School of Computer Science & Engineering

**Trustworthy Systems Group**

# Why Change the Kernel When You Have seL4?

**Gernot Heiser**

gernot@unsw.edu.au
@microkerneldude.bsky.social
https://gernot-heiser.org/

# 3rd Workshop on Kernel Isolation, Safety and Verification (KISV 2025)

**October 13, 2025**
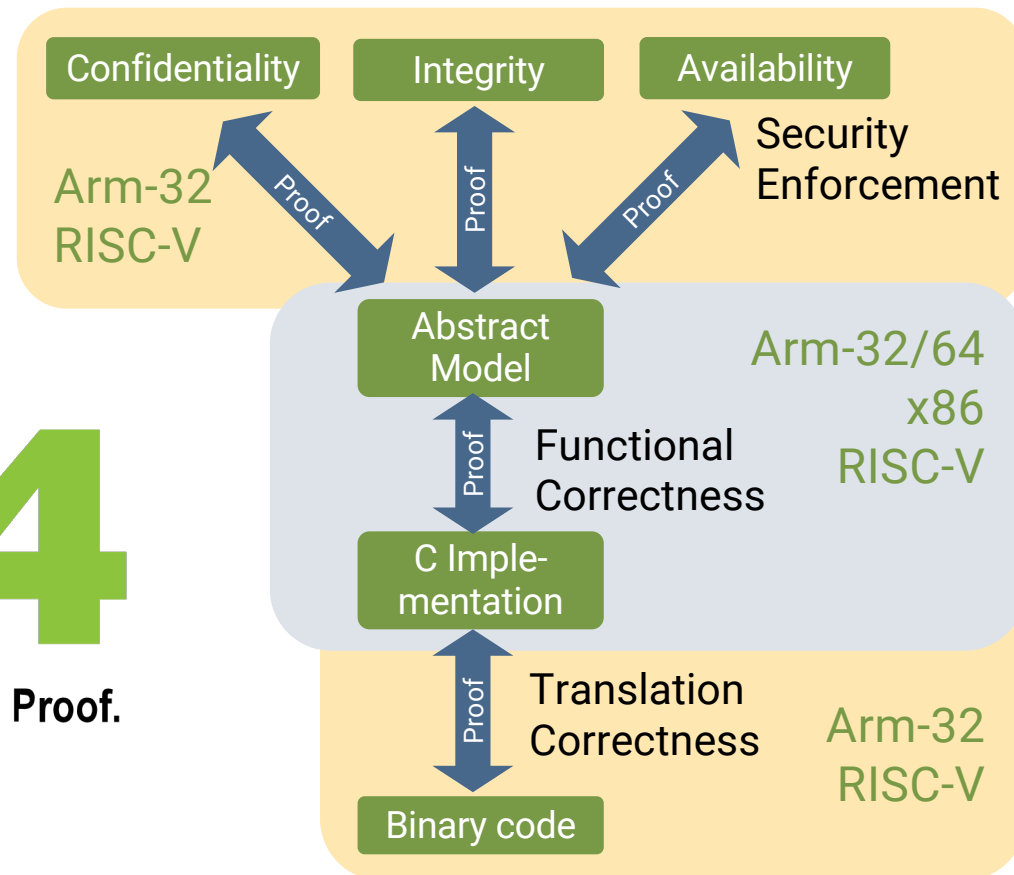**Seoul, Republic of Korea**

In conjunction with the
31st ACM Symposium on Operating Systems Principles (SOSP '25)

This workshop aims to bring together researchers and developers from the field of operating systems, programming languages, security, computer architecture and verification with the goal to accelerate changes in the kernel through a combination of isolation, programming language safety, and formal verification.

**Which kernel?**

**Why?**

UNSW SYDNEY

# Are We Talking About This Kernel?



**sel4**

**Security. Performance. Proof.**

Confidentiality — Integrity — Availability

Arm-32
RISC-V

Security
Enforcement

Proof — Proof — Proof

Abstract
Model

Arm-32/64
x86
RISC-V

Proof — Functional
Correctness

C Imple-
mentation

Proof — Translation
Correctness

Arm-32
RISC-V

Binary code

UNSW
SYDNEY

# Or This One?

# Changing The Kernel: Modules in Rust



User

Kernel

40M SLoC of C

- ✓ Protects rest of kernel from Rust modules
- ❖ Unsafe code?
- ❖ Doesn't protect Rust module from rest of kernel
- ❖ Requires writing modules from scratch

UNSW
SYDNEY

# Changing: Partition Kernel Space

User

Kernel

| 40M SLoC of C | | Partition – restricted address space |
| Partition – restricted address space | | |
| Core kernel – full access | | |

- ✓ Protects partitions from each other
- ❖ Requires HW extensions ⇒ not on older platforms
- ❖ Increases ISA/manufacturer dependence
- ❖ Still fully trust core kernel
- ❖ Privilege revocation?
- ❖ Cost?

UNSW
SYDNEY

# Changing: Partition Kernel Space

User

Kernel

| 40M SLoC of C | | Partition – restricted address space |
|---|---|---|
| | Partition – restricted address space | |
| Monitor – super-privileged | | |

Hypv.

- ✓ Protects partitions from each other
- ❖ Requires HW extensions ⇒ not on older platforms
- ❖ Increases ISA/manufacturer dependence
- ❖ Fully trust monitor (but it does very little)
- ❖ Privilege revocation?
- ❖ Cost?
- ❖ **Squats hypervisor mode – lose virtualisation support?**

UNSW SYDNEY

# Changing: De-Privileging "Kernel"

User



Kernel/
Hypv.

- ✓ Protects **all** modules from each other
- ✓ Requires no special HW
- ✓ Verified kernel
- ✓ Retain virtualisation support
- ❖ High cost?

# Microkernel Overheads

### Round-trip cross-address-space IPC on 64-bit Intel Skylake

| | seL4 | Fiasco.OC aka L4Re | Google Zircon |
|---|---|---|---|
| Latency (cycles) | 986 | 2717 | 8157 |
| Mandatory HW cost* (cycles) | 790 | 790 | 790 |
| Overhead absolute (cycles) | 196 | 1972 | 7367 |
| Overhead relative | 25% | 240% | 930% |

**Smaller is better**

**\*:** The Cost of SYCALL + 2 × SWAPGS + SYSRET = 395 cycles, times 2 for round-trip

**Source:**

Zeyu Mi, Dingji Li, Zihan Yang, Xinran Wang, Haibo Chen: "SkyBridge: Fast and Secure Inter-Process Communication for Microkernels", EuroSys, April 2019

UNSW SYDNEY

# Microkernel Overheads

High syscall rate = 61k/s

seL4 round-trip address-space switch = 1k cy

Conservative IMHO

Assume average 2 R-T AS switches / syscall:

  Switch O/H = 2 × 61k/k × 1kcy = 122M cy/s

Assume 3GHz clock:

  O/H = 122M cy/s / 3Gcy/s = 122/3k = 4%

Assume 4-core CPU:

  O/H = 4%/4 = 1% of CPU!

Assume Linux max CPU load = 25%

  relative O/H = 4 × 1% = 4%

Why would anyone care?

UNSW SYDNEY

# But Is This Real?

Test bed: LionsOS

- Simple, from-scratch seL4-based OS
- Highly modular design,
  strict separation of concerns
- Adaptable "Lego® kit" approach
- Designed for embedded / cyber-
  physical systems

# Underneath https://sel4.systems/

# Networking Layer

Client can be a VM

IP stack is library – not in system's TCB!

Tx Virt encapsulates traffic-shaping policy

Strict separation of concerns!

Client

lwIP

Copy

Handles broadcasts

ARP

Client

lwIP

Copy

Copier for security (if needed)

Tx Virt

Rx Virt

Driver

NIC

Virtualiser shares device, incl address mapping, cache maintanance

Translates HW-specific device interface to HW-independent device-class interface

UNSW
SYDNEY

# Networking Layer



**Zero-copy communication:**
- Lock-free, single-producer, single-consumer, bounded queues
- Synchronised by semaphores

**Benefits:**
- simple components
- **location transparency**

# Packet Round-Trip Context Switches



15 context switches per packet!

#12: Completion IRQ – switch Idle→Driver

#11: Driver→Idle

#13: Driver→Tx Virt

#9: Client→Tx Virt

#10: Tx Virt→Driver

Client

lwIP

Copy

#15: Client→Idle

#14: Tx Virt→Client

ARP

Tx Virt

Rx Virt

Driver

NIC

Tx

Rx

Client

lwIP

Copy

#4: Copy→Client

#8: Driver→Client

#3: Rx Virt→Copy

#5: Client→Copy

#2: Driver→Rx Virt

#6: Copy→Virt

#1: Rx IRQ – switch Idle→Driver

#7: Virt→Driver

**Must return free buffers!**

UNSW SYDNEY

# Comparing to Linux

Completion IRQ –
No context switch

#1: Rx IRQ –
switch Idle→Client

#2: Client→Idle

Client

Idle

IP Stack ↔ NIC Driver

NIC

Client

IP Stack

Tx Virt

Copy

Driver

ARP

Client

Copy

Rx Virt

seL4

NIC

- **LionsOS:** 30 mode switches, 15 context switches
- **Linux:** 6 mode switches, 2 context switches

UNSW
SYDNEY

# Comparing Performance: Setup



- External load generator
- Measures throughput, latency
- Client echoes packets

# What Do We Expect?

Ethernet packet size = 1.5kB

Assume Linux mode switch = half context switch

LionsOS O/H = 12/pk × 0.5k cy = 6k cy/pkt

Max packet rate for 1Gb/s NIC:

  rate = 1Gb/s / 1.5kB = 1Gb/s / 12kb = 833k/s

Worst-case O/H for 1Gb/s NIC:

  O/H = 6k cy/pkt * 833k pkt/s = 0.5G cy/s

Assume 3GHz clock:

  rel O/H = 0.5G cy/s / 3G cy/s = 17% of core

UNSW
SYDNEY

# However, There's Batching

- Each component will process everything in its queue before signalling another component
- No component will ever busy-poll!

Client

lwIP — Copy

ARP

Client

lwIP — Copy

Tx Virt

Rx Virt

Driver

Tx

Rx

NIC

- Dramatically reduces context switches under load!
- Measure 5–10 pkt/IRQ!

UNSW SYDNEY

# What Do We Expect?

Ethernet packet size = 1.5kB

Assume Linux mode switch = half context switch

LionsOS O/H = 12/pk × 0.5k cy = 6k cy/pkt

Max packet rate for 1Gb/s NIC:

  rate = 1Gb/s / 1.5kB = 1Gb/s / 12kb = 833k/s

Worst-case O/H for 1Gb/s NIC:

  O/H = 6k cy/pkt * 833k pkt/s = 0.5

Assume 3GHz clock:

  rel O/H = 0.5G cy/s / 3G cy/s = 17% of core

At 100Mb/s, packet spacing = 1/(83k/s) = 12μs

  rel O/H = 17%/10 = 1.7% of core

> Highly pessimistic due to natural batching!

> Avoid batching by spacing packets!

# Performance: i.MX8M, 1Gb/s Eth, UDP



Single-core configuration

# Performance: Processing Cost per Byte

# Performance: Round-Trip Times

# Performance: i.MX8M, 1Gb/s Eth, UDP



Multicore configuration

# Performance: x86, 10Gb/s Eth, UDP



Single-core configuration

# Performance: x86, 10Gb/s Eth, UDP



Multicore configuration

UNSW
SYDNEY

# Syscall cost simulations (x86)



Legend:
- LionsOS (green dashed)
- "free" IPC (black dashed)

Subtract 500cy per kernel entry

Single-core configuration

Y-axis: CPU Utilisation (%) — 0, 20, 40, 60, 80
X-axis: Applied Load (Gb/s) — 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10

UNSW
SYDNEY

# Why Is LionsOS Faster Than Linux?

**Linux:**
- NW driver: 3k lines
- NW system total: 1M lines

**Microkernel overheads are in the noise!**

**LionsOS executes less code!**

**LionsOS:**
- NW driver: 400 lines
- Virtualiser: 160 lines
- Copier: 80 lines
- IP stack: much simpler, client library
- shared NW system total < 1,000 lines

UNSW
SYDNEY

# Why Is LionsOS So Simple?

Helps development **and** correctness!

**Radical simplicity:**

- Fine-grained modularity, strict separation of concerns
- Event-driven programming model, strictly sequential modules
- Static architecture
- Use-case-specific policies

Concurrency by distributing modules across cores

Matches embedded space – little dynamic resource management

Use-case diversity by replacing components

UNSW SYDNEY

# But I Want A **Real** OS!

# Cost Of A Dynamic OS

- More complexity, larger code size

  *Might affect cache footprint?*

- Double book-keeping, multiple server invocations

  *IPC overheads in the noise*

- Higher startup times due to dynamic resource allocation

  *fork() will be the test!*

- Resource revocation may require indirection

  *seL4 caps can be revoked without*

- "Universal" policies are complex & costly

  *Do we need them?*

UNSW SYDNEY

# Do We Need "Universal" Policies?

**Claim:**
- Systems rarely change policies on-the-fly
- Can change policy by replacing policy module

Keep configuration complexity off-line!

**LionsOS experiment:**
- Reload component with new policy implementation
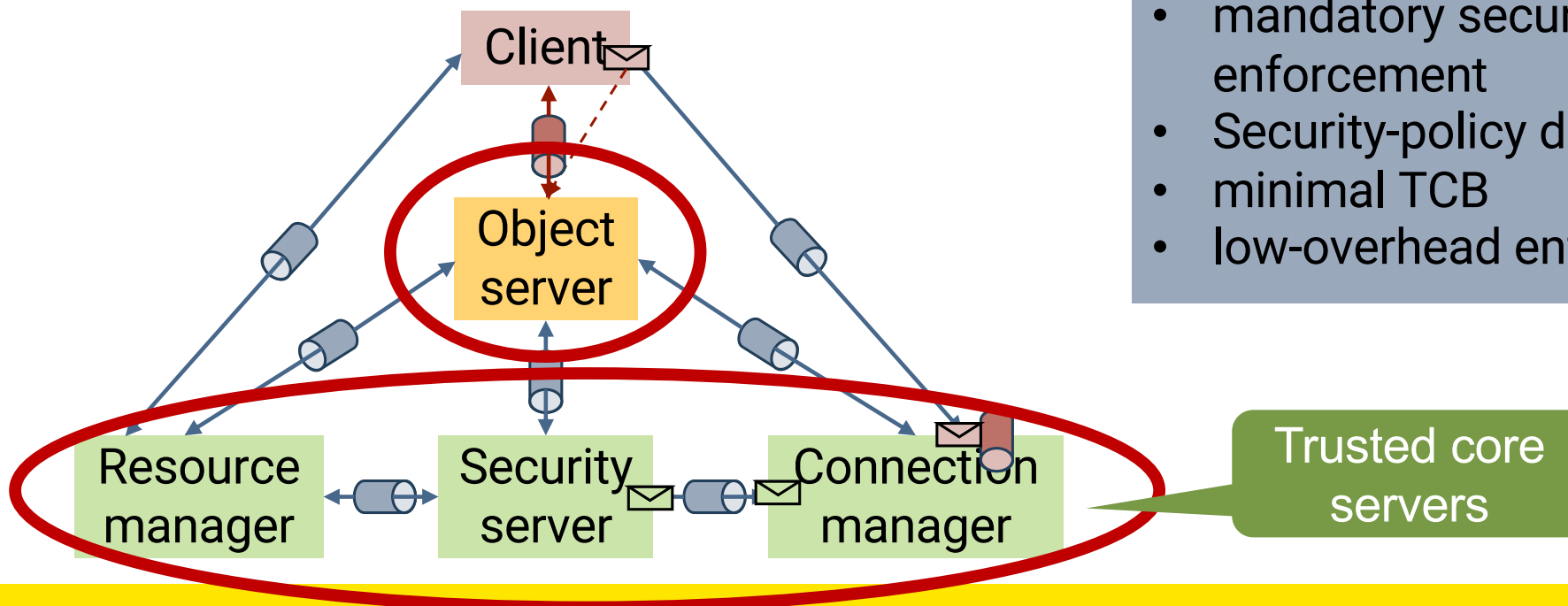- Cost: **17µs** on i.MX8M

# Djawula: PoC Of General-Purpose OS

**Aim:** General-purpose OS that **provably** enforces a general security policy

**Requires:**
- mandatory security-policy enforcement
- Security-policy diversity
- minimal TCB
- low-overhead enforcement

Client

Object server

Resource manager

Security server

Connection manager

Trusted core servers

UNSW
SYDNEY

# Core Ideas: Dynamic Enforcement

Invoke for security-sensitive operations

Consult if no cached decision

Consult & update for dynamic policy

**Client**

**Object Server**

Policy Enforce-ment

Service

Decision making

**Security Server**

Security State

Policy store

UNSW
SYDNEY

# Core Ideas: Resource Donation



**Request stateful service**

Client$_1$

**Client$_x$ state**

**Client$_1$ state**

Object Server

**Early days:**
- Working on framework, details of model
- Adopt components from LionsOS

Rc_O   Rc_C   Rc_x

Resource Manager

Resource containers

UNSW SYDNEY

# My View Of "Changing The Kernel"

# "Changing The Kernel"

| 40M SLoC of C | Module – separate address space | Module – separate address space |
|---|---|---|
| | Module – separate address space | |

**User**
············································

**Kernel / Hypv.**

sel4 Microkernel/Hypervisor

- ✓ Protects **all** modules from each other
- ✓ Requires no special HW
- ✓ Verified kernel
- ✓ Retain virtualisation support
- ✓ **Cost is in the noise with the right design!**

UNSW SYDNEY

# https://trustworthy.systems



**We're hiring!**
Operating-systems faculty