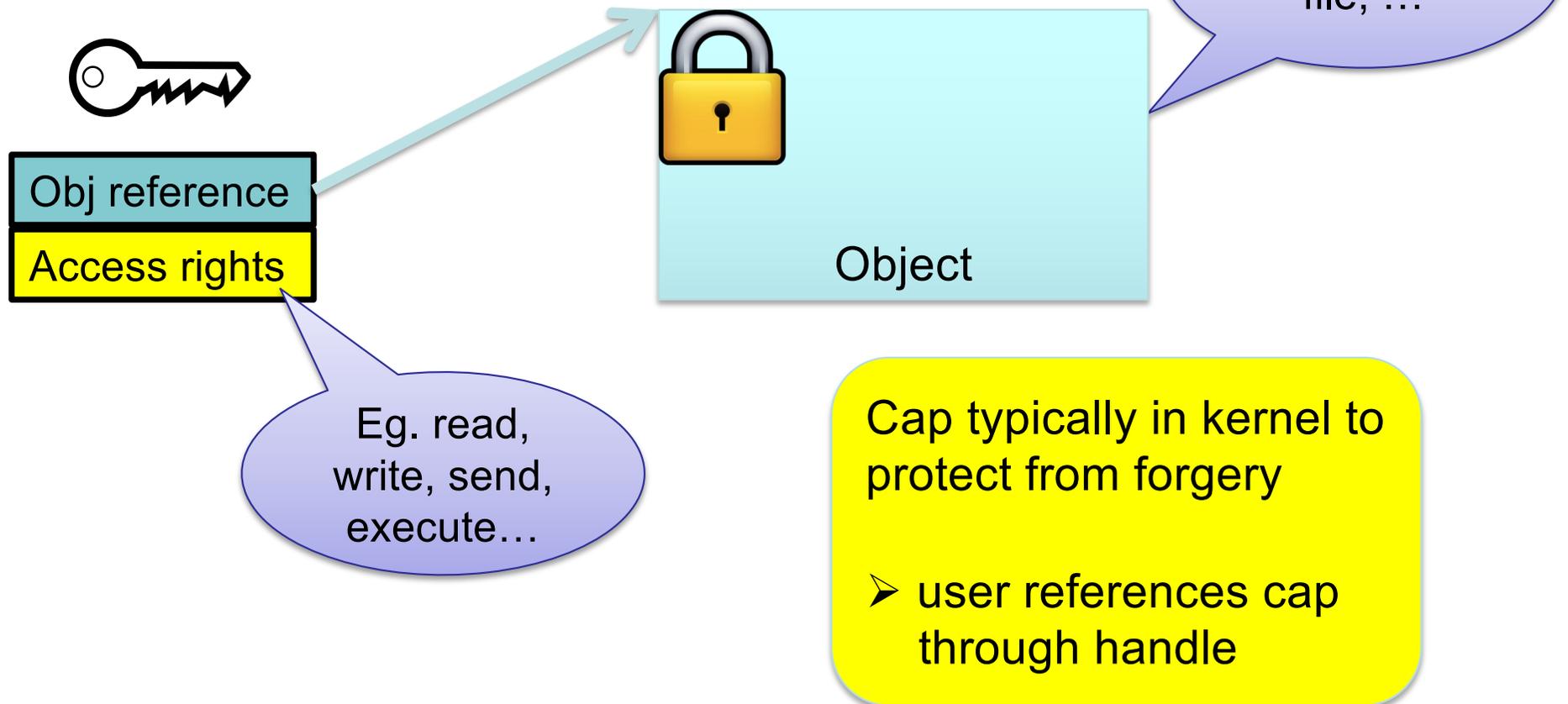
4. No concurrency in the kernel

- Interrupts never enabled in kernel
- Interruption points to bound latencies
- Clustered multikernel design for multicores

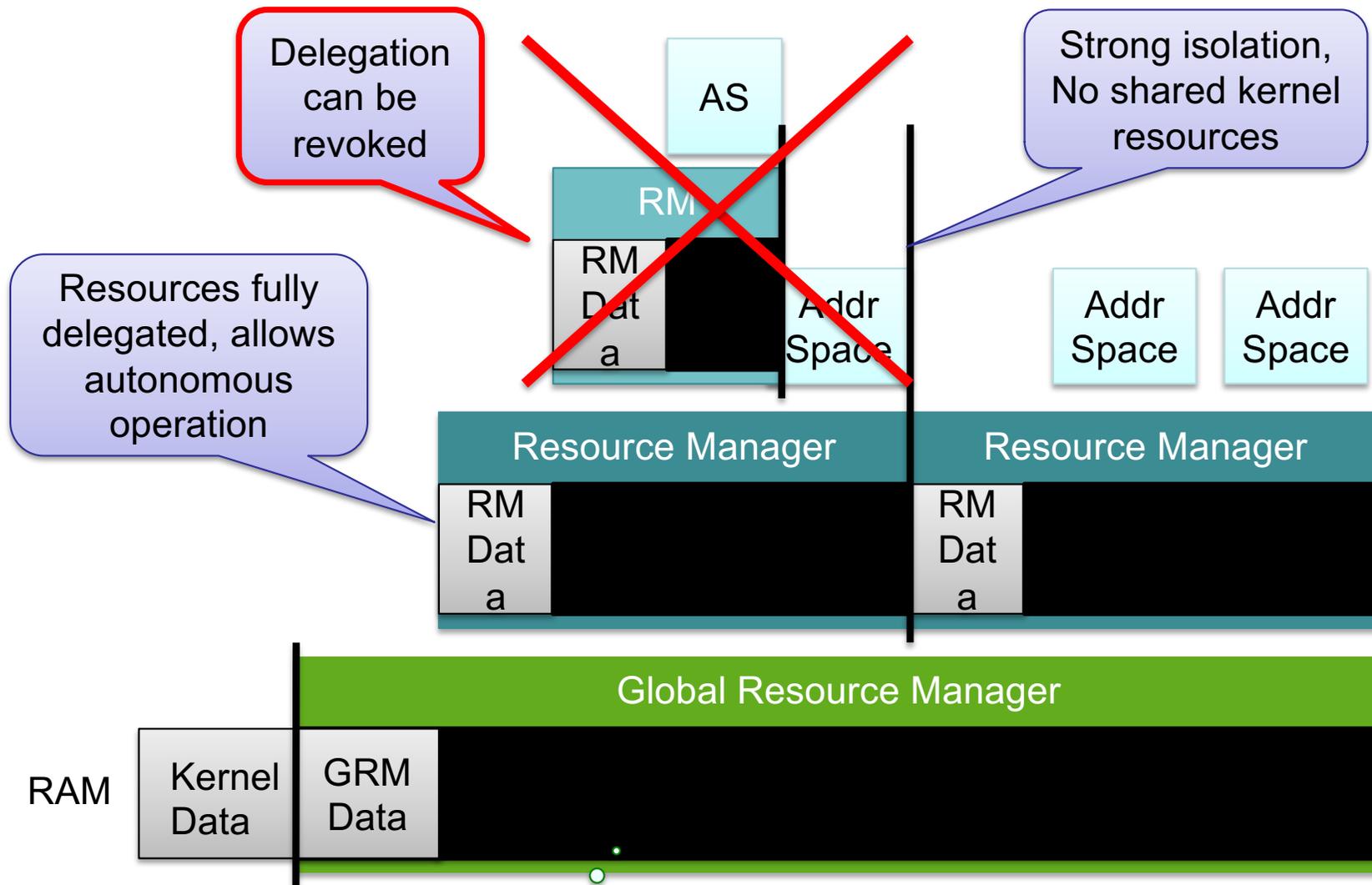


# What are Capabilities?

**Cap = Access Token**

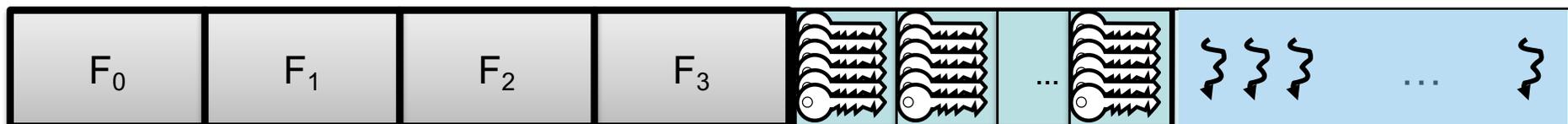
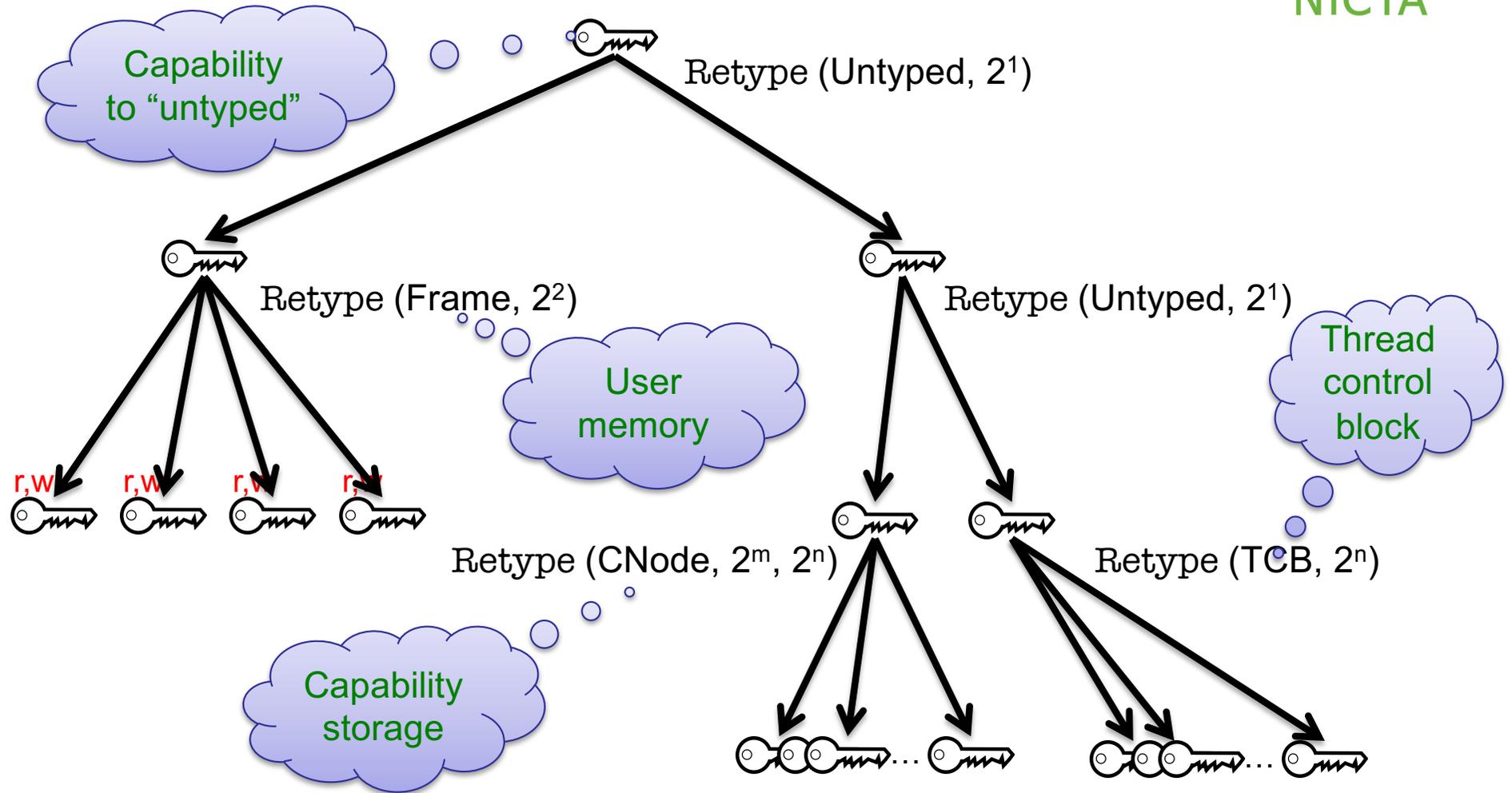


# seL4 User-Level Memory Management



**“Untyped” (unallocated) memory**

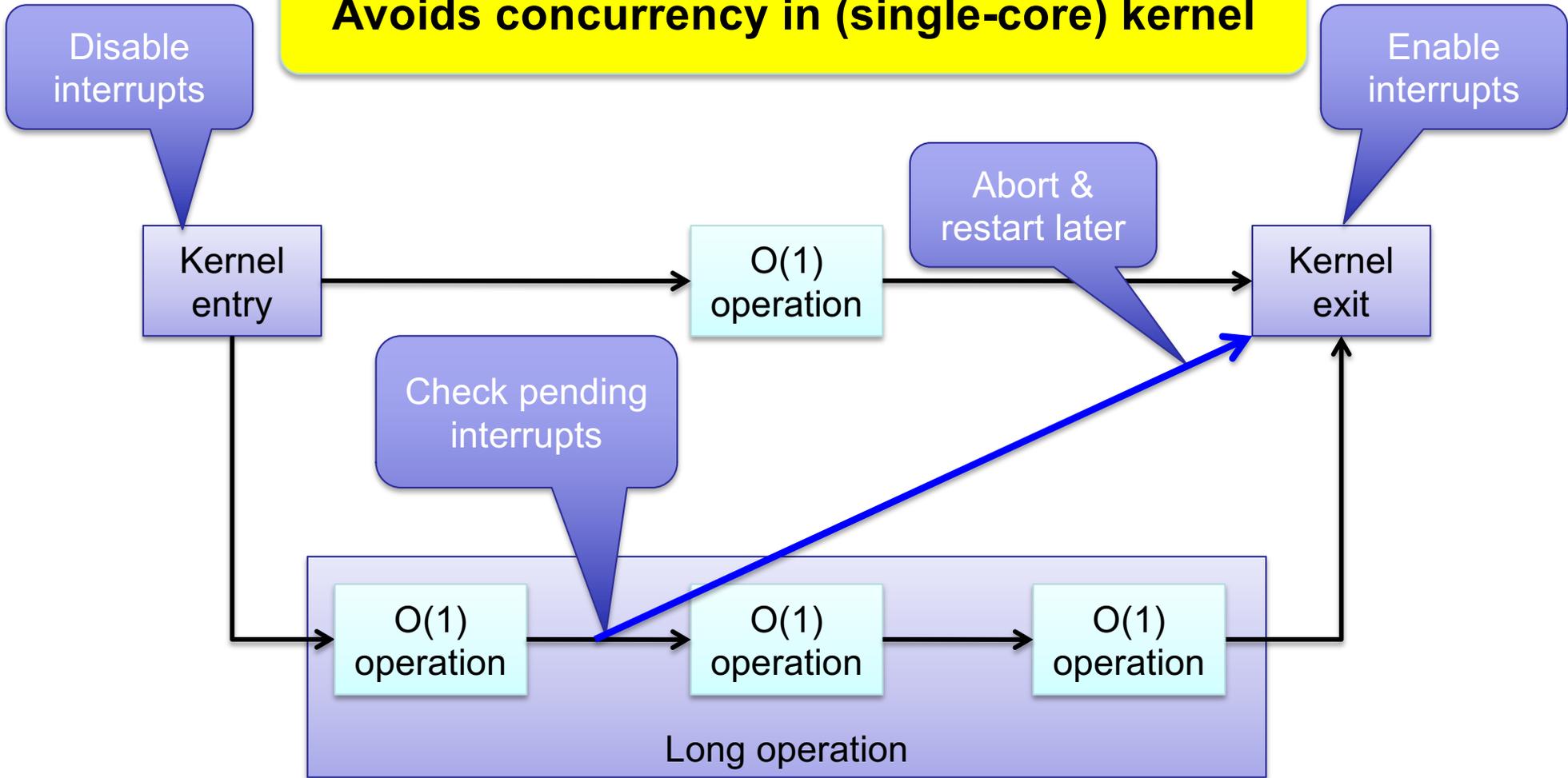
# seL4 Memory Management Mechanics: Retype



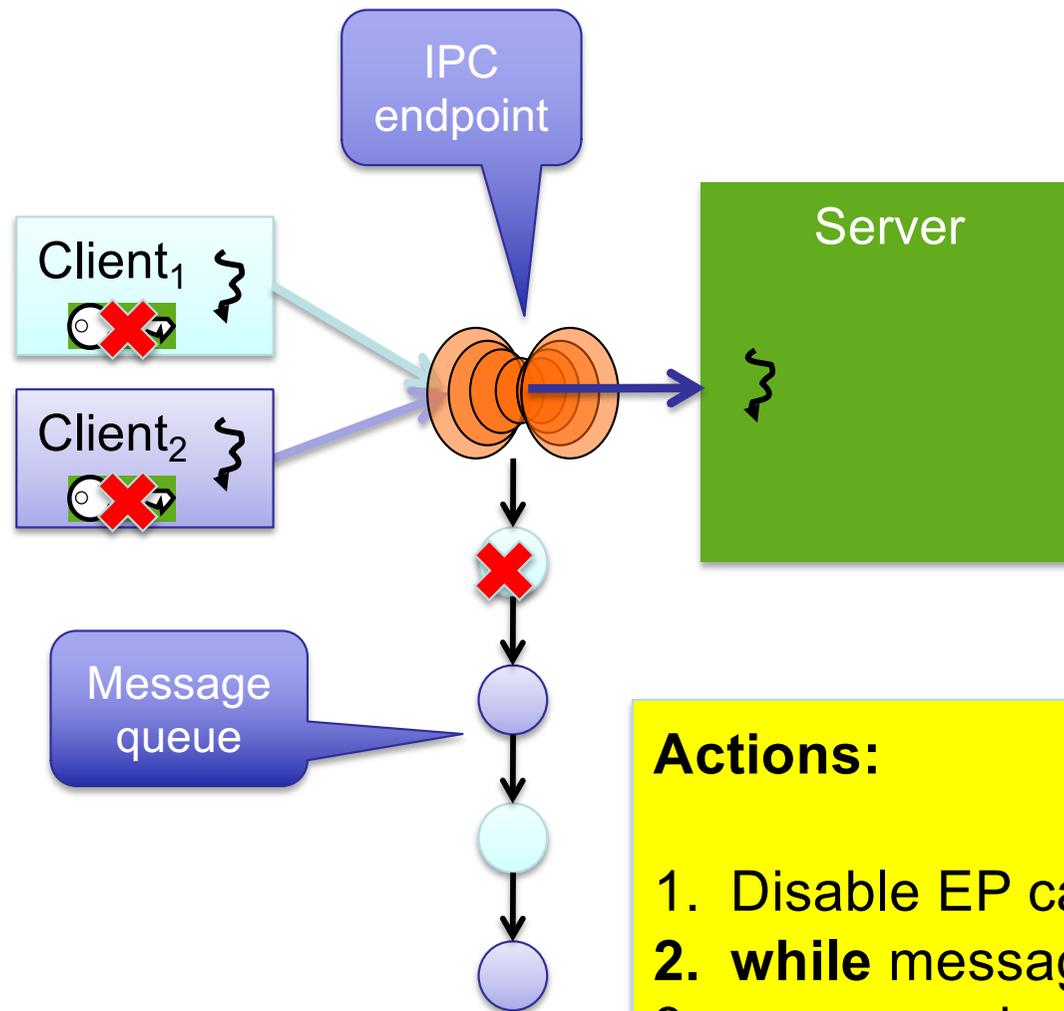
# Incremental Consistency



**Avoids concurrency in (single-core) kernel**



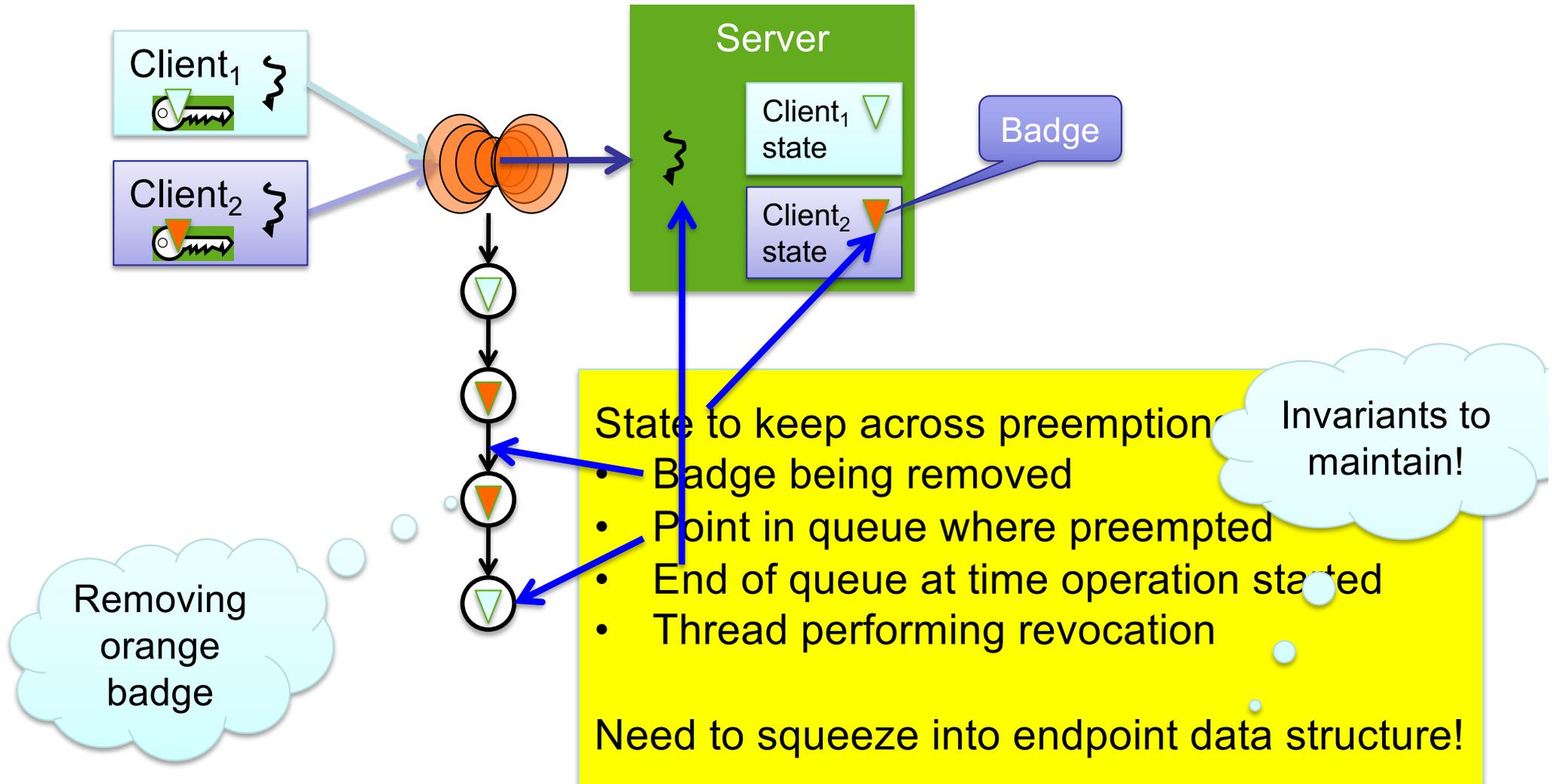
# Example: Destroying IPC Endpoint



## Actions:

1. Disable EP cap (prevent new messages)
2. **while** message queue not empty **do**
3.     remove head of queue (abort message)
4.     check for pending interrupts
5. **done**

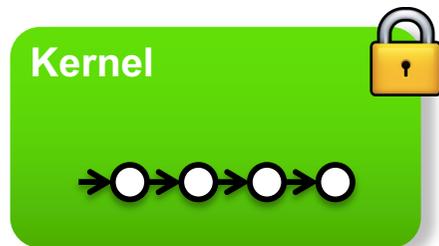
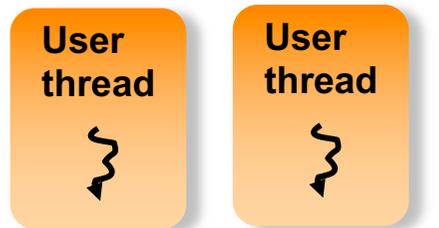
# Difficult Example: Revoking IPC “Badge”



# Approaches for Multicore Kernels

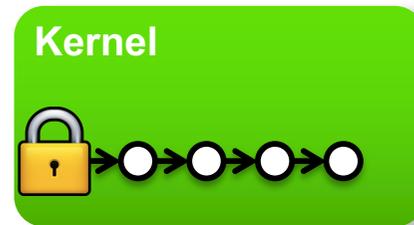
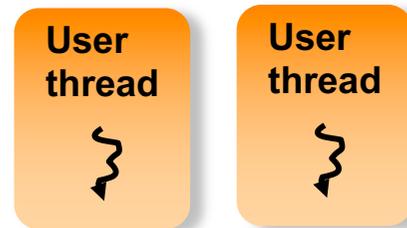


**SMP  
big lock**



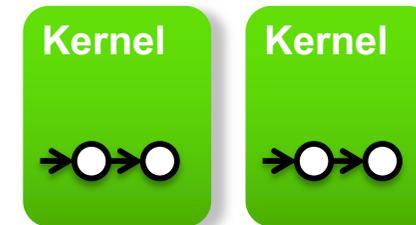
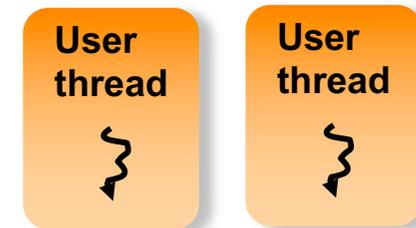
**Core**   **Core**

**SMP  
fine-grained locks**



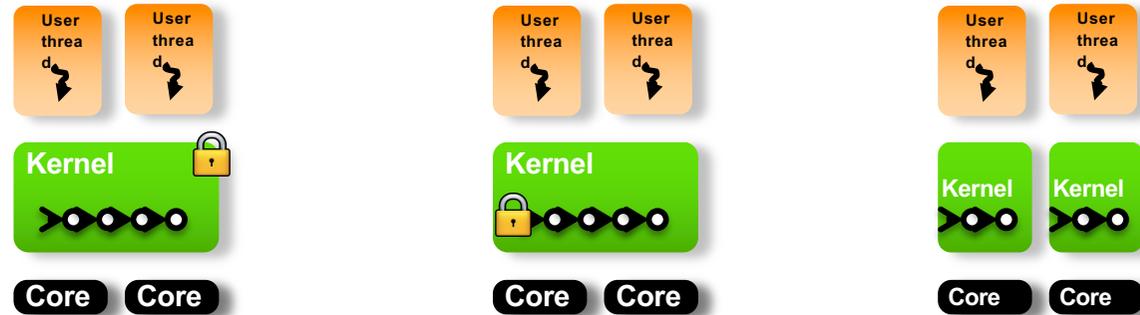
**Core**   **Core**

**Multikernel  
no locks**



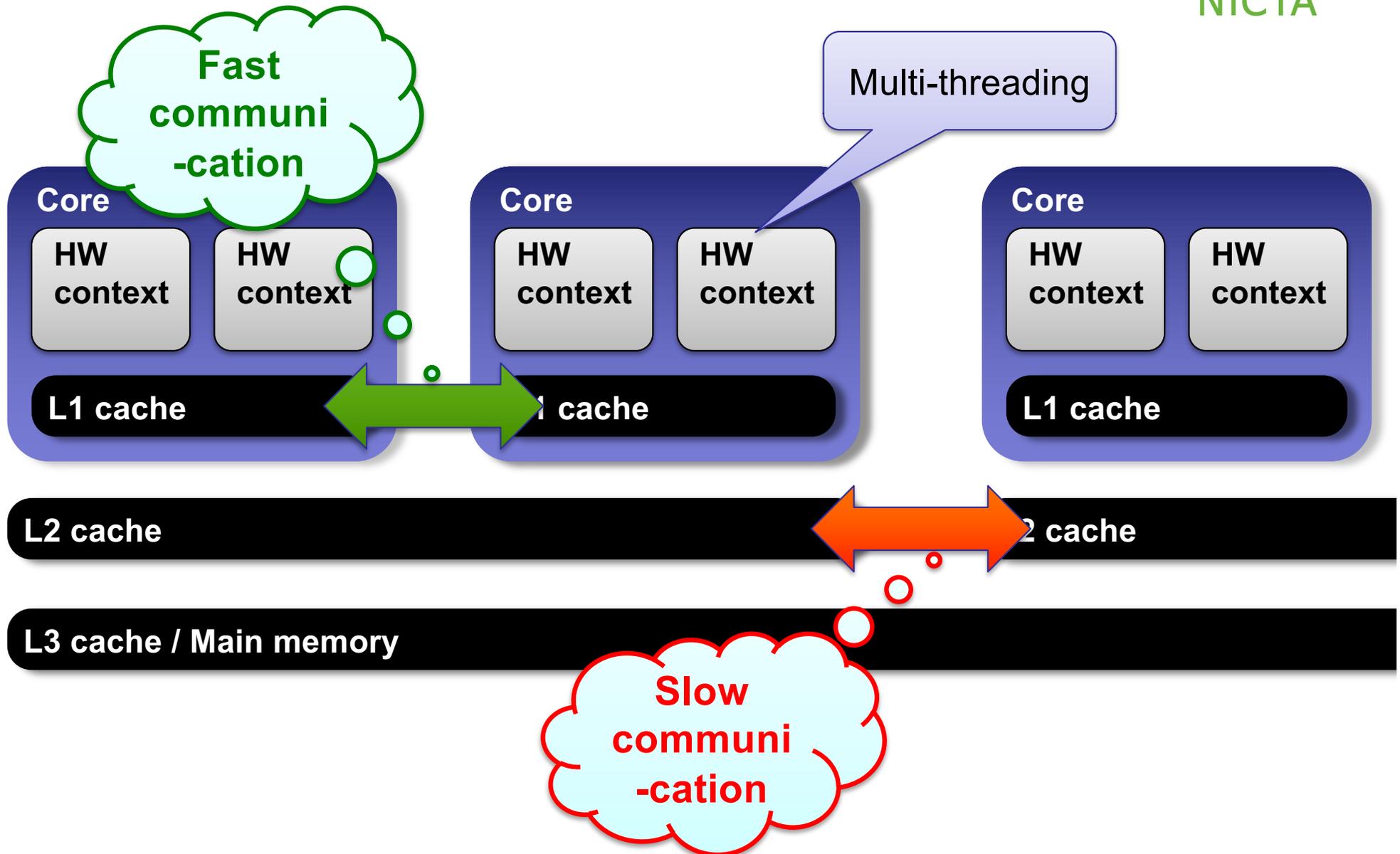
**Core**   **Core**

# Multicore Kernel Trade-Offs



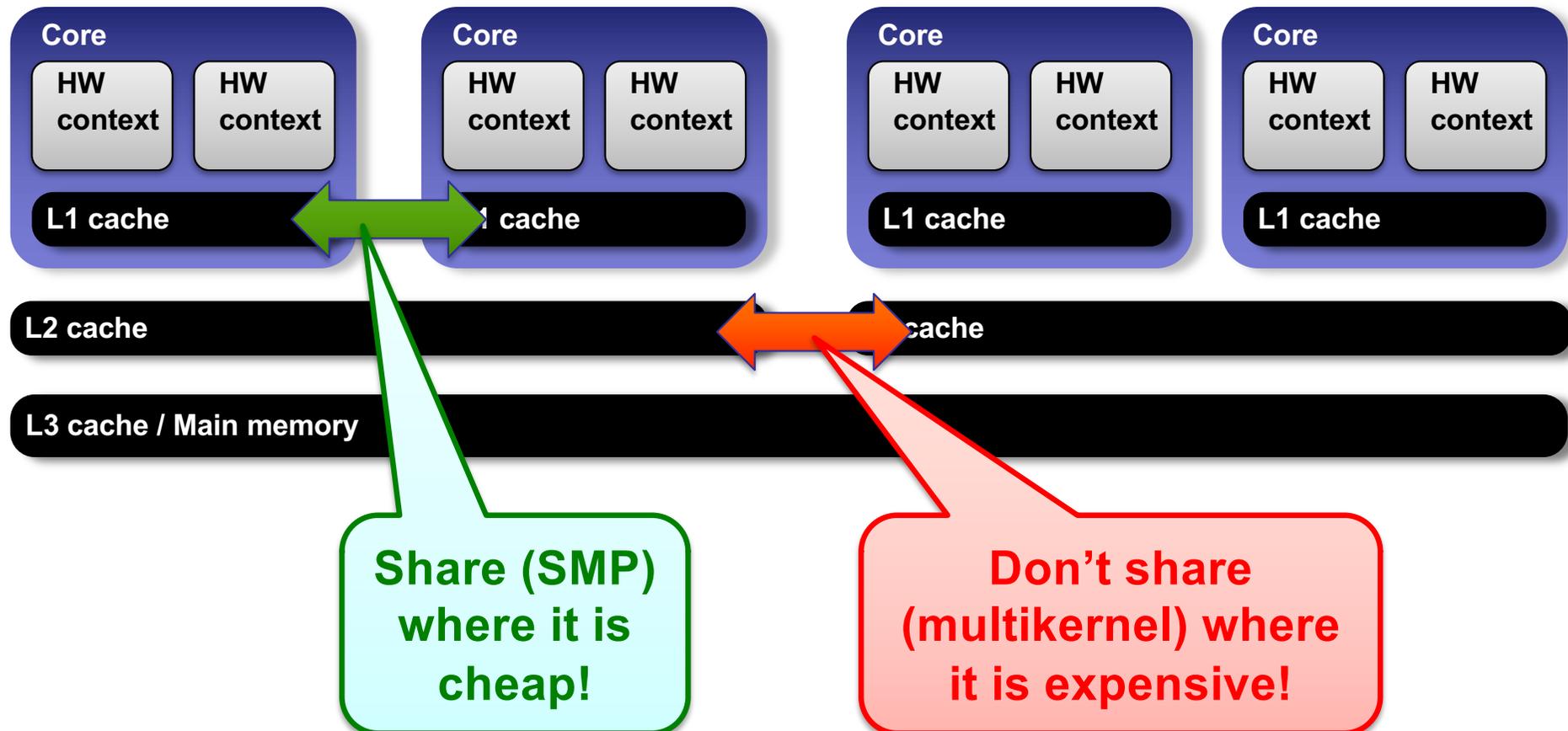
Property	Big Lock	Fine-grained Locking	Multikernel
Data structures	shared	shared	distributed
Scalability	poor	good	excellent
Concurrency in kernel	zero	high	zero
Kernel complexity	low	high	low
Resource management	centralised	centralised	distributed

# Reality of Multicore is NUMA!

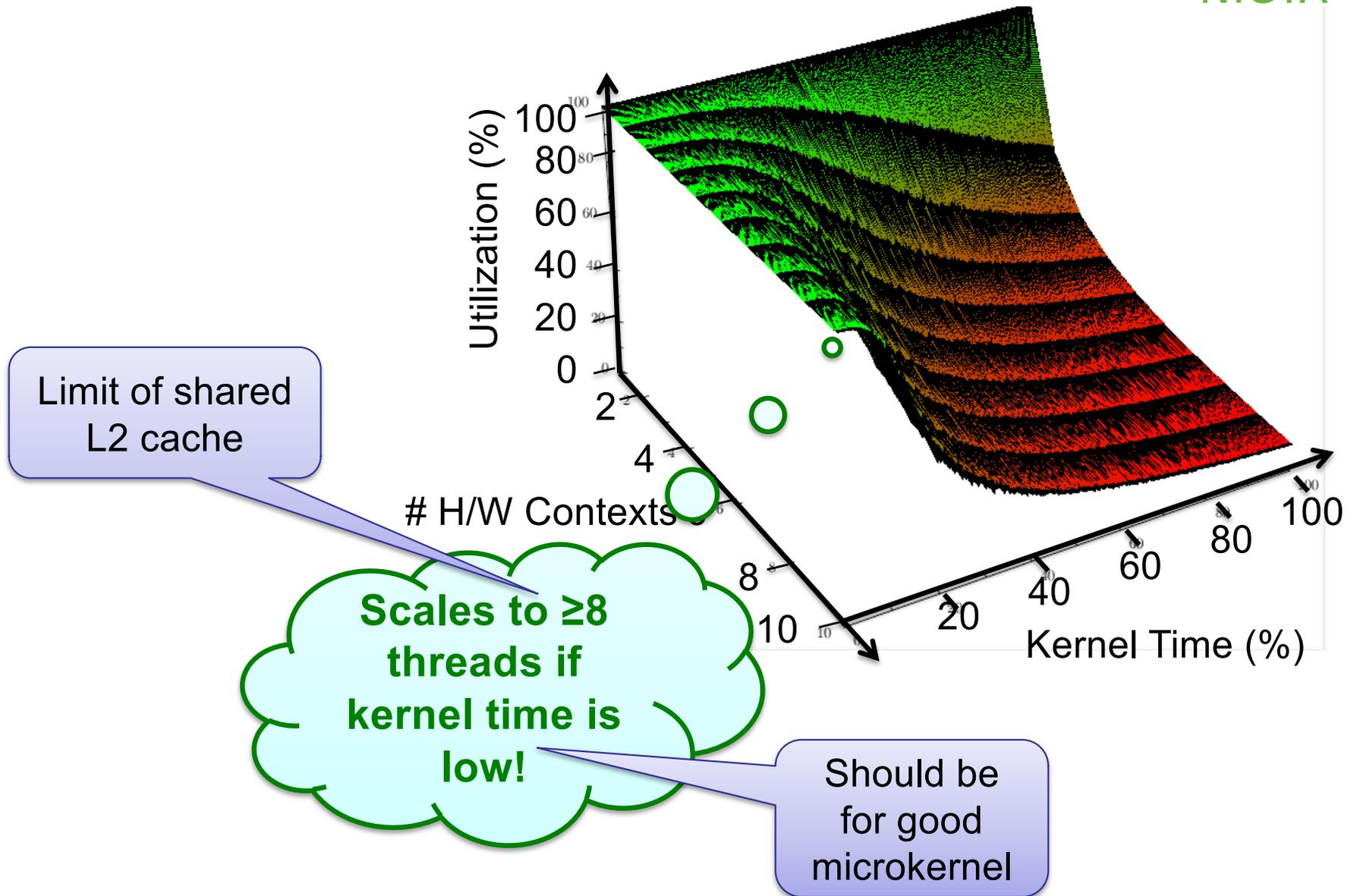


# Microkernel Principle: Policy Freedom

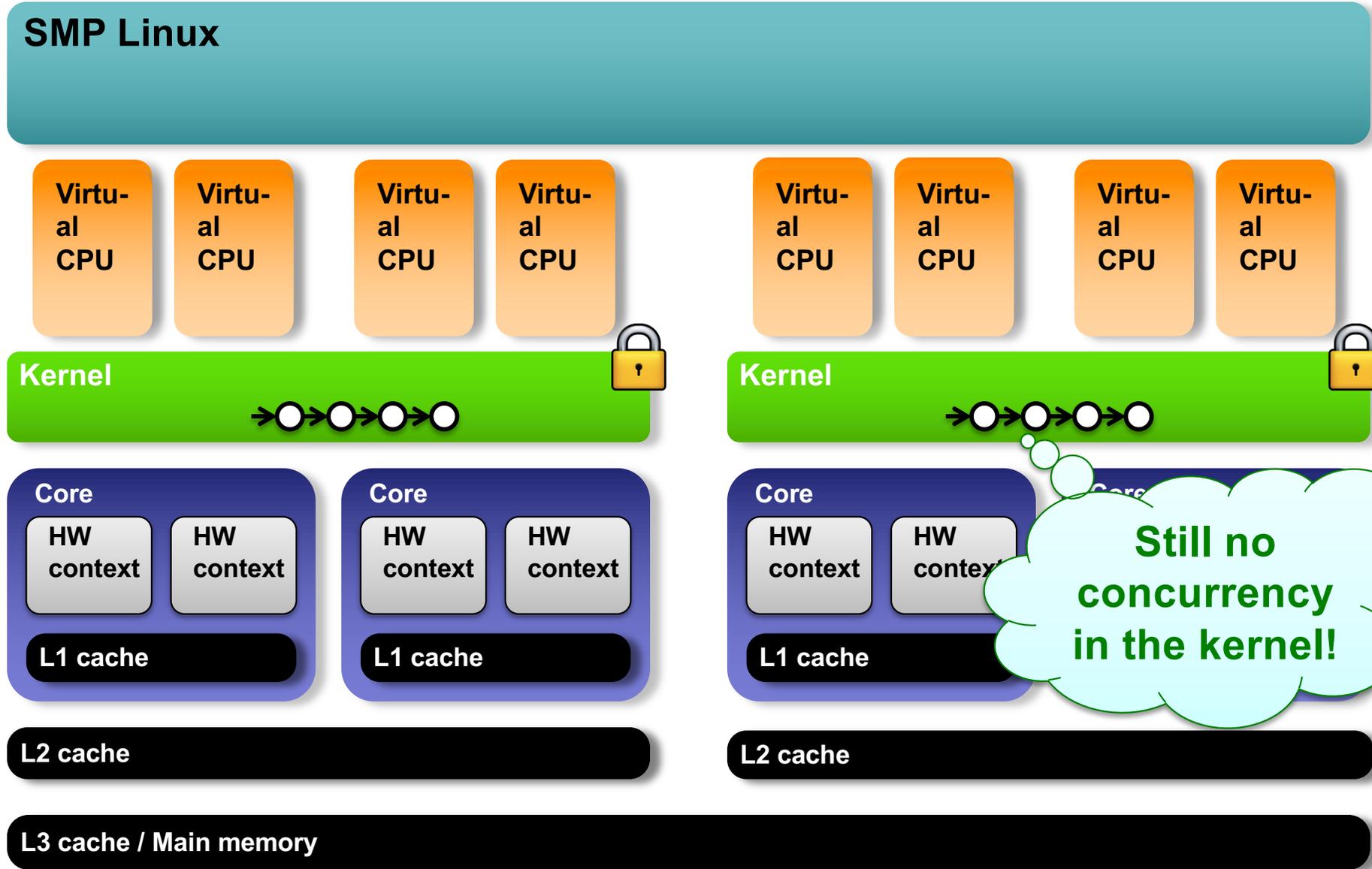
- Kernel must not dictate policy
- Kernel must not introduce avoidable overhead



# Performance of Big Kernel Lock



# Resulting Design: Clustered Multikernel



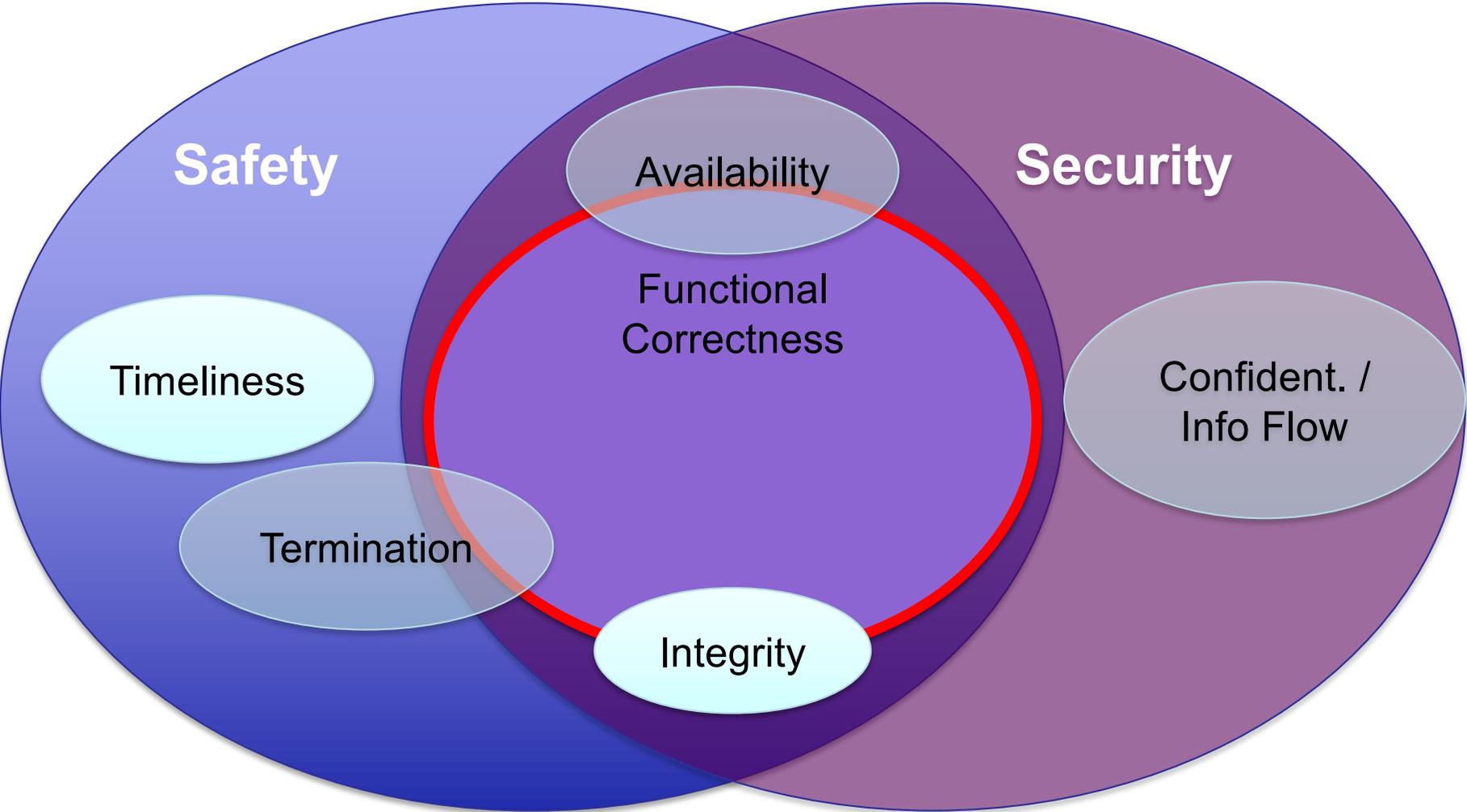
# Agenda

---

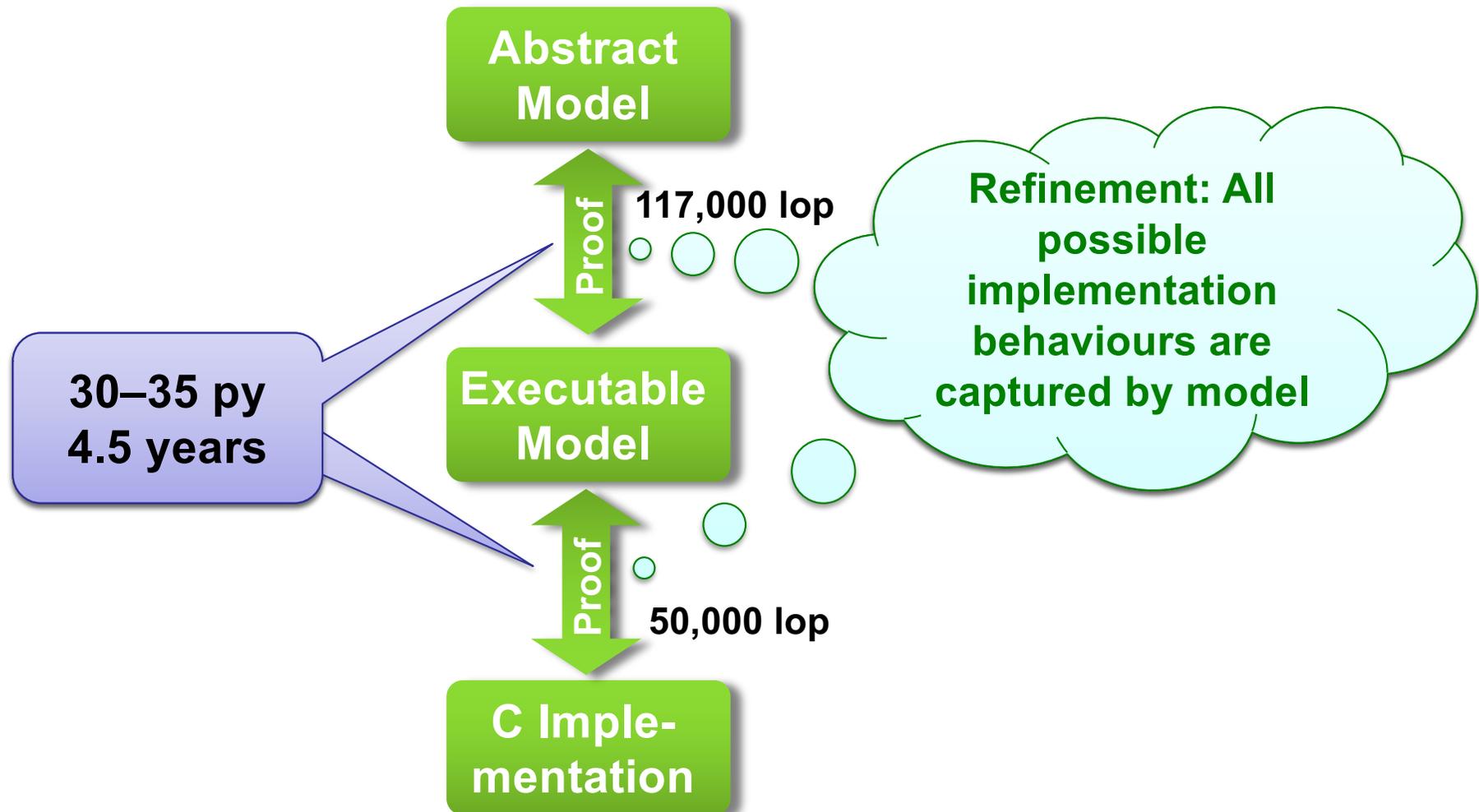


- Motivation
- What is a microkernel, and what is L4?
- seL4 – designed for trustworthiness
- **Establishing trustworthiness**
- From kernel to system
- Sample system 1: Secure access controller
- Sample system 2: RapiLog

# seL4 as Basis for Trustworthy Systems



# Proving Functional Correctness



# Correctness



```
datatype
  rights = Read
        | Write
        | Grant
        | Create
```

```
record cap =
  entity :: entity_id
  rights :: rights
```

```
record constdefs
  schedule :: "unit s_monad"
type "schedule ≡ do"
```

```
lemma isolation:
  "[sane s;
   s' ∈ execute cmds s;
   isEntityOf s es;
   isEntityOf s e;
   entity c = e;
   c :=> subSysCaps s es]
  ⇒ c :=> subSysCaps s' es"
```

```
schedule :: Kernel ()
schedule = do
  action ← getSchedulerAction
```

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

  tptr->tcbPriority = prio;
}

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

Specification

(L)

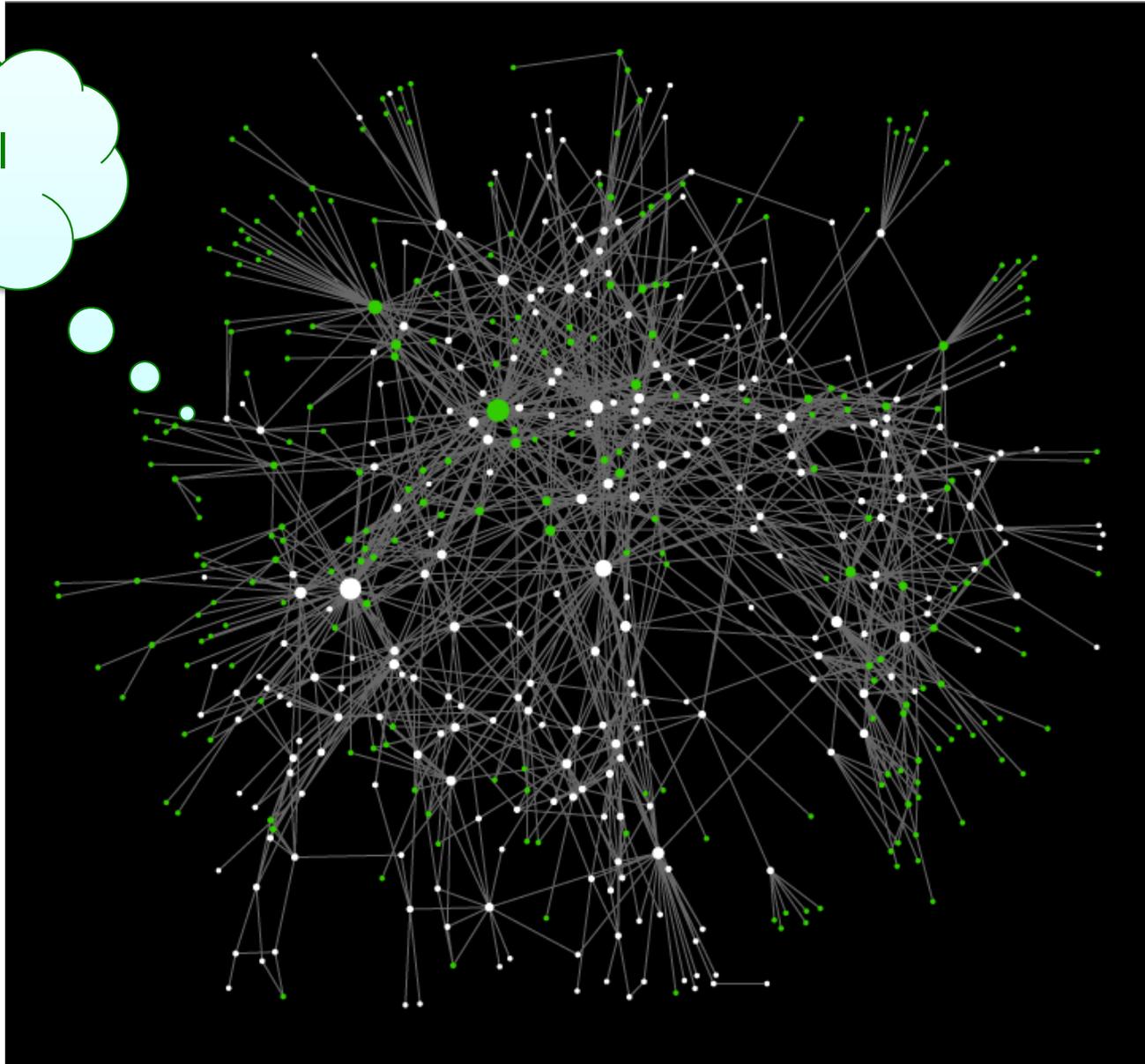
Thread

5,700

Performance Implementation  
(C/asm)  
Model

# Why So Long for 9,000 LOC?

seL4 call graph



# Costs Breakdown



Haskell design	2 py
C implementation	2 weeks
Debugging/Testing	2 months
Kernel verification	12 py
Formal frameworks	10 py
<b>Total</b>	<b>25 py</b>
Repeat (estimated)	6 py
Traditional engineering	4–6 py

## Did you find bugs???

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

# seL4 Formal Verification Summary

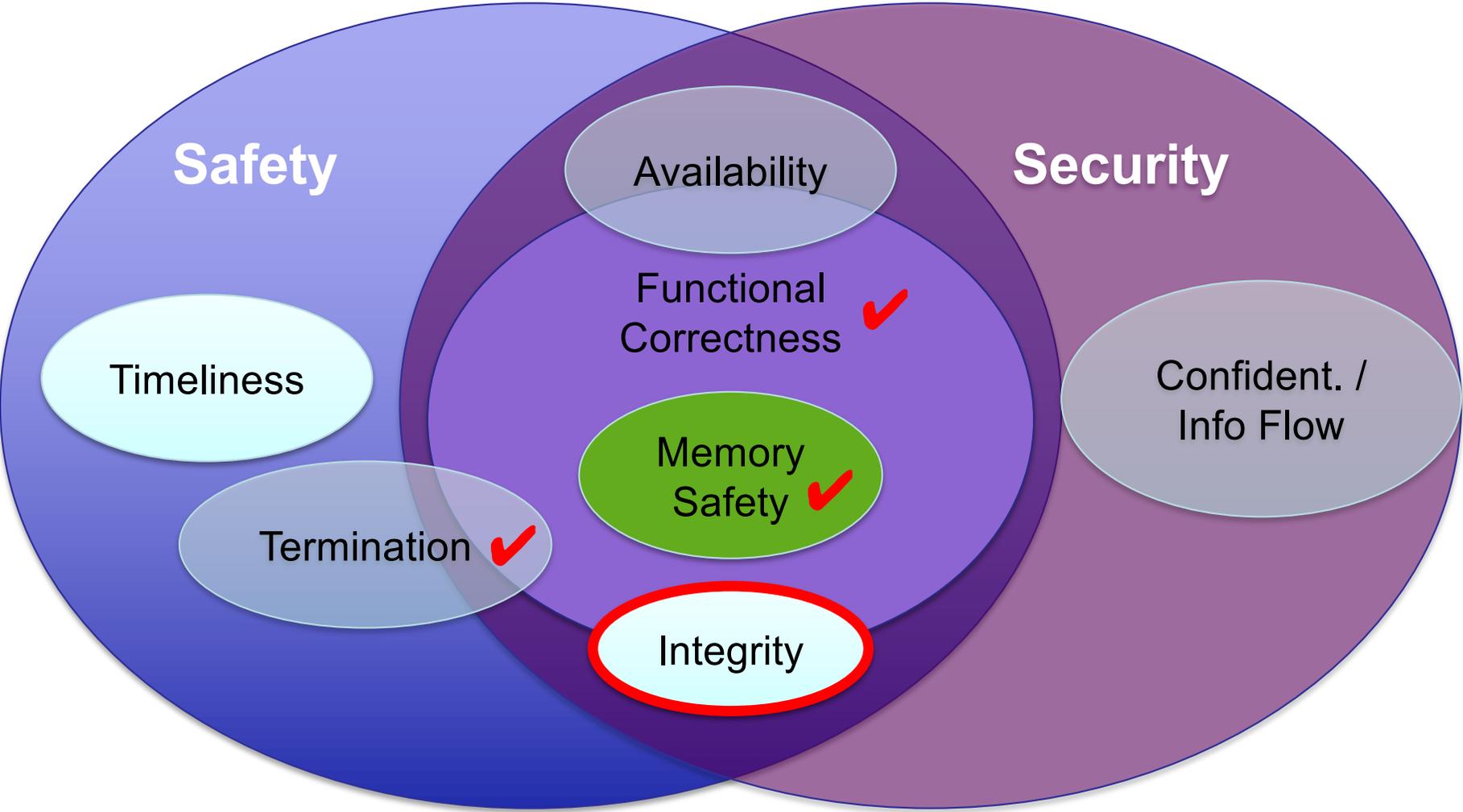
---

## Kinds of properties proved

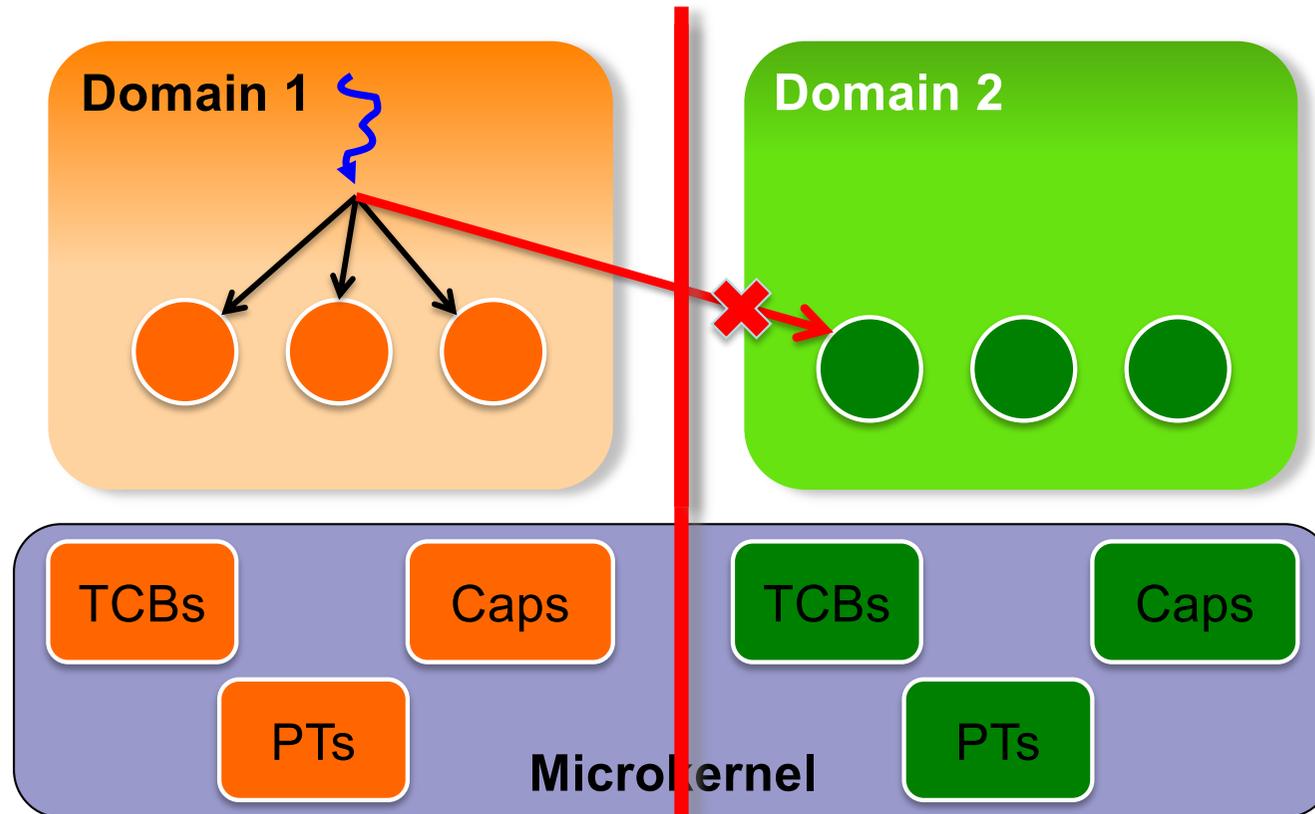
- Behaviour of C code is fully captured by abstract model
- Behaviour of C code is fully captured by executable model
- Kernel never fails, behaviour is always well-defined
  - assertions never fail
  - will never de-reference null pointer
  - cannot be subverted by malformed input
- All syscalls terminate,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

Can prove further properties on abstract level!

# seL4 as Basis for Trustworthy Systems



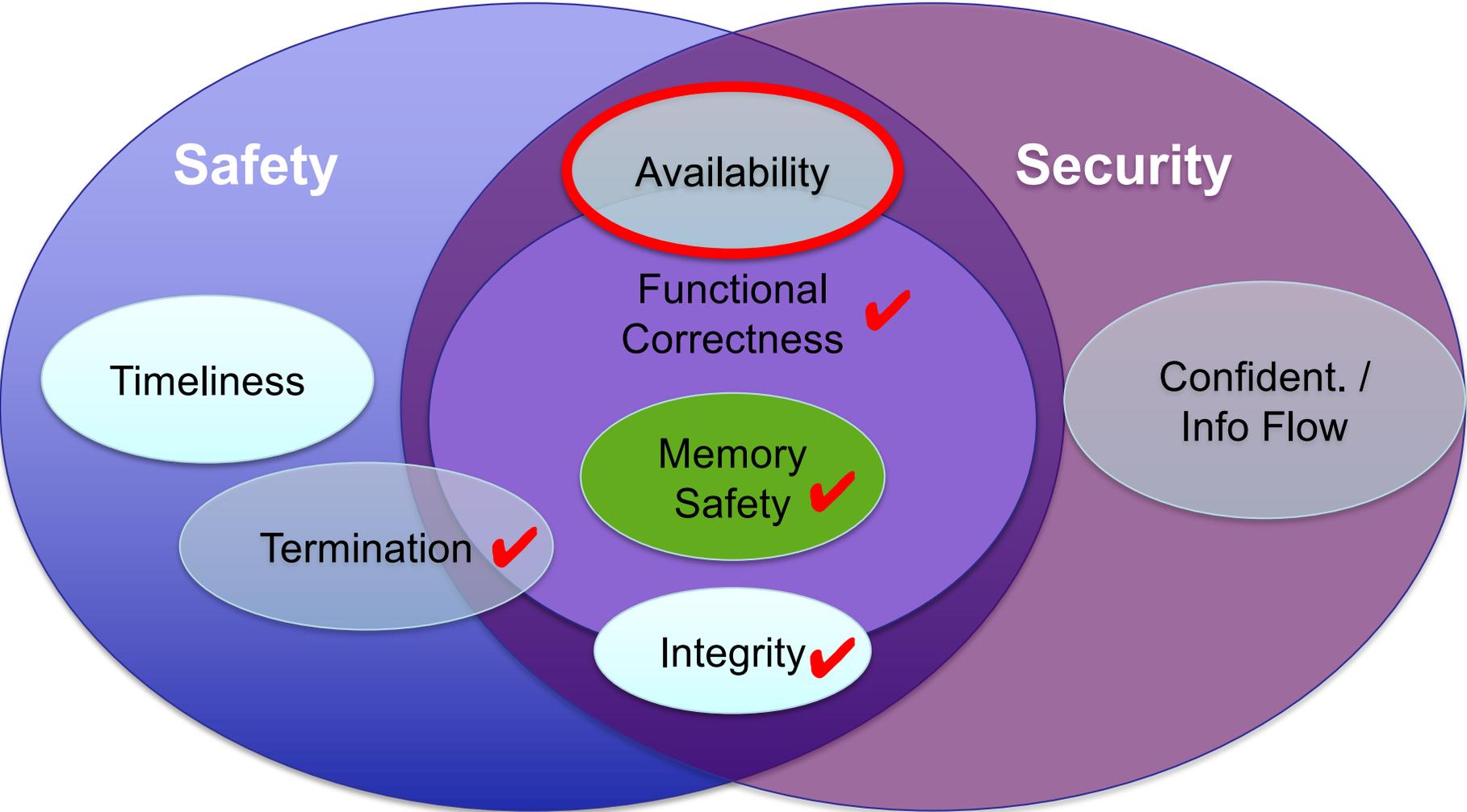
# Integrity: Limiting Write Access



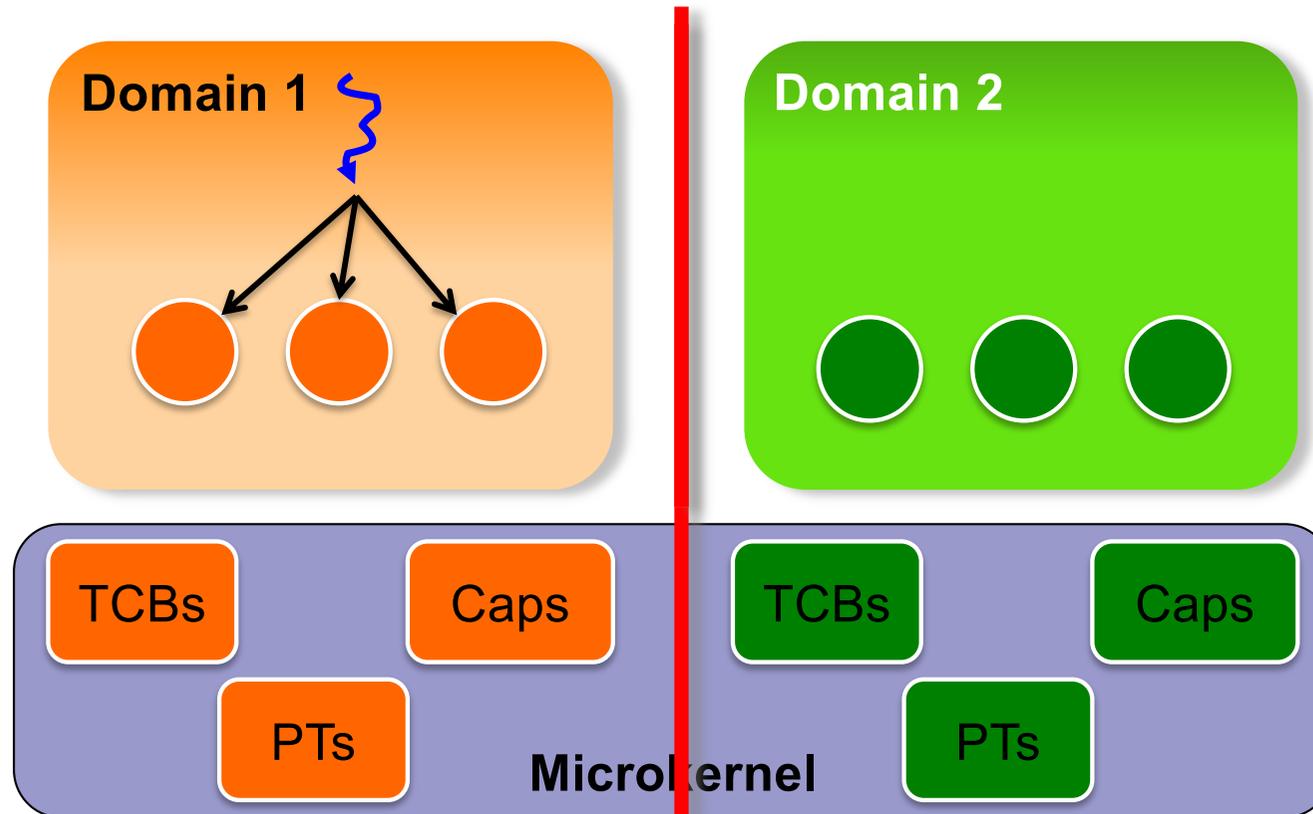
## To prove:

- Domain-1 doesn't have write *capabilities* to Domain-2 objects  
⇒ no action of Domain-1 agents will modify Domain-2 state
- Specifically, *kernel does not modify on Domain-1's behalf!*
  - Prove kernel only allows write upon capability presentation

# seL4 as Basis for Trustworthy Systems

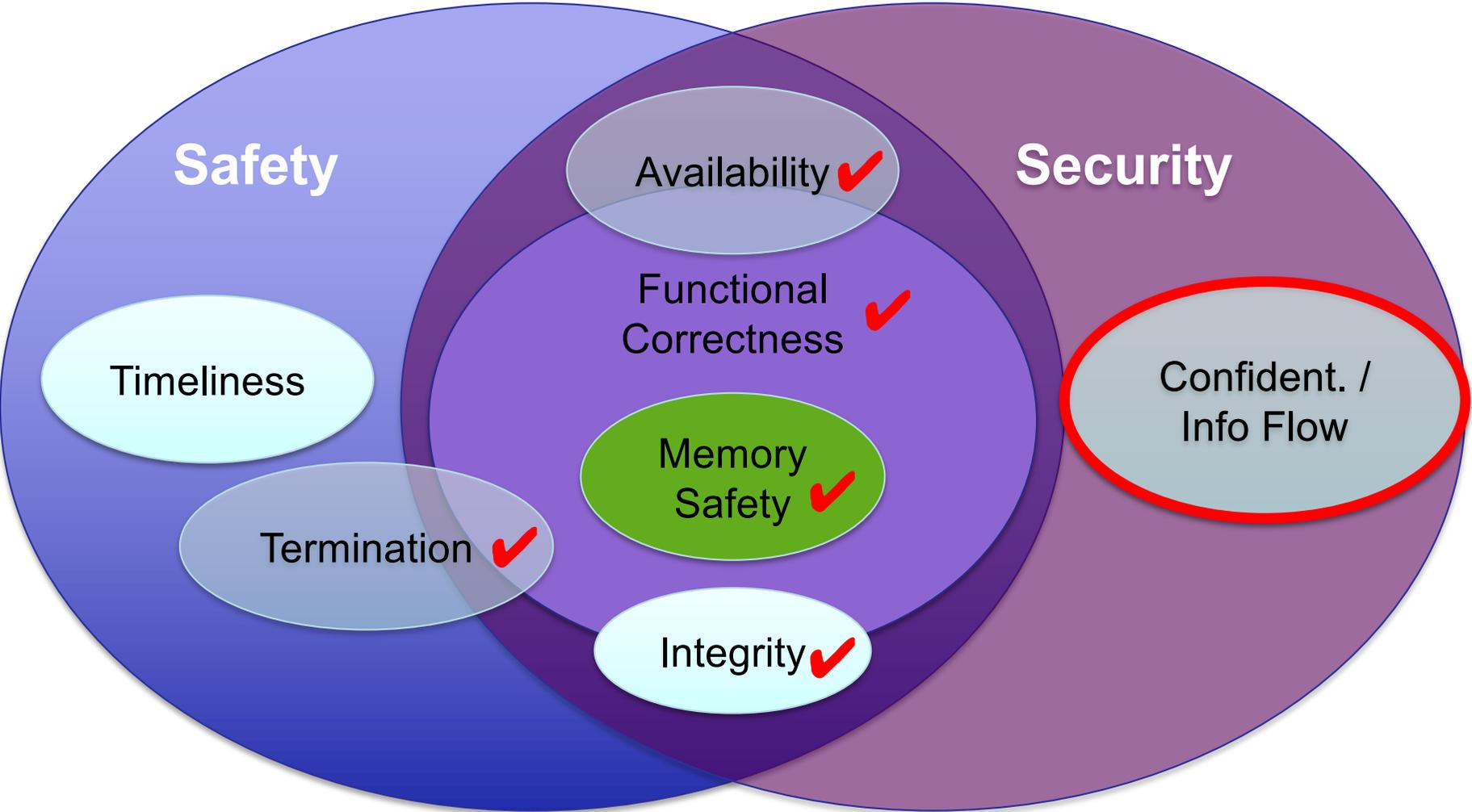


# Availability: Ensuring Resource Access

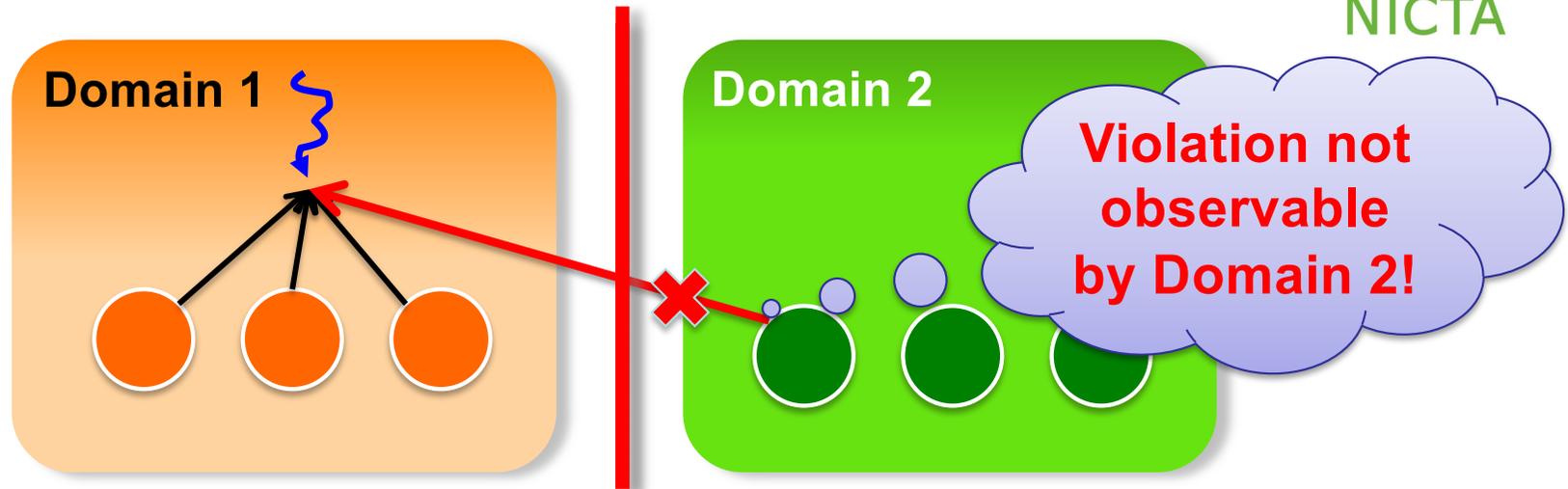


- Strict separation of kernel resources  
⇒ agent cannot deny access to another domain's resources

# seL4 as Basis for Trustworthy Systems



# Confidentiality: Limiting Read Accesses



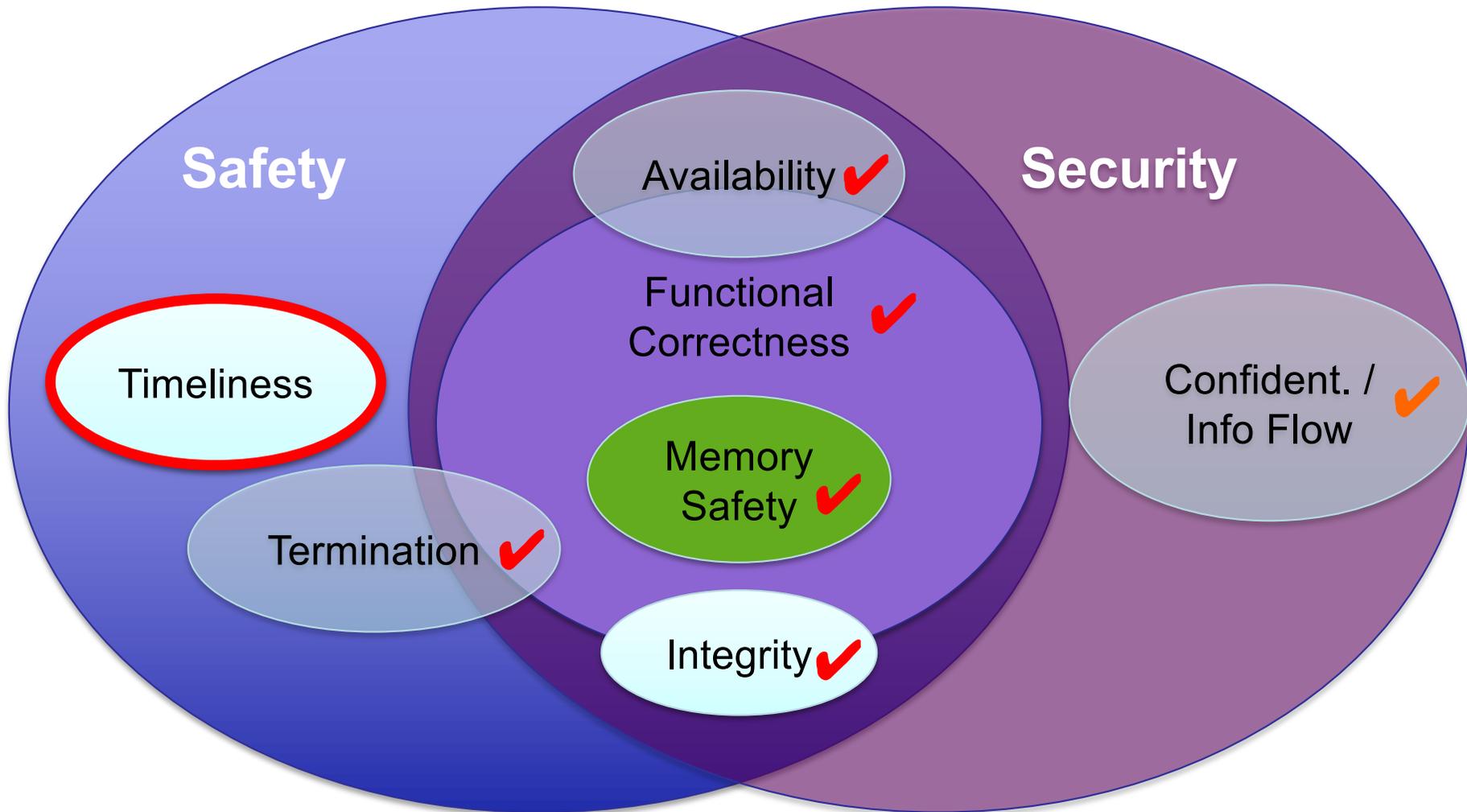
## To prove:

- Domain-1 doesn't have read capabilities to Domain-2 objects  
⇒ no action of any agents will reveal Domain-2 state to Domain-1

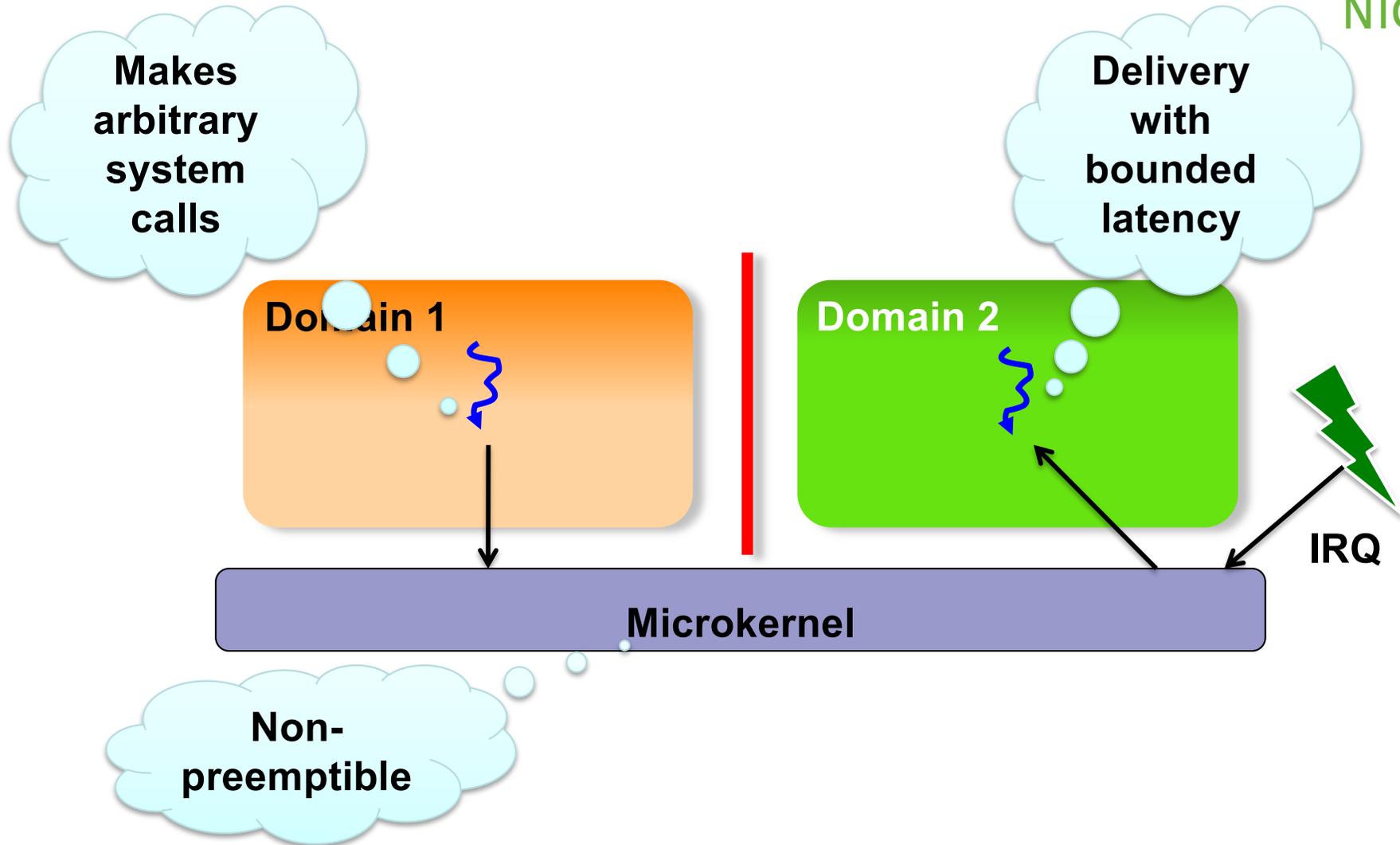
## Non-interference proof in progress:

- Evolution of Domain 1 does not depend on Domain-2 state
- Presently cover only overt information flow

# seL4 as Basis for Trustworthy Systems

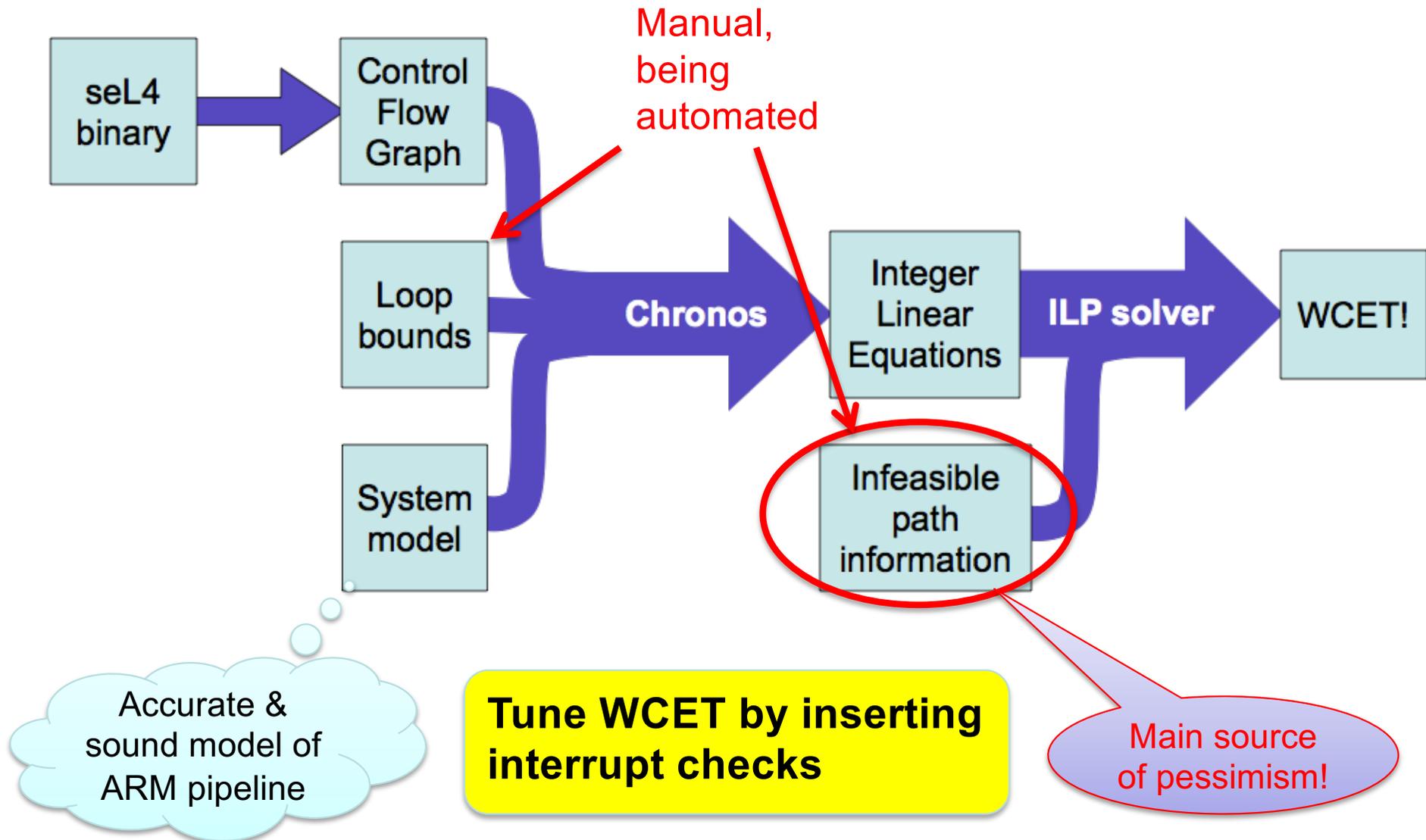


# Timeliness

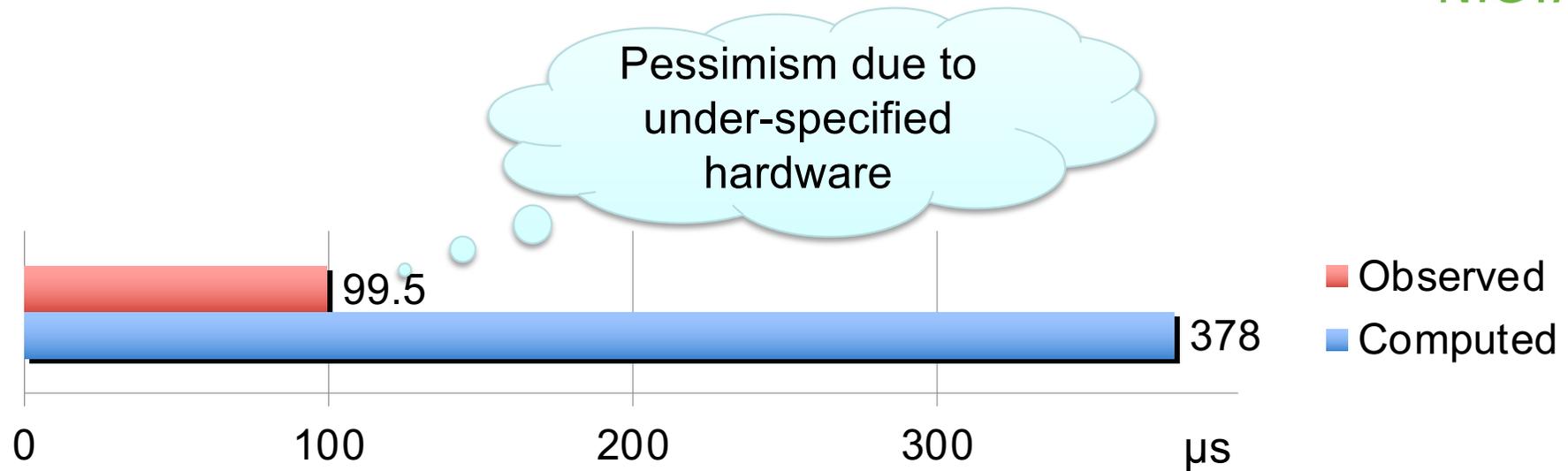


**Need worst-case execution time (WCET) analysis of kernel**

# WCET Analysis Approach



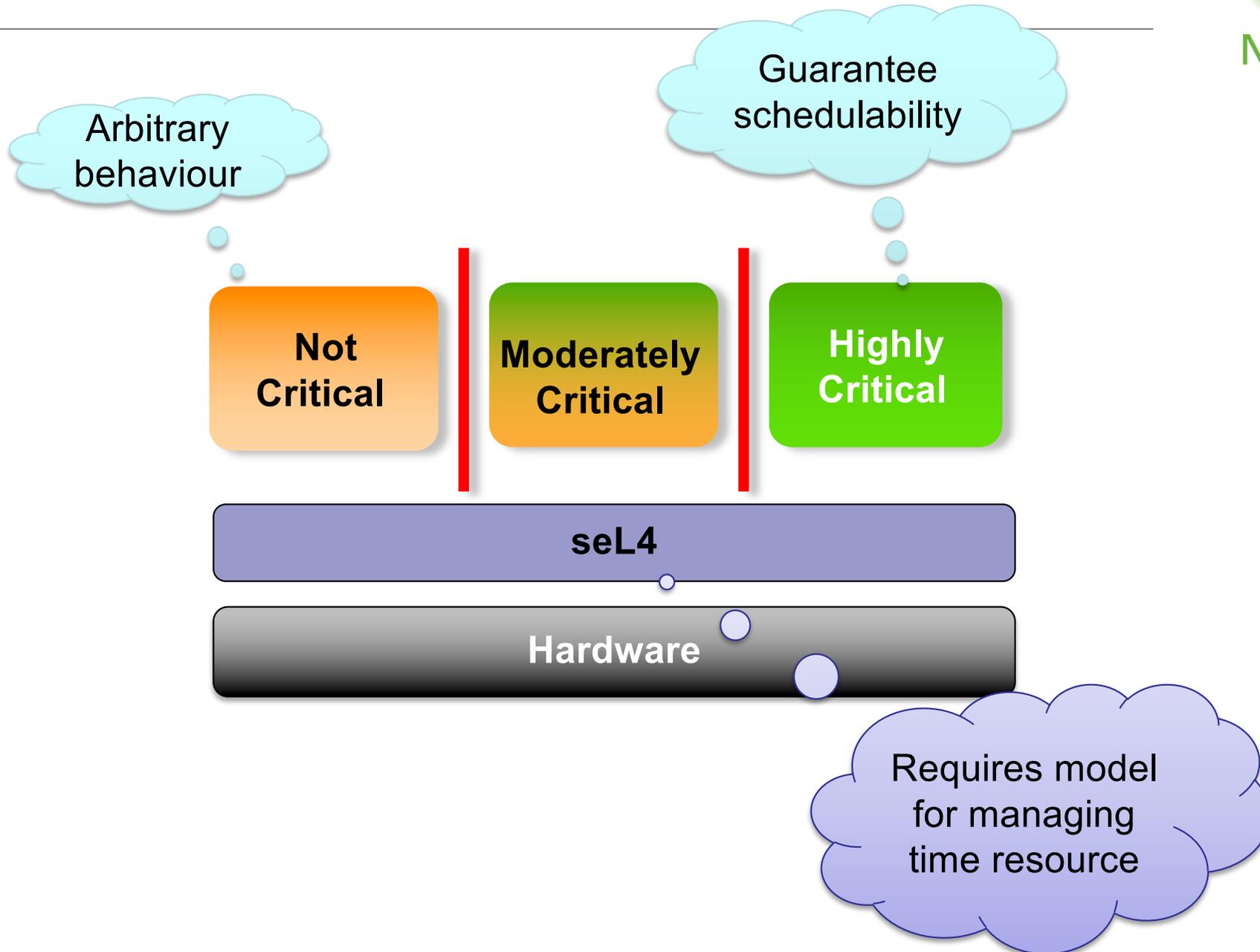
# Result



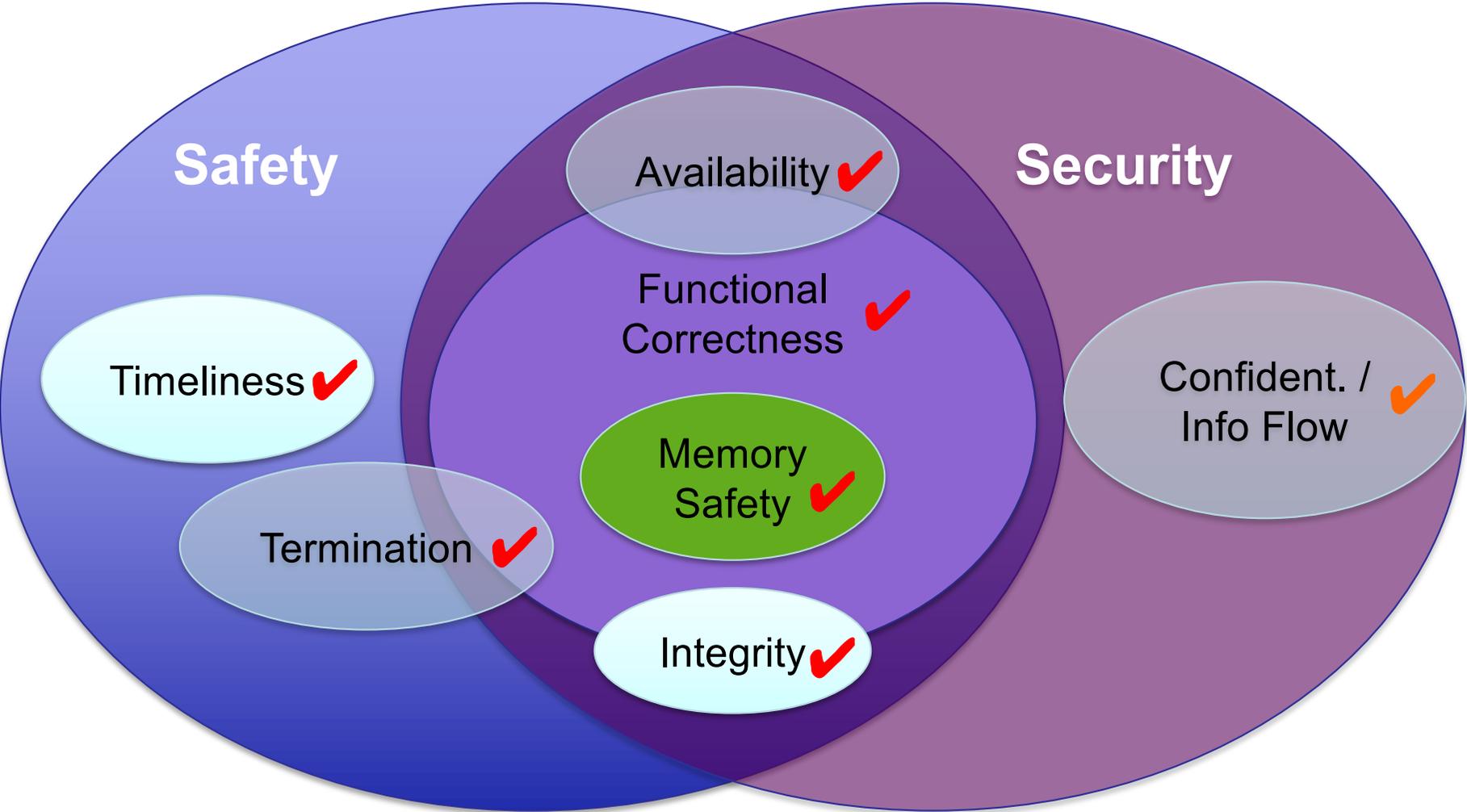
**WCET presently limited by verification practicalities**

- 10 μs seem achievable

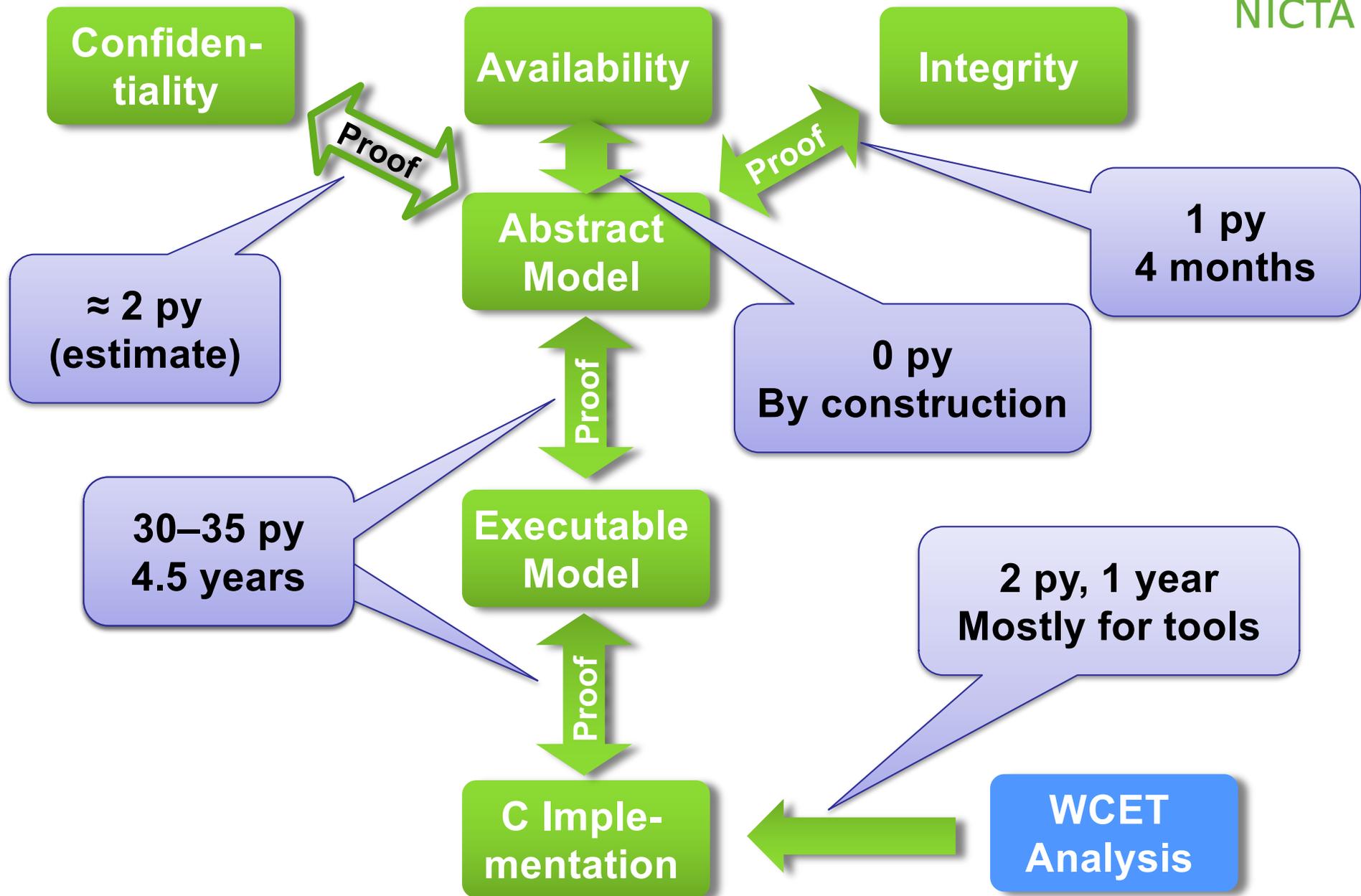
# Future: Whole-System Schedulability



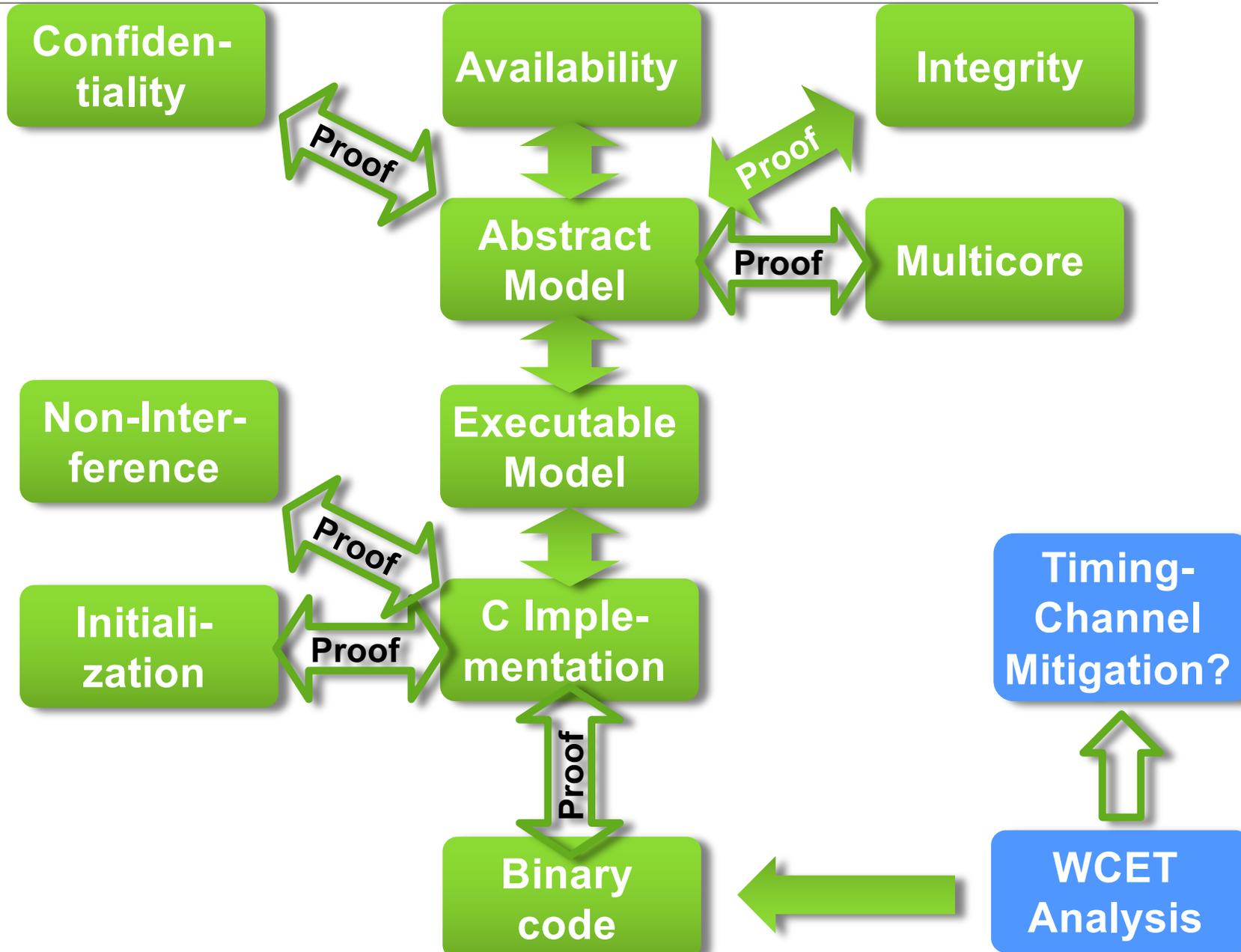
# seL4 as Basis for Trustworthy Systems



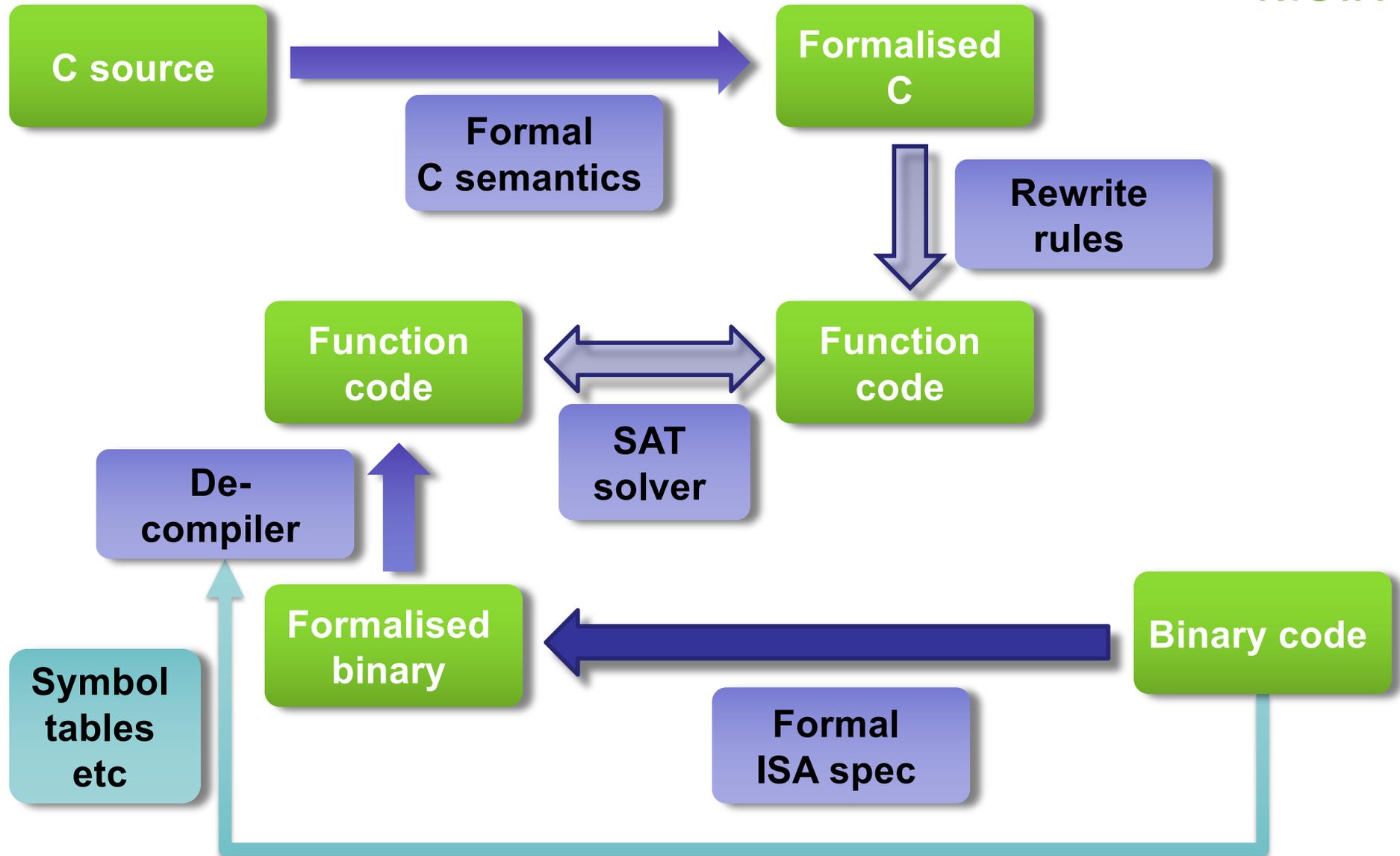
# Proving seL4 Trustworthiness



# seL4 – the Next 24 Months



# Binary Code Verification (In Progress)



# Multikernel Verification

- By definition, multikernel images execute independently
  - except for explicit messaging



- To prove:
  - isolated images are initialised correctly
  - images maintain isolation at run time

Essentially non-interference

# Agenda

---



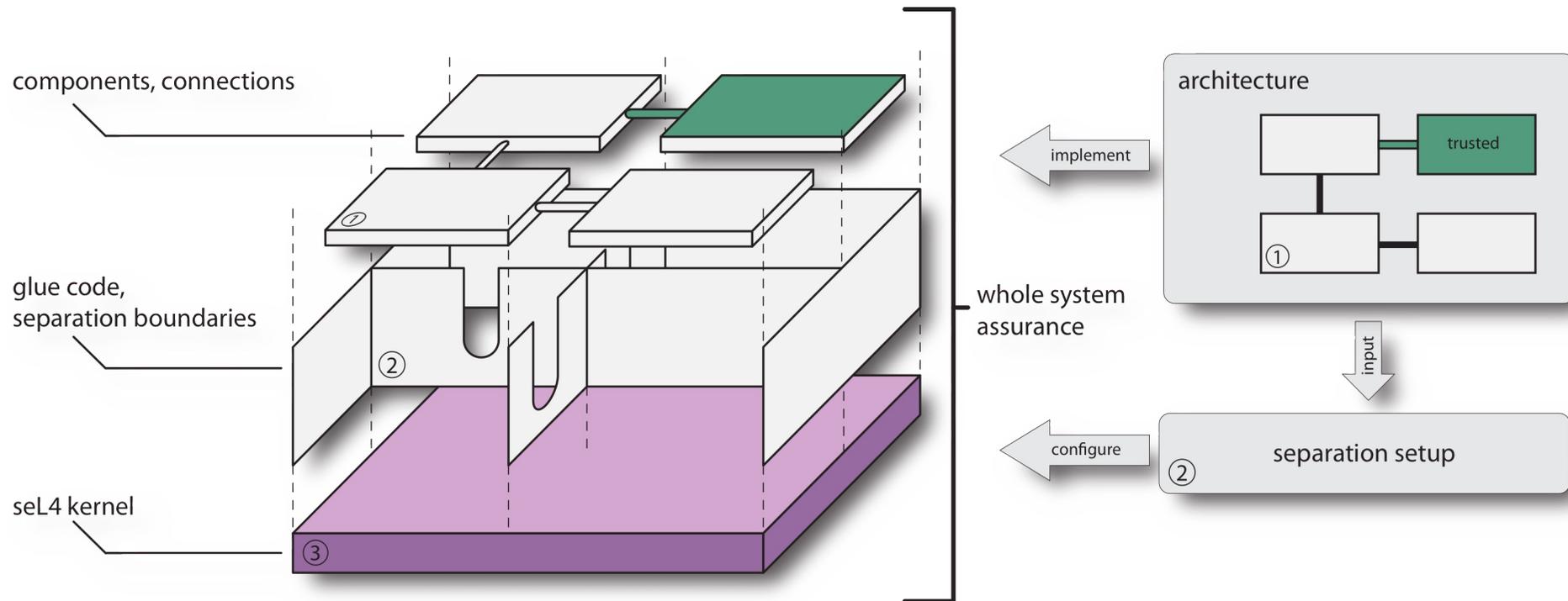
- Motivation
- What is a microkernel, and what is L4?
- seL4 – designed for trustworthiness
- Establishing trustworthiness
- **From kernel to system**
- Sample system 1: Secure access controller
- Sample system 2: RapiLog

# Phase Two: Full-System Guarantees

- Achieved: Verification of microkernel (8,700 LOC)
- Next step: Guarantees for real-world systems (1,000,000 LOC)

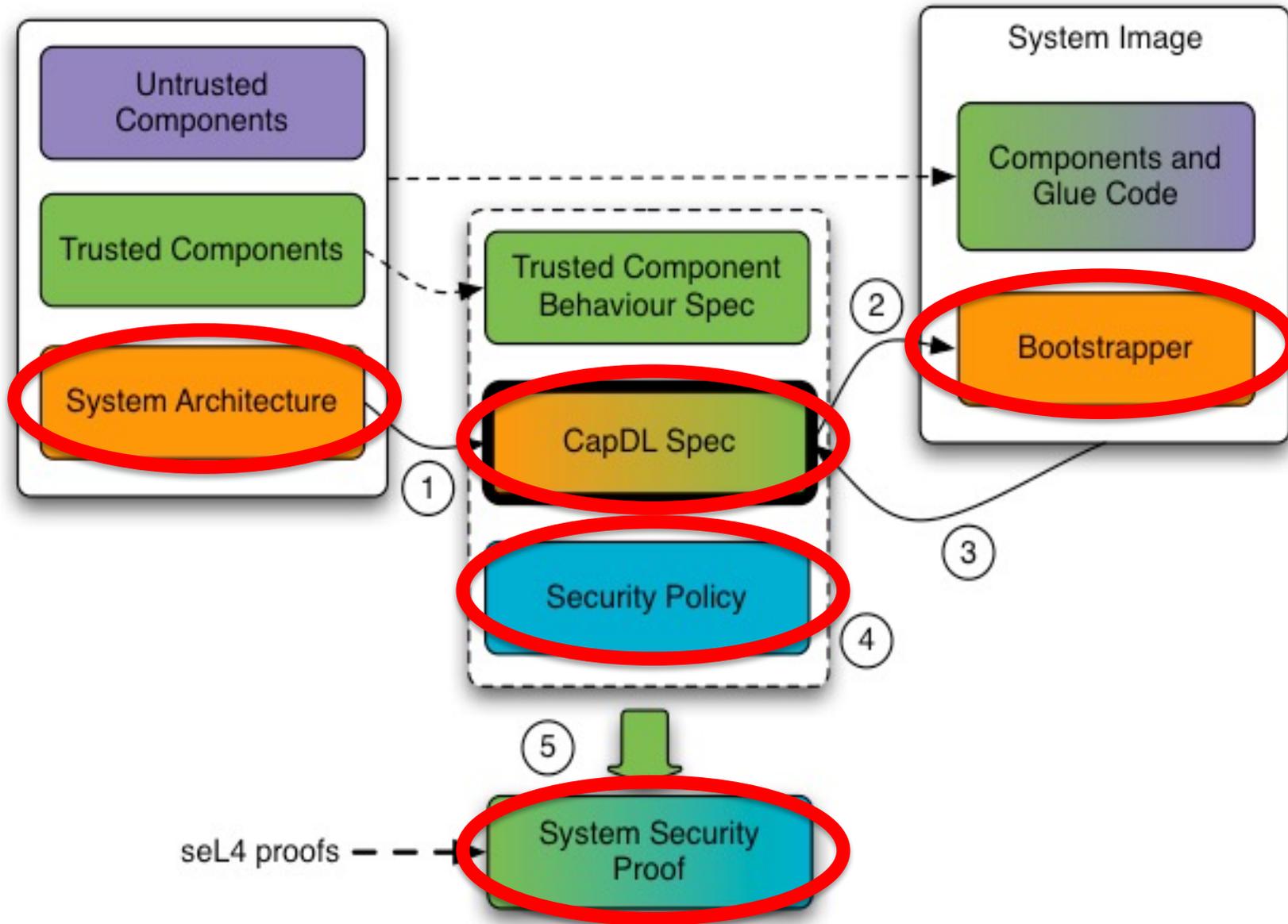


# Overview of Approach

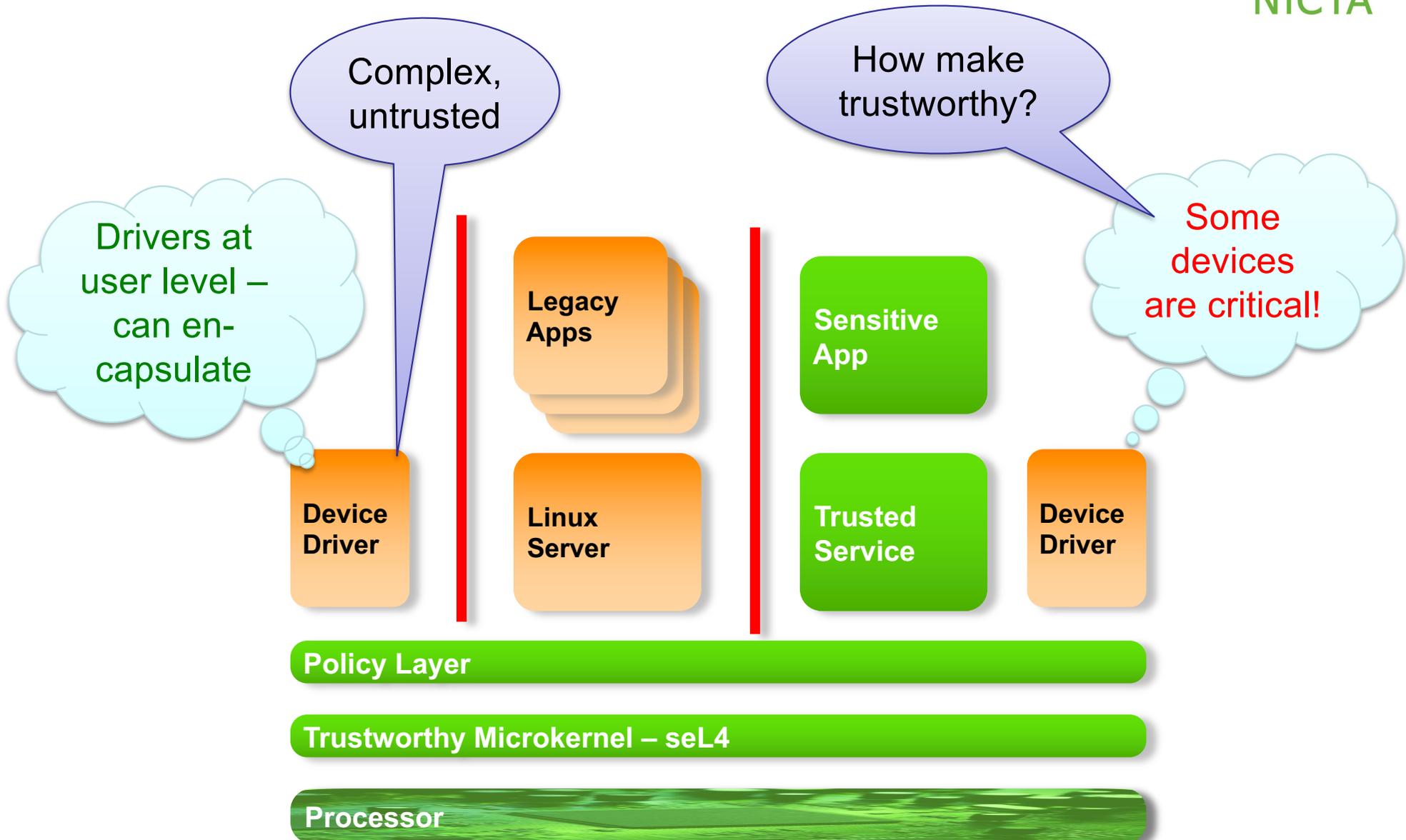


- Build system with minimal TCB
- Formalize and prove security properties about architecture
- Prove correctness of trusted components
- Prove correctness of setup
- Prove temporal properties (isolation, WCET, ...)
- Maintain performance

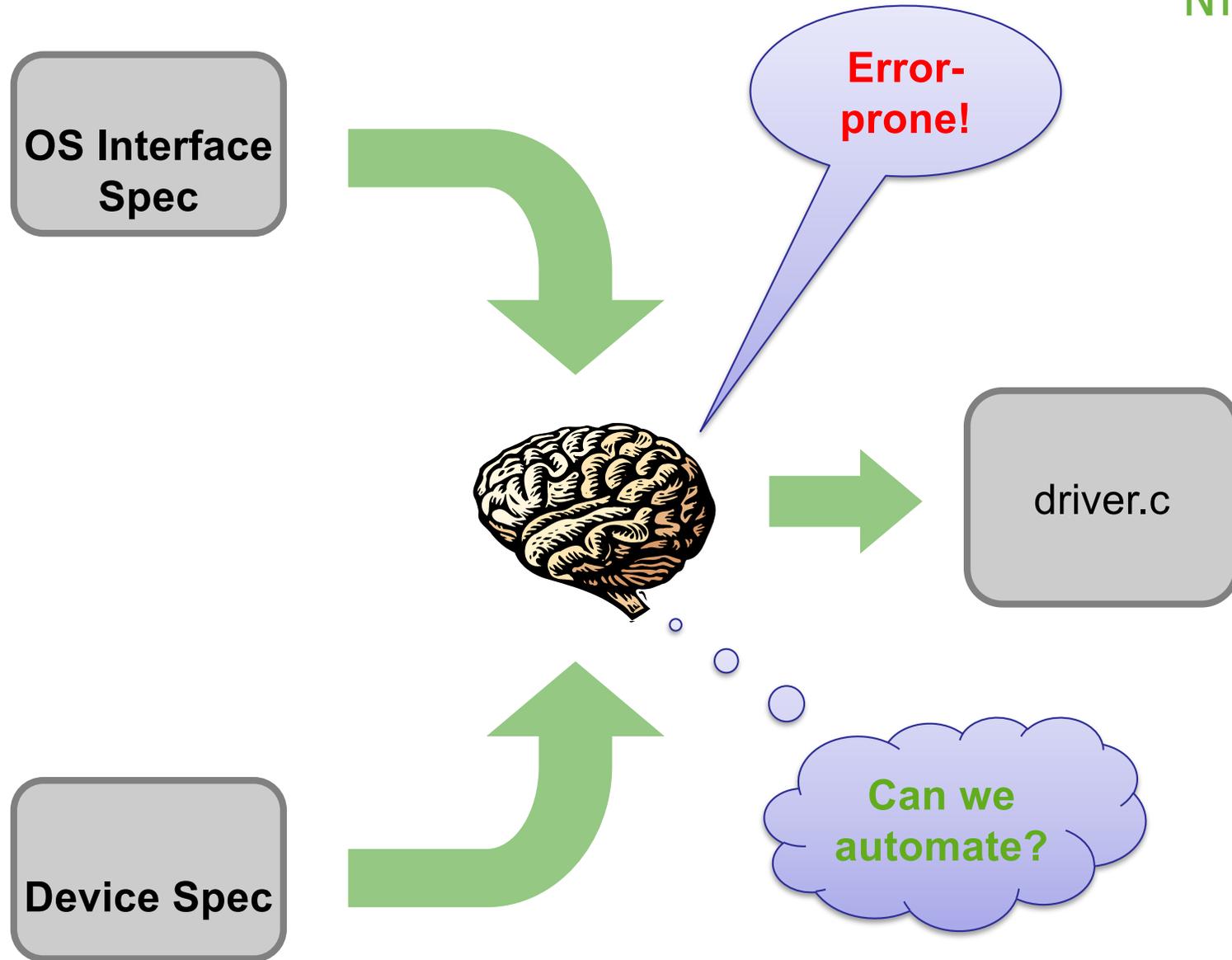
# Specifying Security Architecture



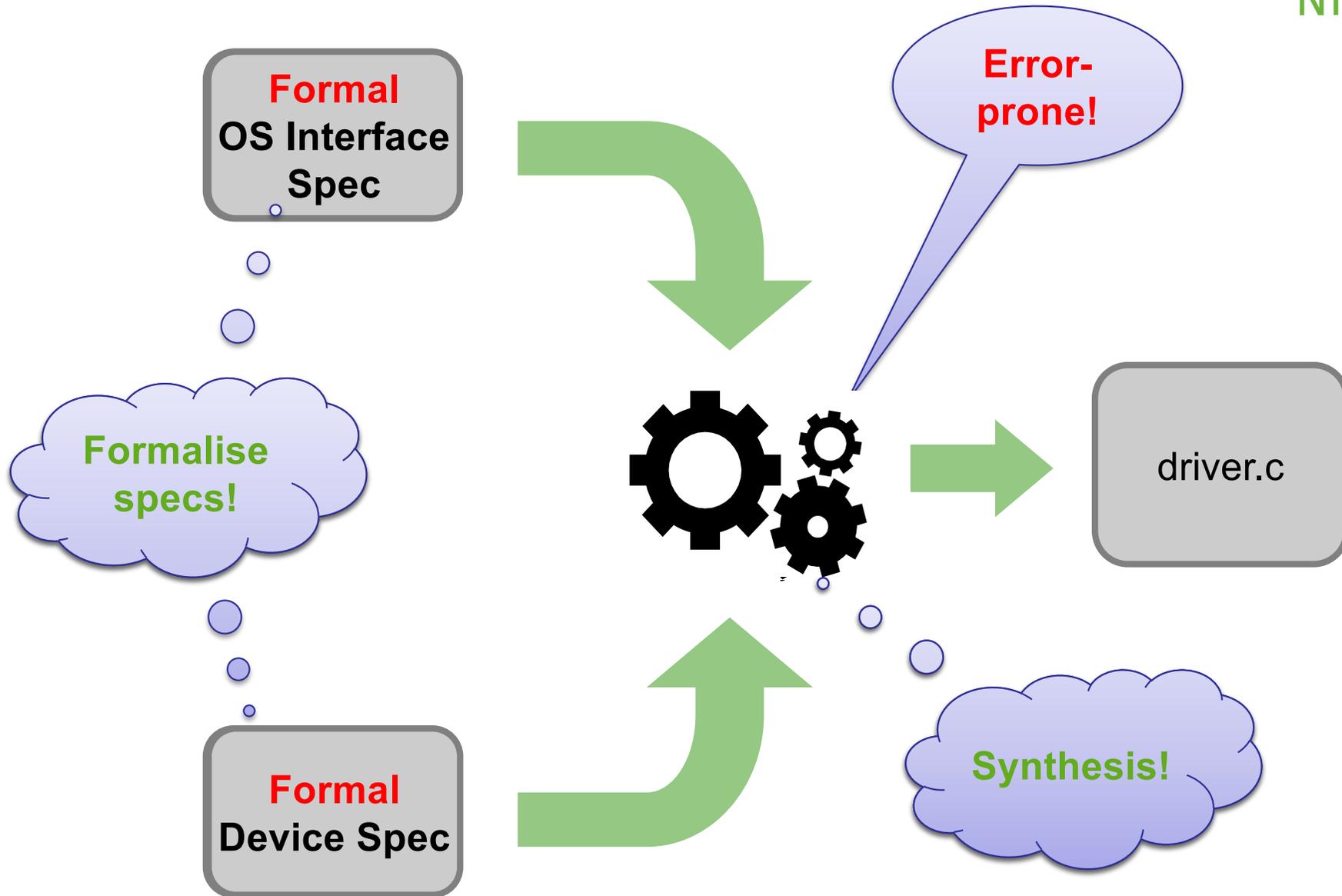
# Device Drivers



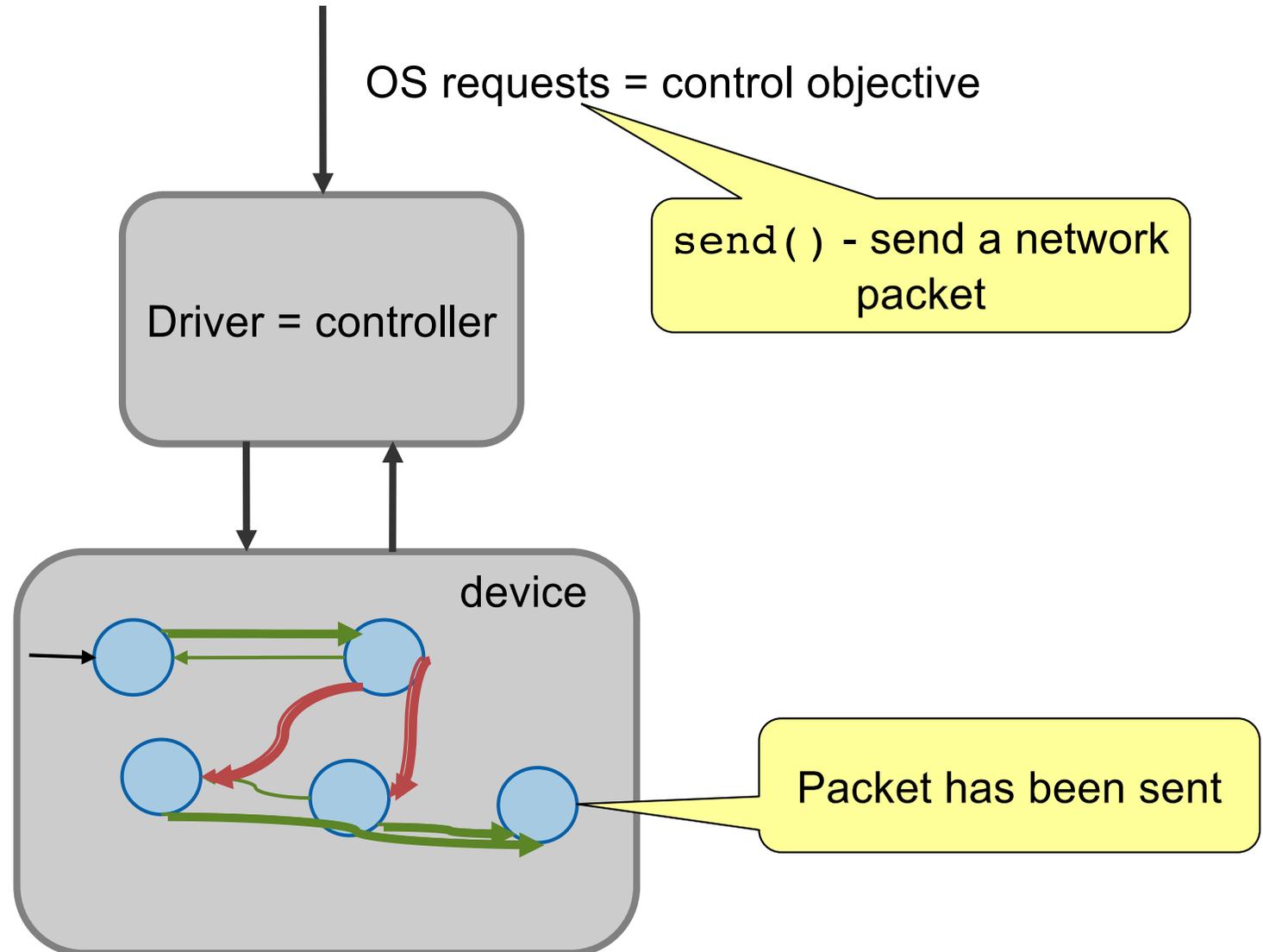
# Driver Development



# Driver Development

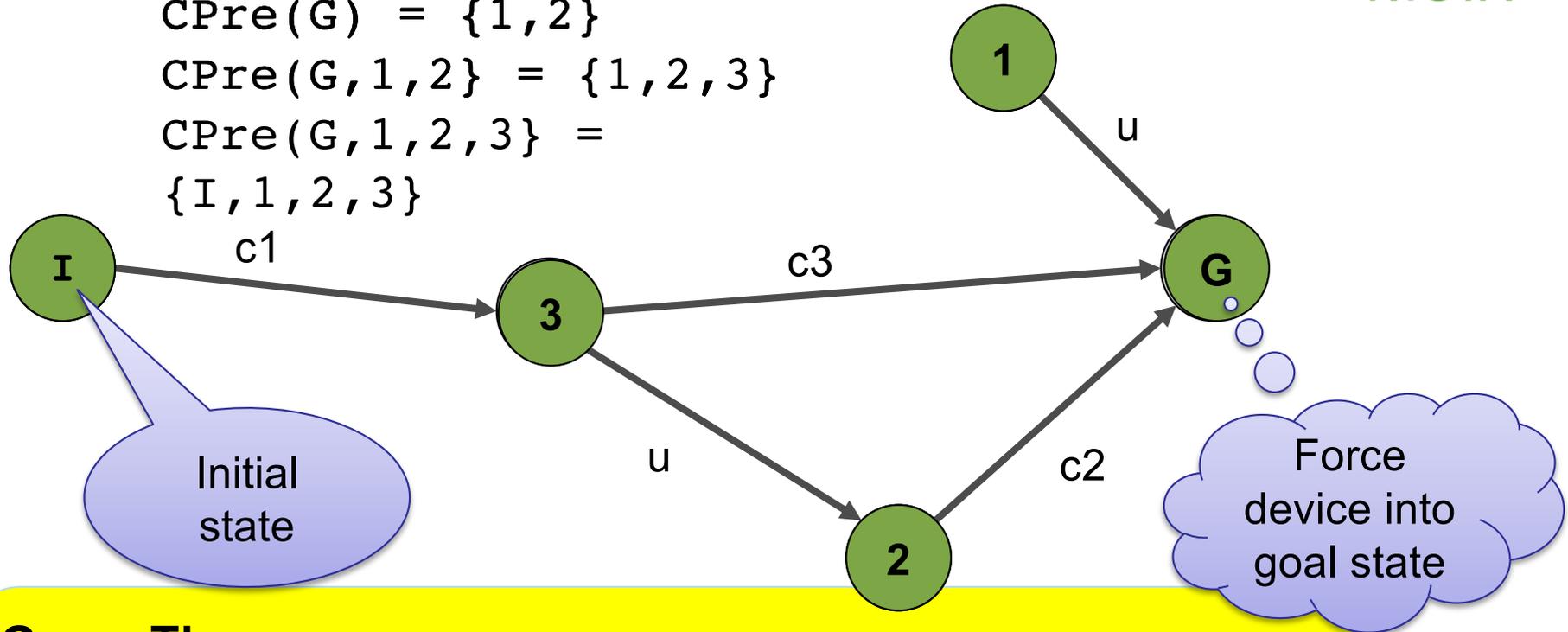


# Driver Synthesis as Controller Synthesis



# Synthesis Algorithm (Main Idea)

$CPre(G) = \{1, 2\}$   
 $CPre(G, 1, 2) = \{1, 2, 3\}$   
 $CPre(G, 1, 2, 3) = \{I, 1, 2, 3\}$



## Game Theory

- Framework for verification and synthesis of reactive systems
- Provides classification of games and complexity bounds
- Provides algorithms for winning strategies!

Device driver!

# Drivers Synthesised (To Date)



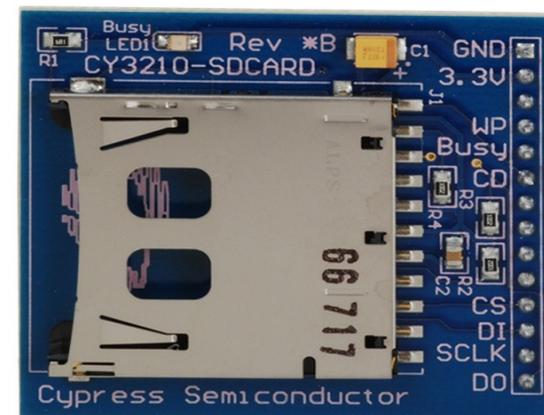
IDE disk controller



W5100 Eth shield

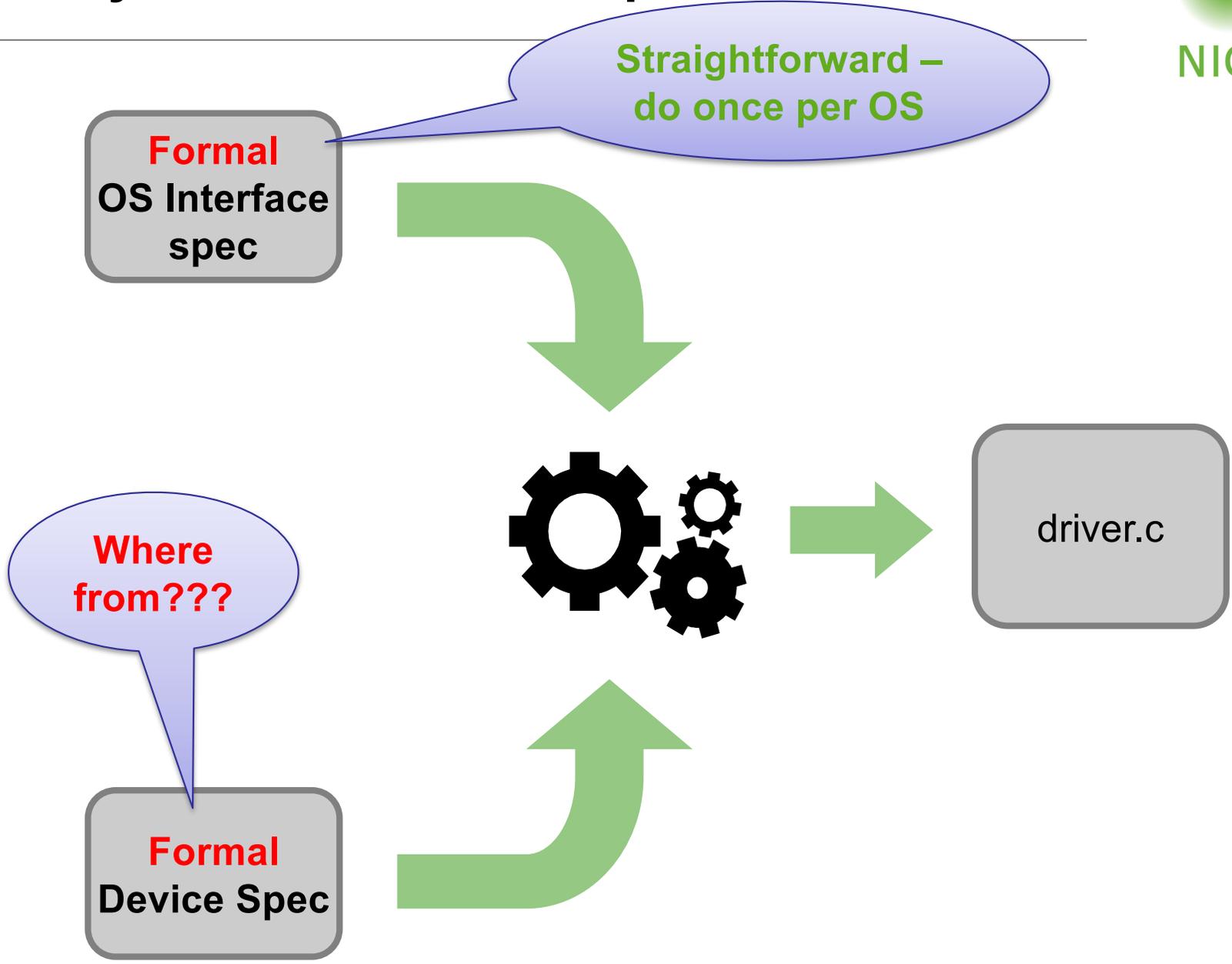


Asix AX88772  
USB-to-Eth adapter

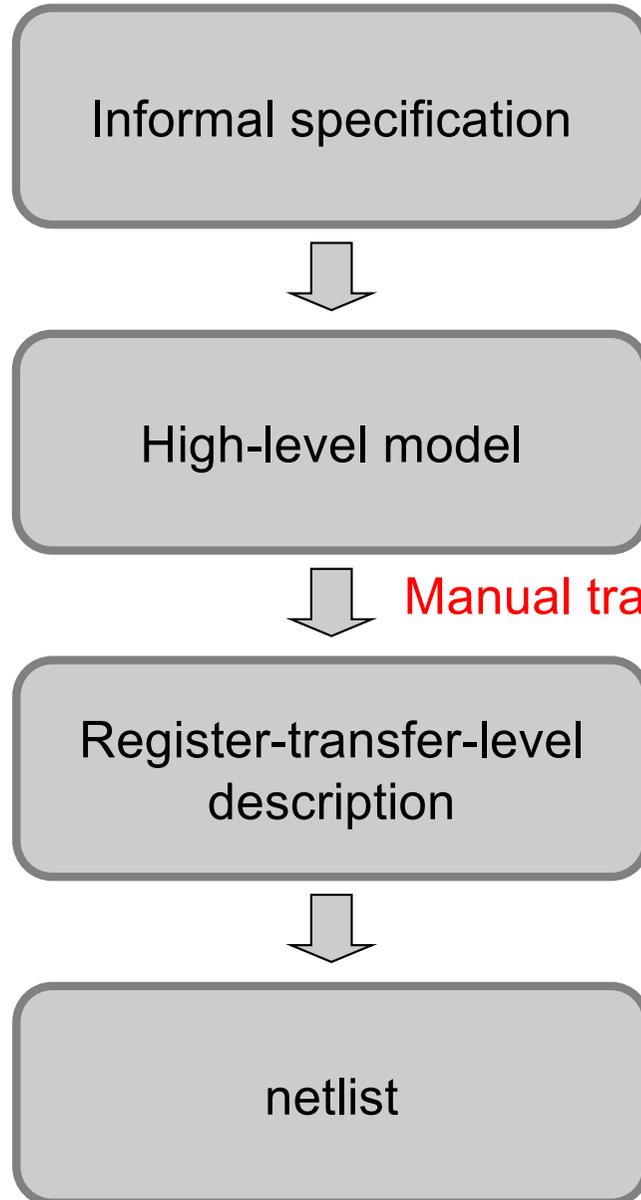


SD host controller

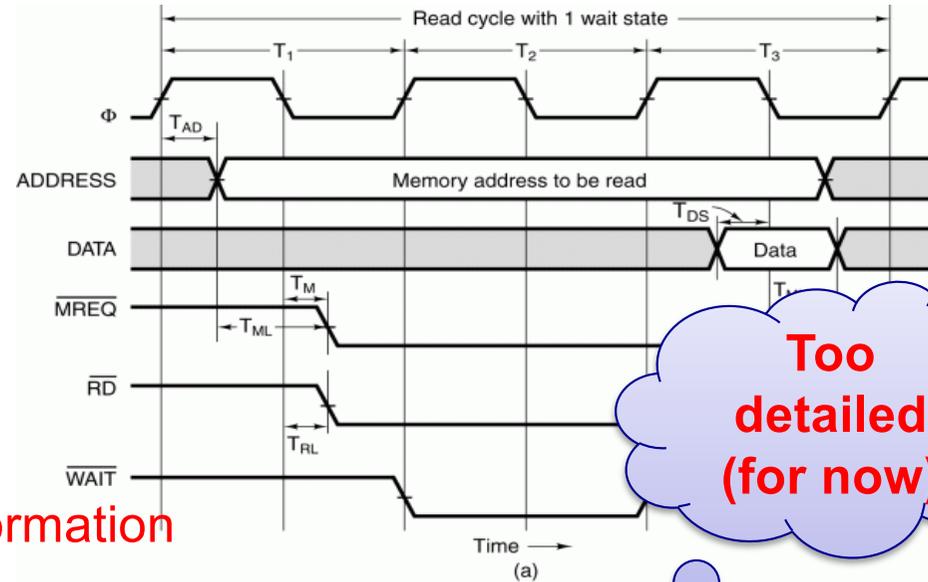
# Driver Synthesis: Interface Specs



# Hardware Design Workflow



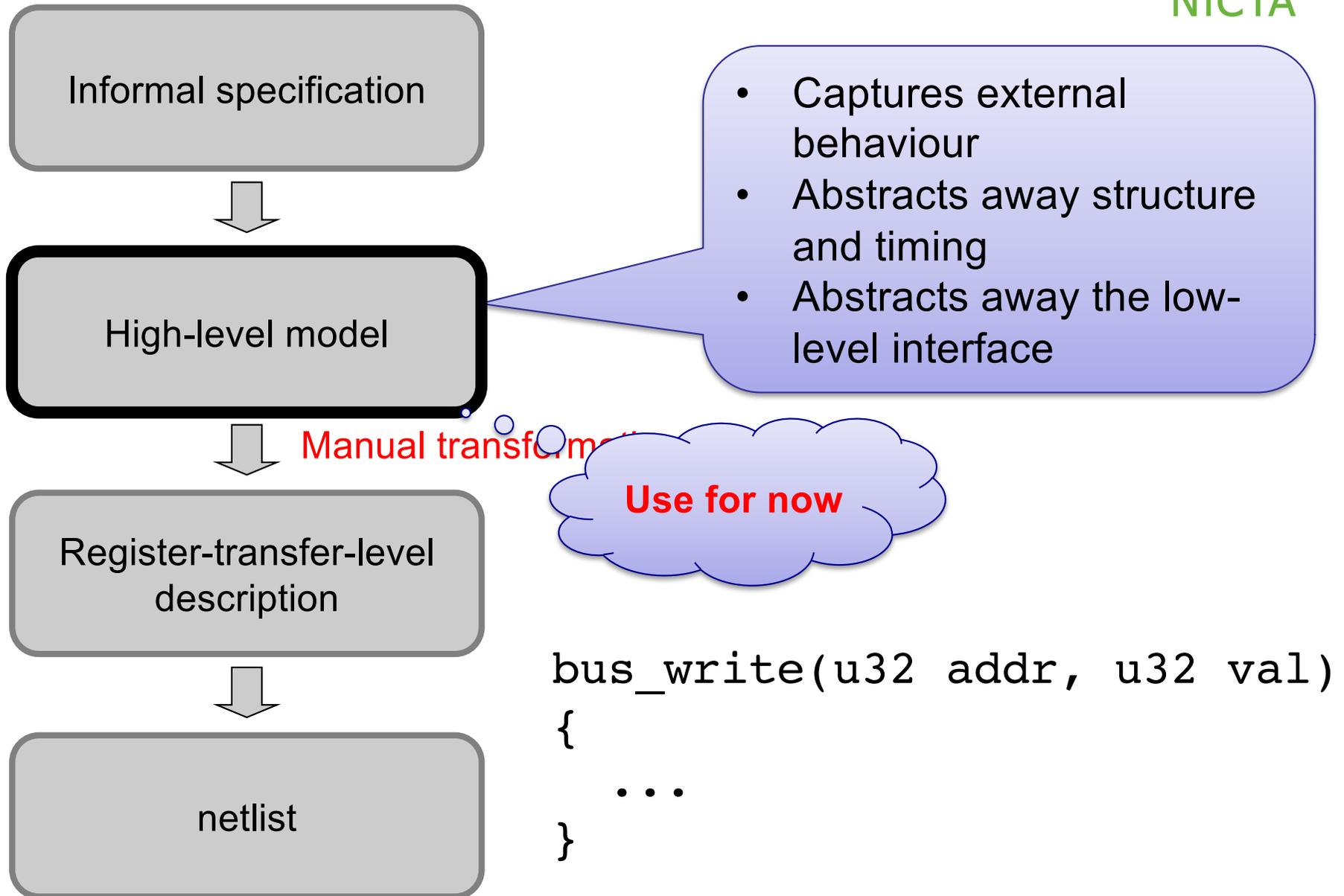
Manual transformation



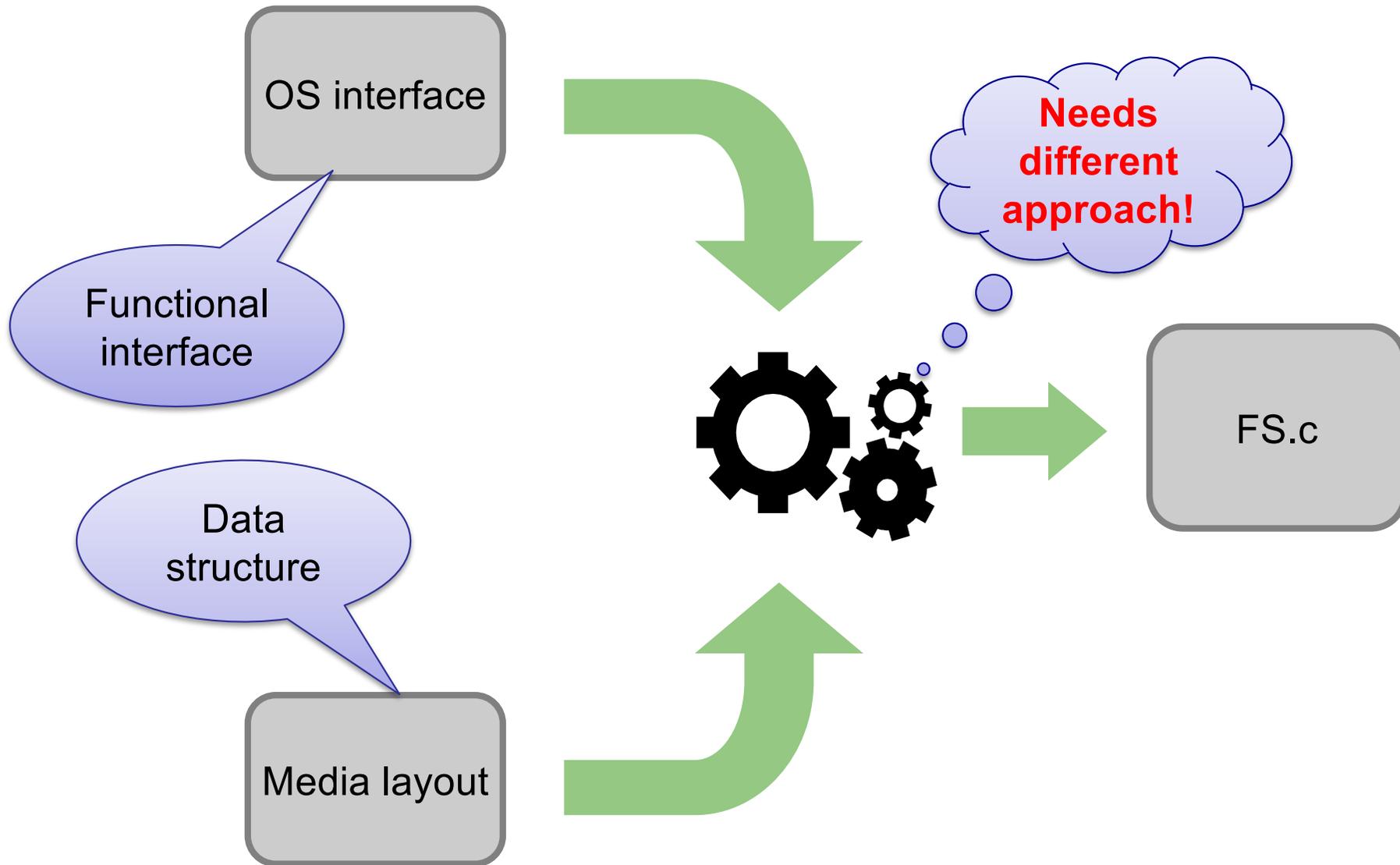
Too detailed (for now)

- Low-level description: registers, gates, wires.
- Cycle-accurate
- Precisely models internal device architecture and interfaces
- “Gold reference”

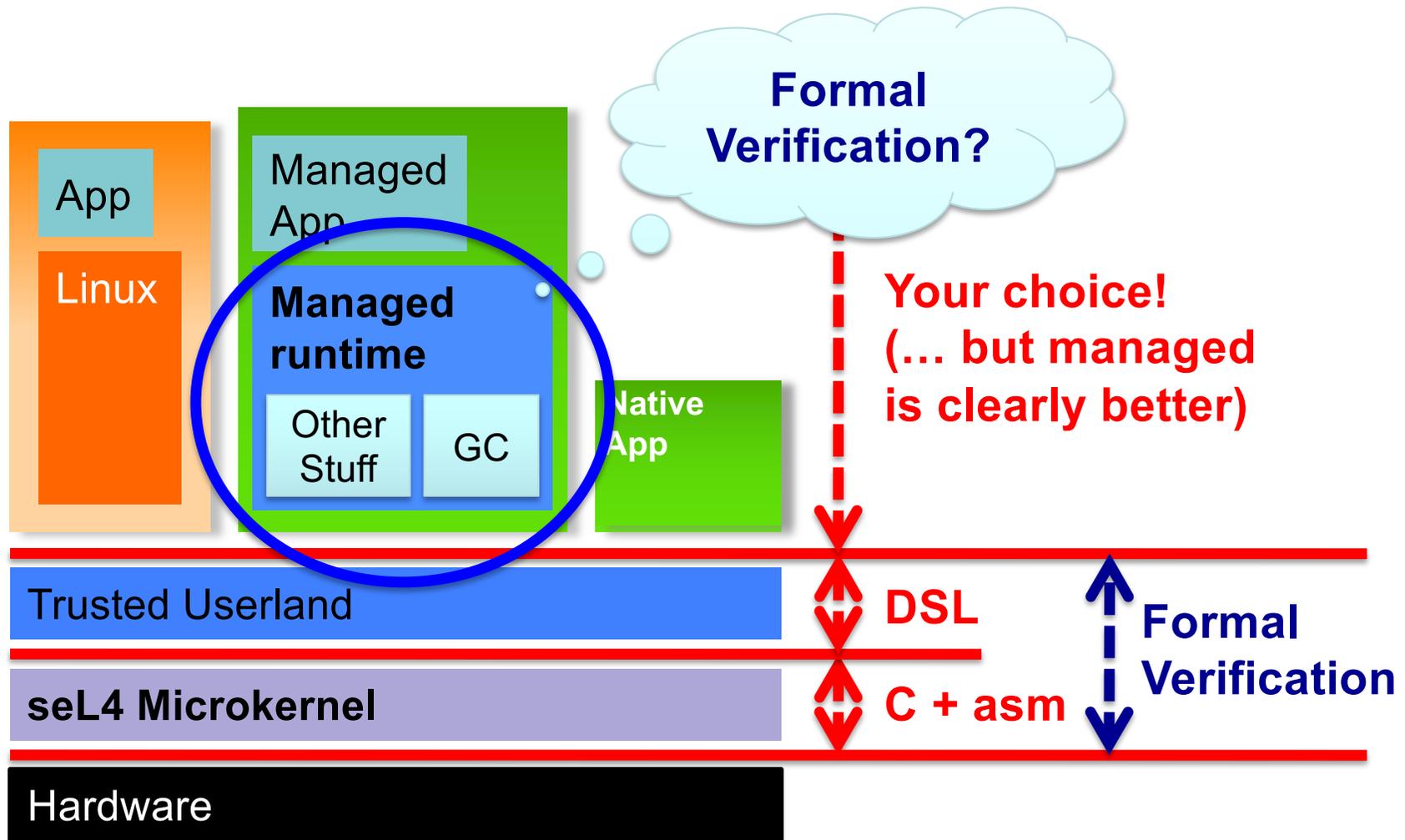
# Hardware Design Workflow



# From Drivers to File Systems?



# Building Secure Systems: Long-Term View



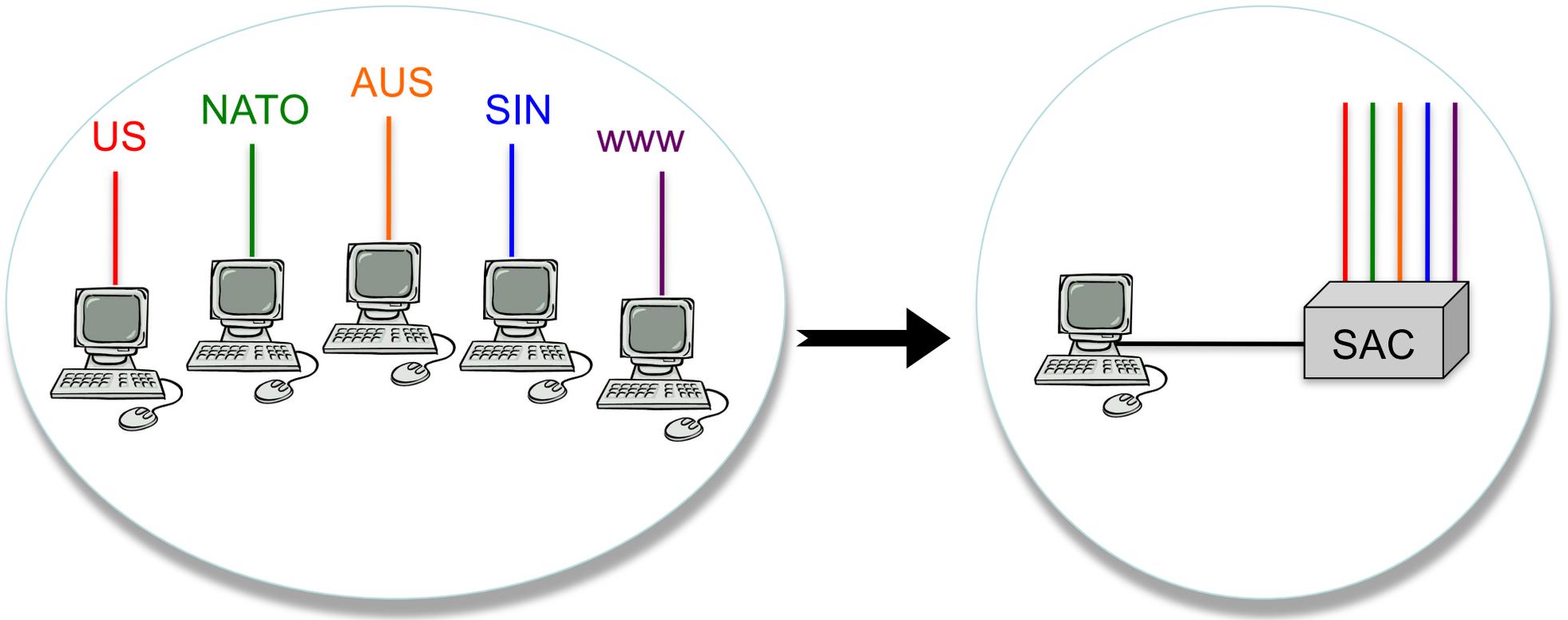
# Agenda

---

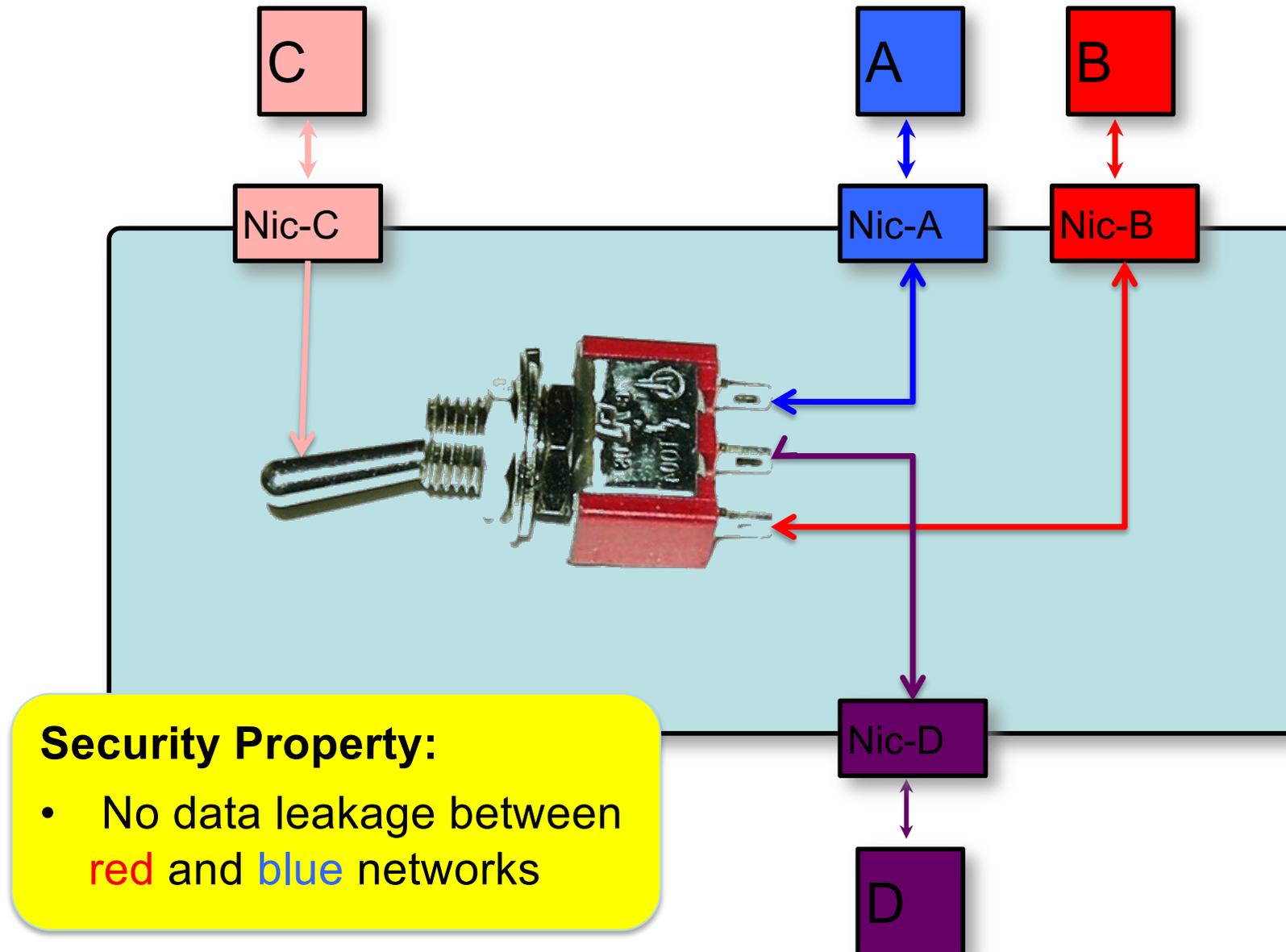


- Motivation
- What is a microkernel, and what is L4?
- seL4 – designed for trustworthiness
- Establishing trustworthiness
- From kernel to system
- **Sample system 1: Secure access controller**
- Sample system 2: RapiLog

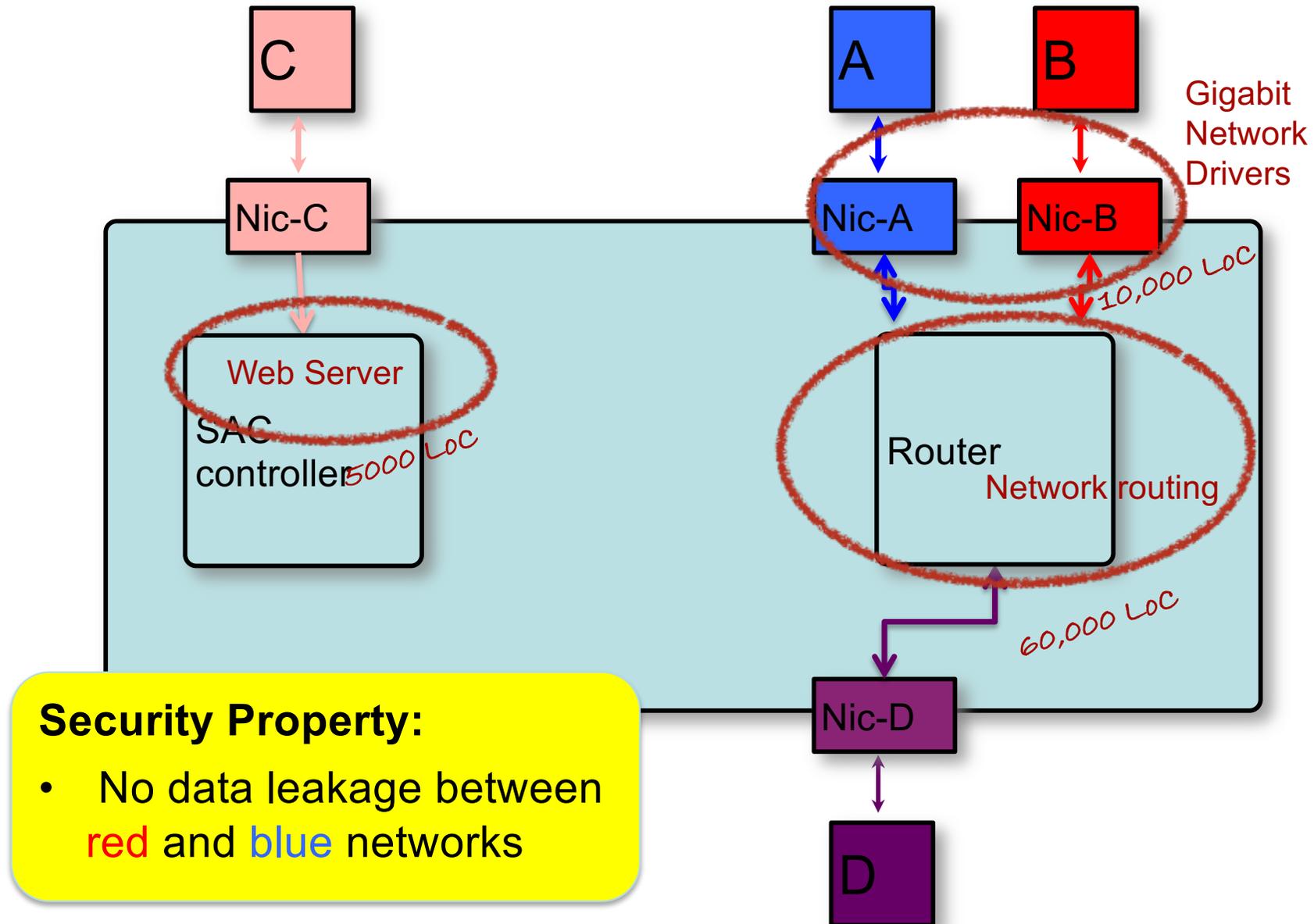
# Proof of Concept: Secure Access Controller



# Logical Function

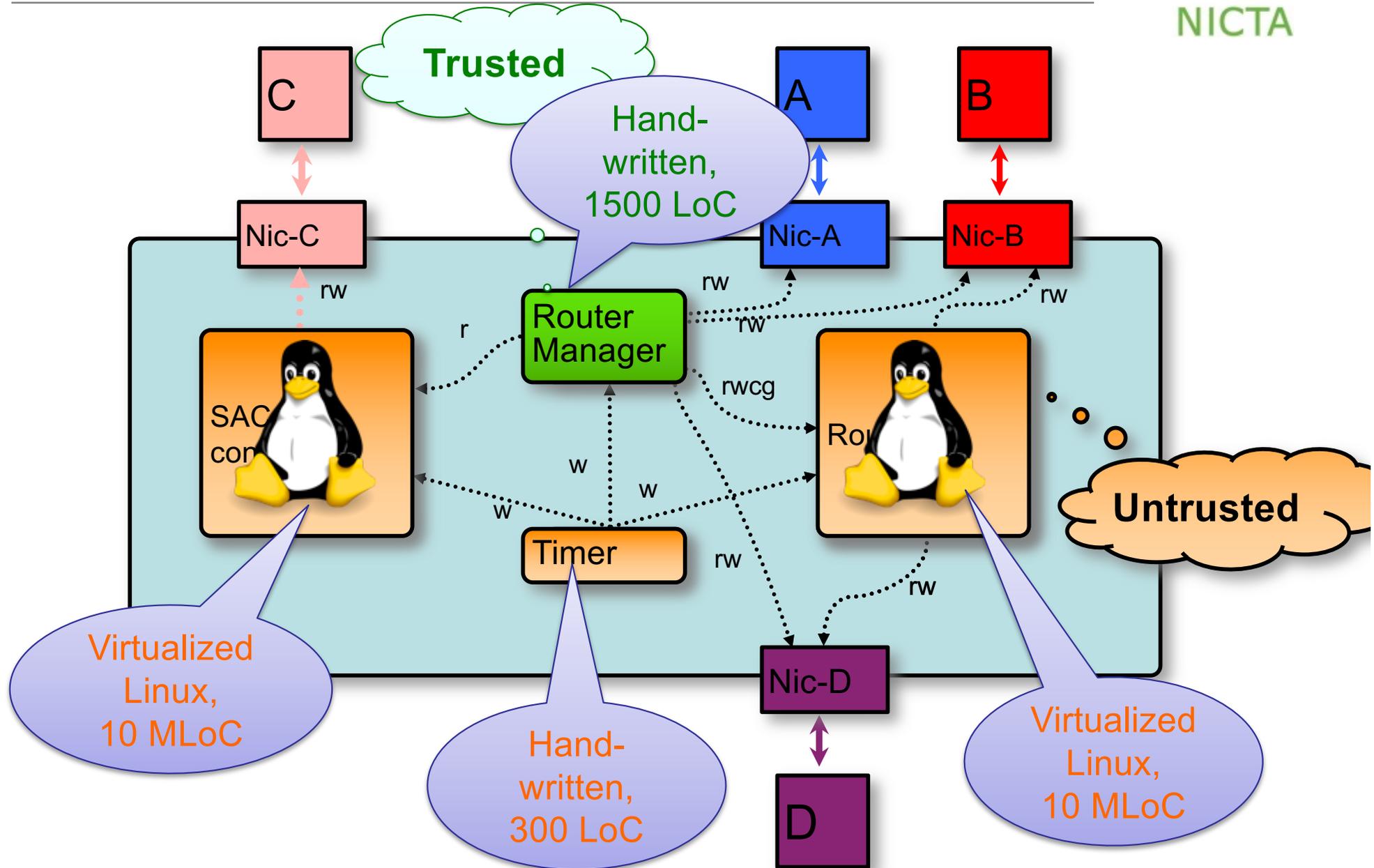


# Logical Function





# Implementation



# Agenda

---

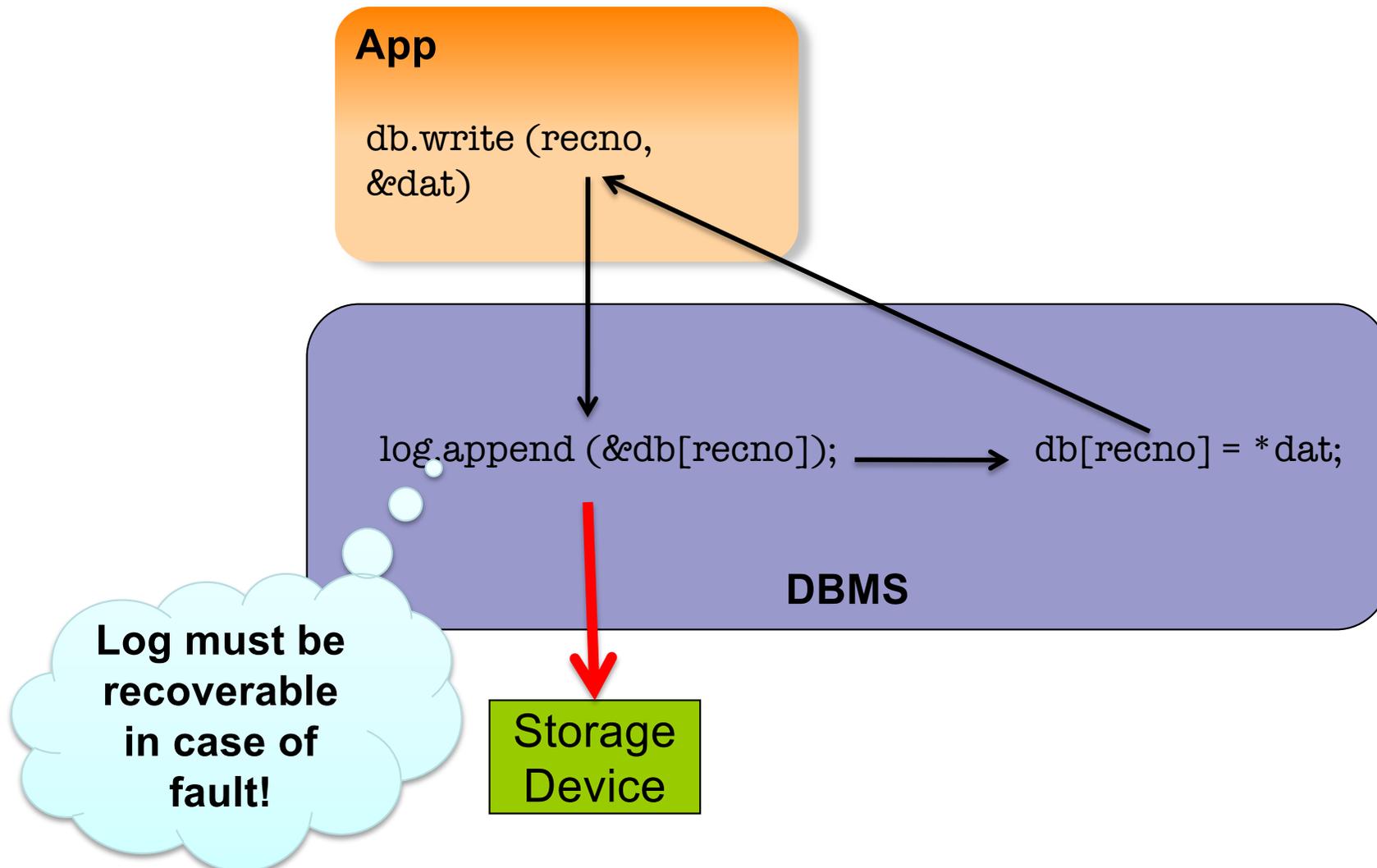


- Motivation
- What is a microkernel, and what is L4?
- seL4 – designed for trustworthiness
- Establishing trustworthiness
- From kernel to system
- Sample system 1: Secure access controller
- **Sample system 2: RapiLog**

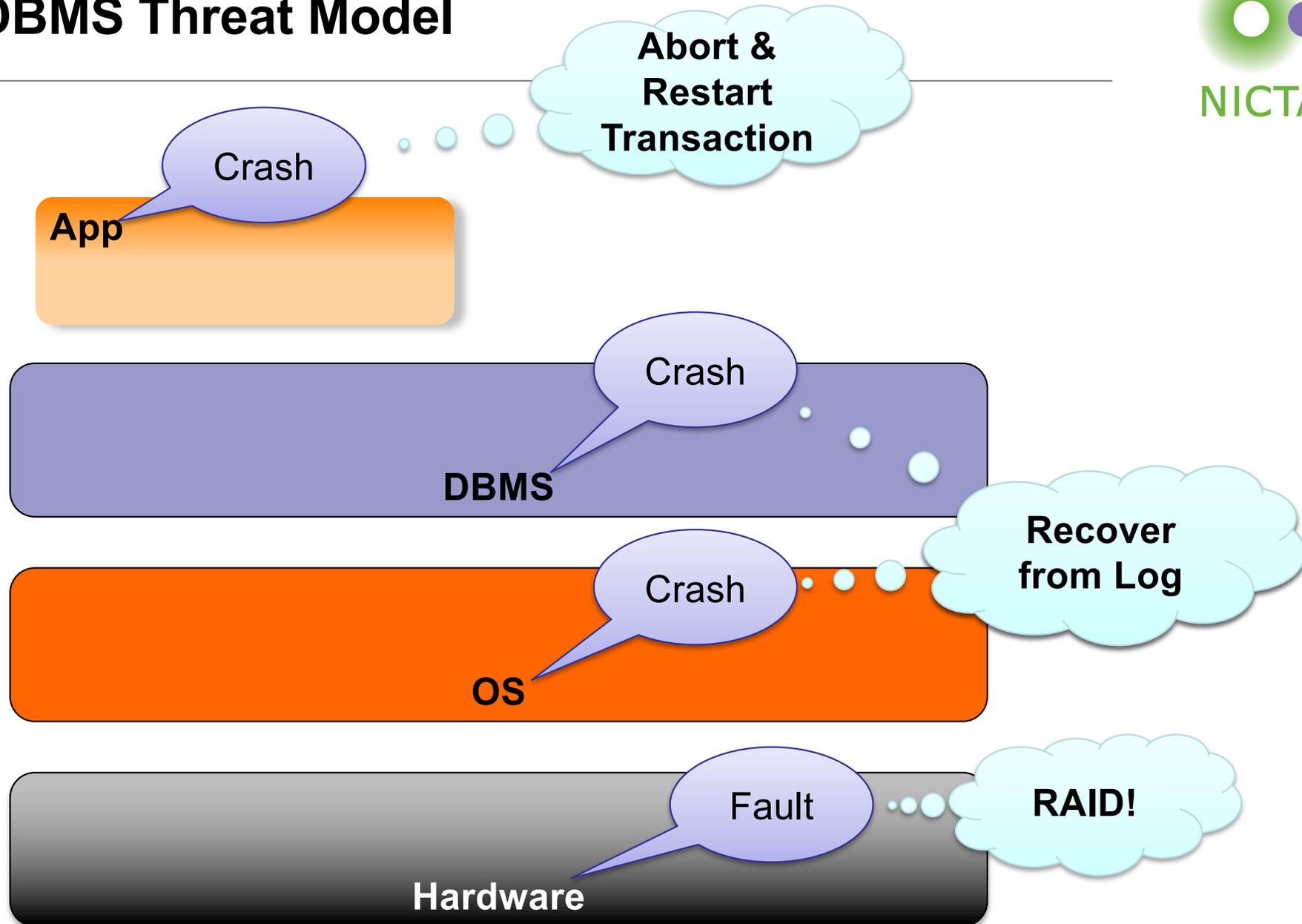
# Database Transactions



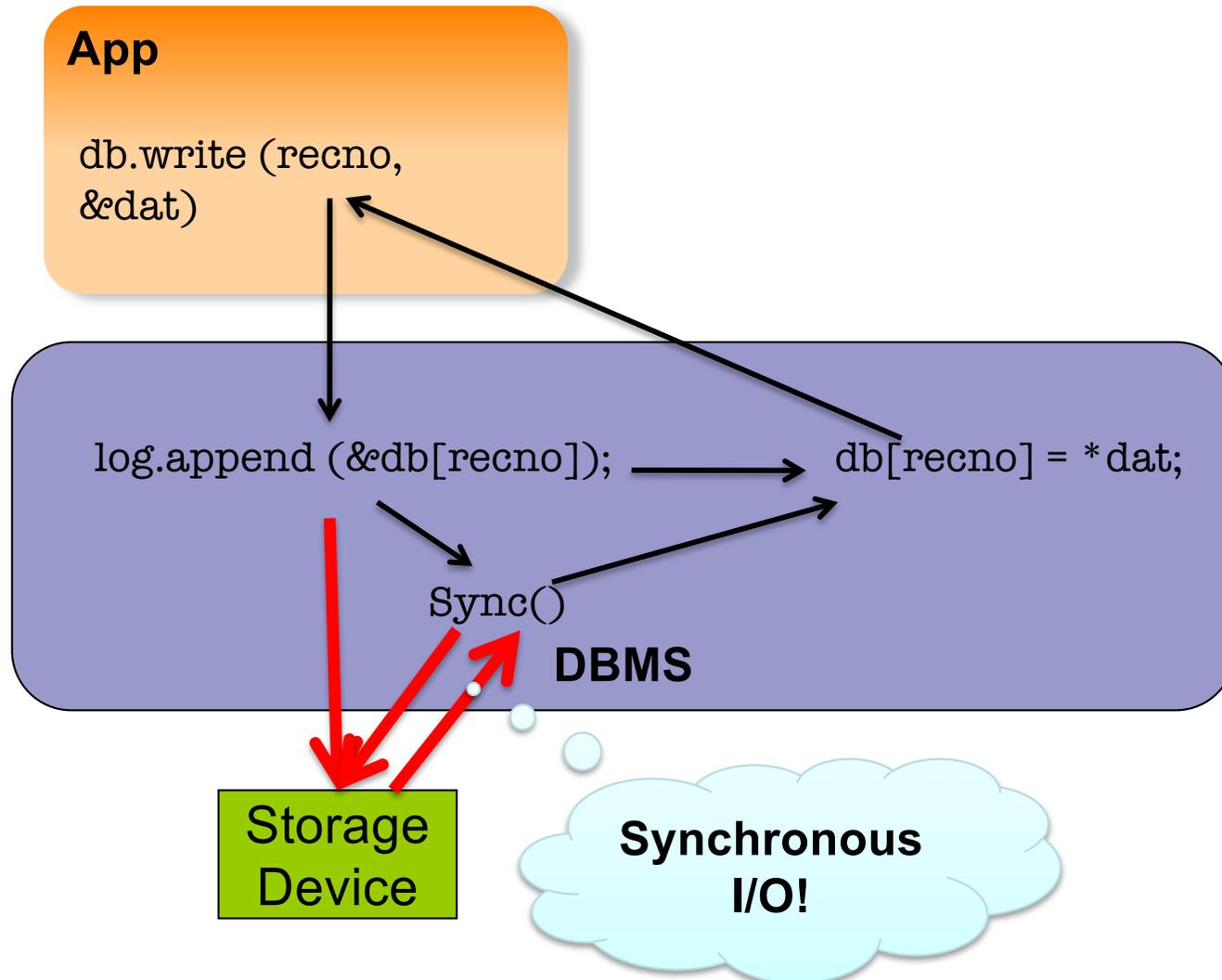
Various approaches, but today usually *write-ahead logging*:



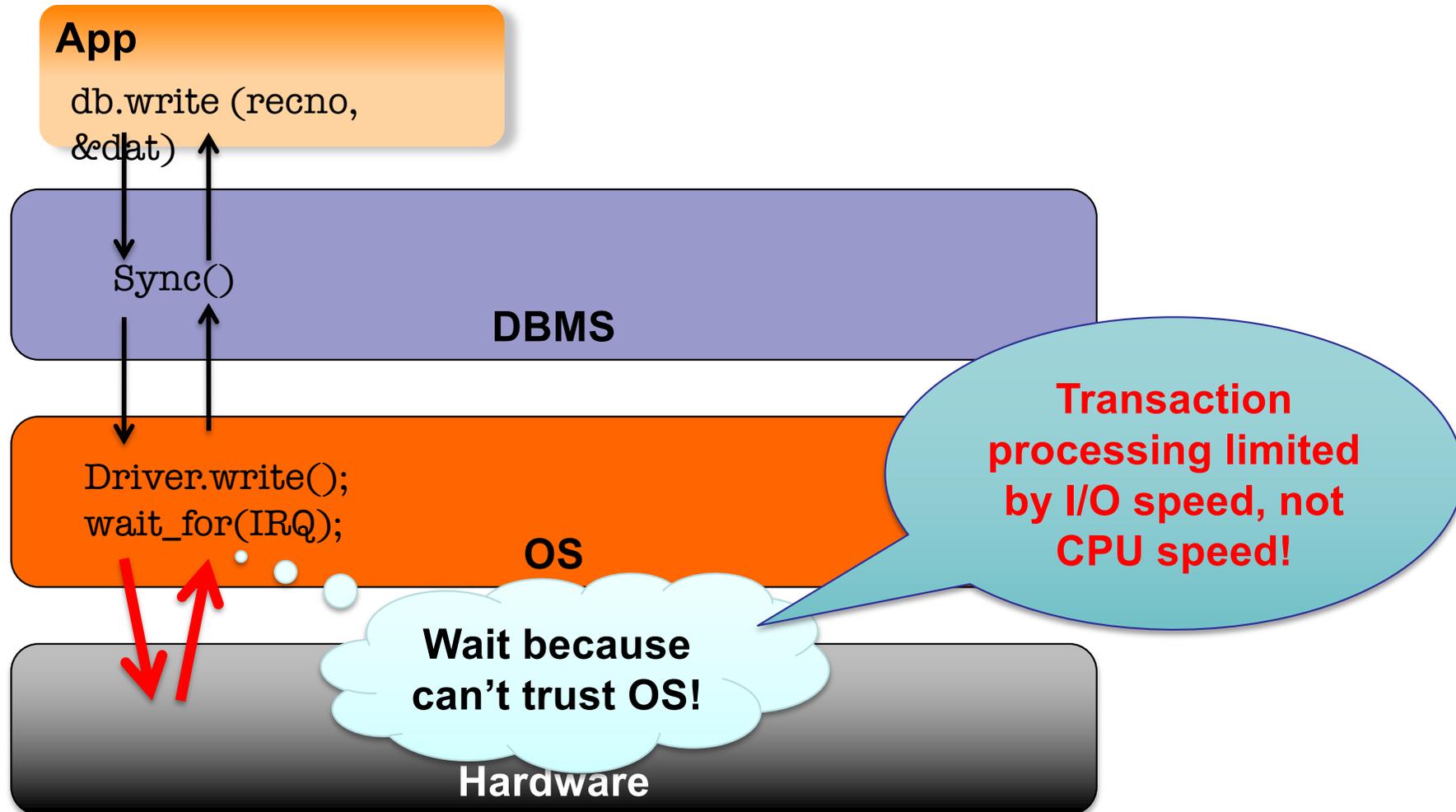
# DBMS Threat Model



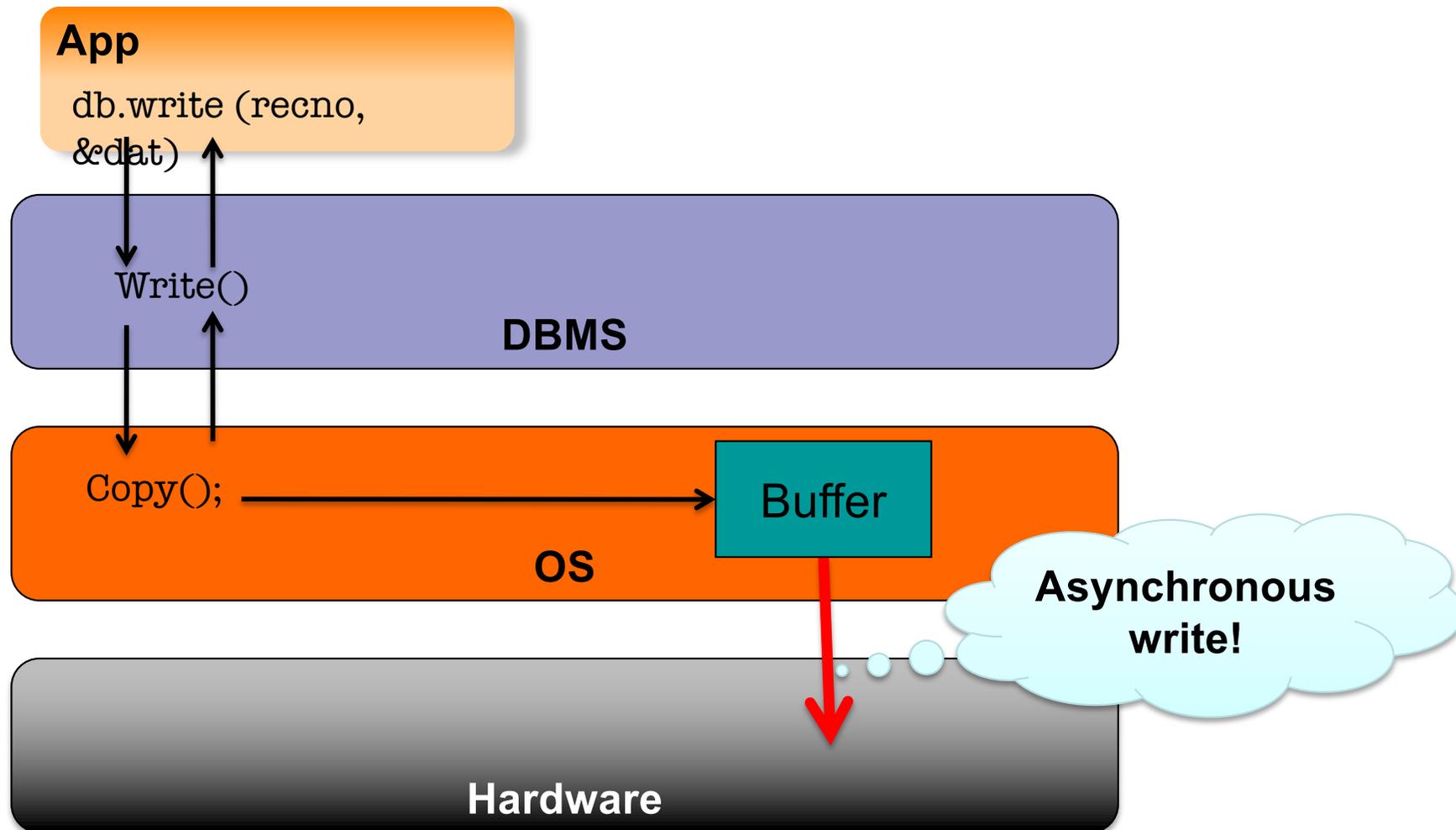
# Log Data Must Be Recoverable!



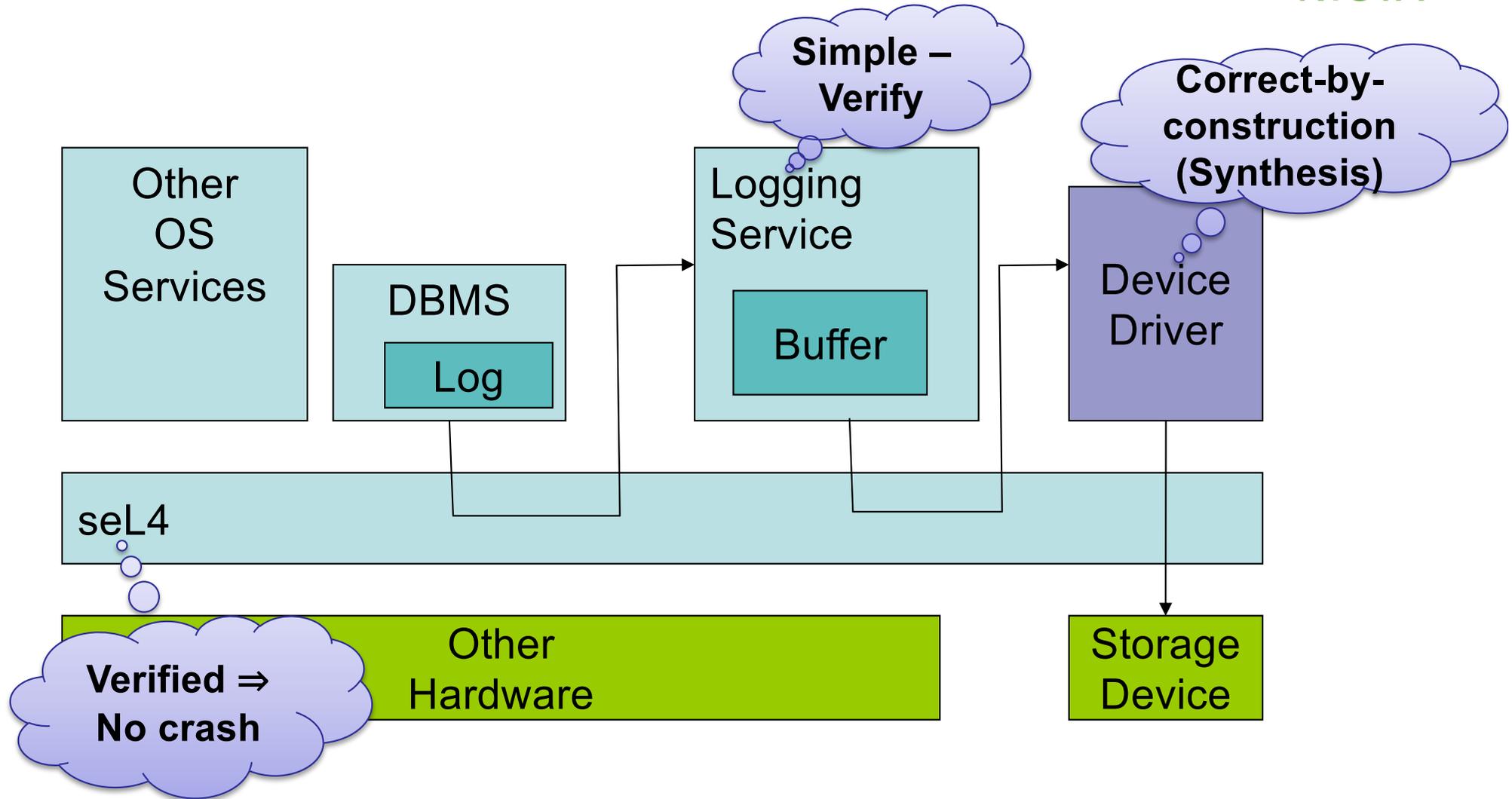
# Log Data Must Be Recoverable!



# What If We Could Trust the OS?

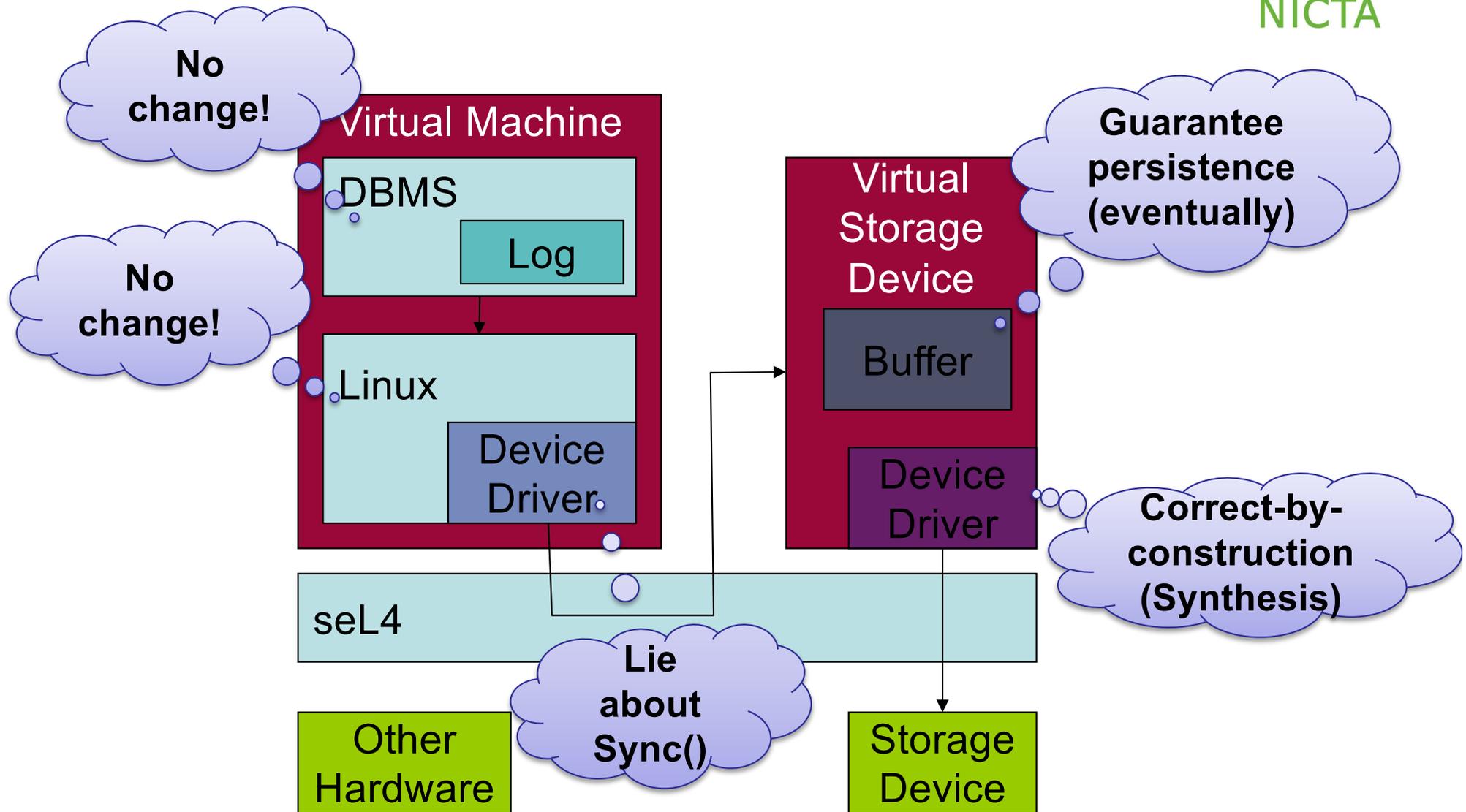


# But We Can Trust seL4!

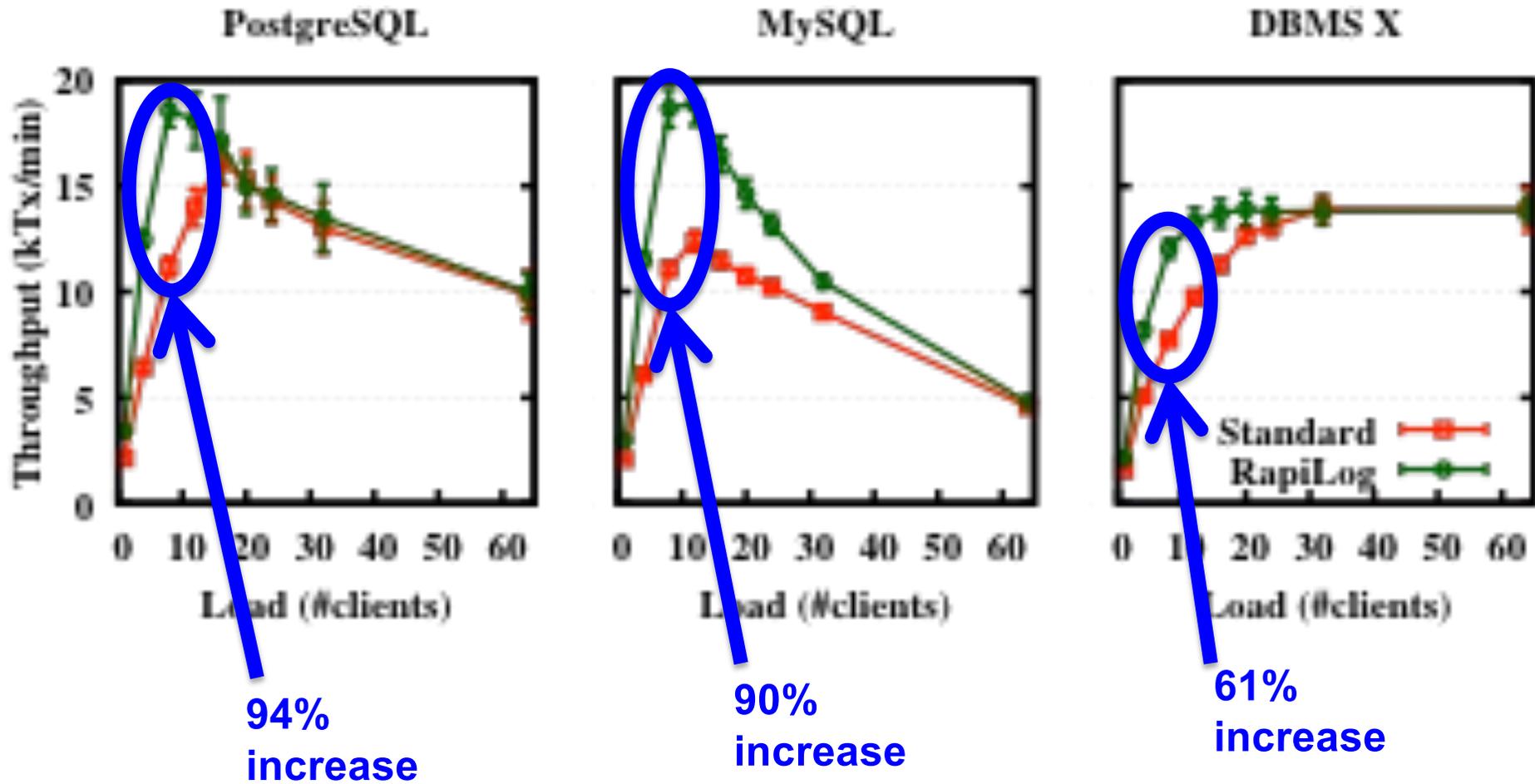


**Problem: Needs DBMS re-write**

# RapiLog: Use Virtualization

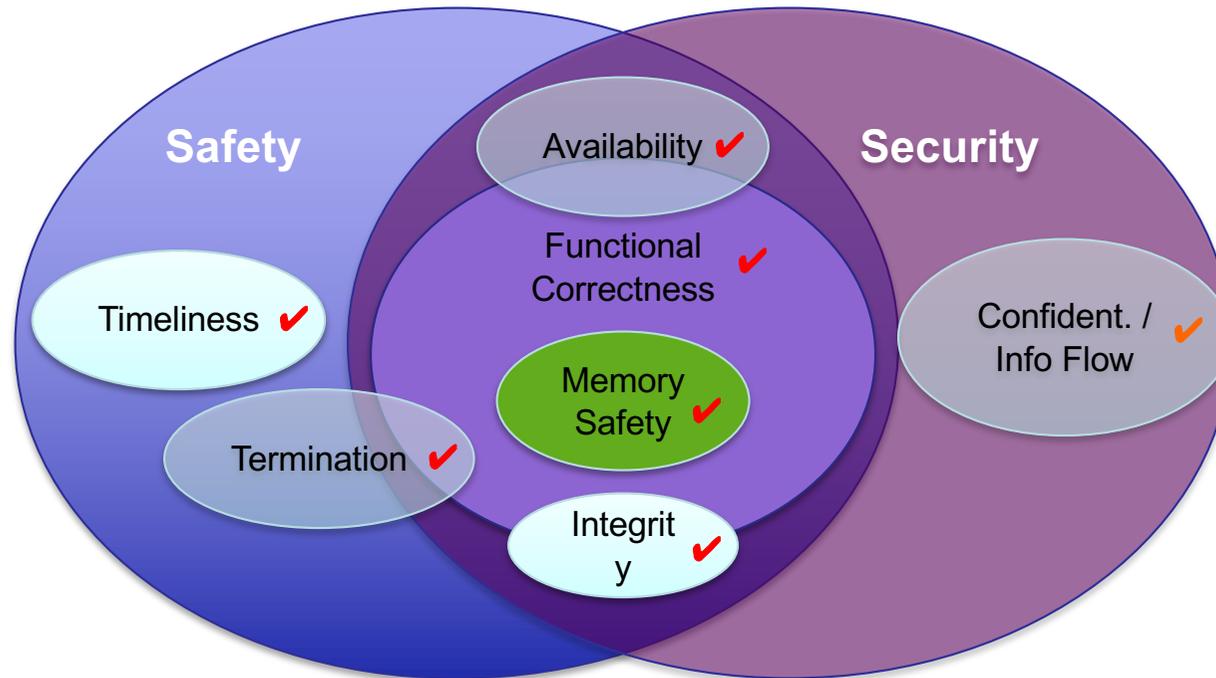


# Performance



**Also maintain durability on power failure!**

# Trustworthy Systems – We’ve Made a Start!



## Thank You!

<mailto:gernot@nicta.com.au>

Twitter @GernotHeiser

Google: “nicta trustworthy systems”