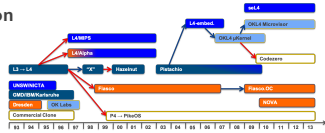


2019 T2 Week 05b
Microkernel Design & Implementation
The 25-year quest for the right API
@GernotHeiser



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>

L4 Microkernels – Deployed by the Billions



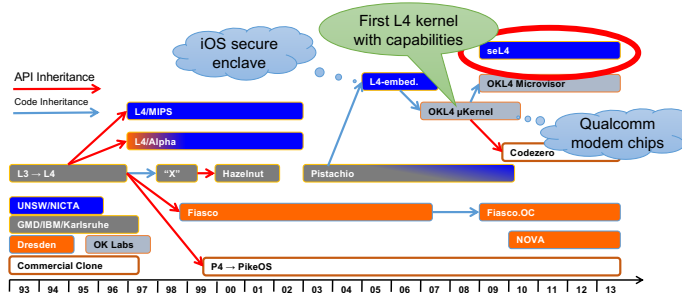
L4: The Quest for a Real Microkernel

L4: The Quest for a Real Microkernel



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. [Liedtke, SOS'95]

L4: 25 Years High Performance Microkernels



L4 IPC Performance Over the Years

Name	Year	Processor	MHz	Cycles	µs
Original	1993	i486	50	250	5.00
Original	1997	Pentium	160	121	0.75
L4/MIPS	1997	R4700	100	86	0.86
L4/Alpha	1997	21064	433	45	0.10
Hazelnut	2002	Pentium 4	1,400	2,000	1.38
Pistachio	2005	Itanium	1,500	36	0.02
OKL4	2007	XScale 255	400	151	0.64
NOVA	2010	i7 Bloomfield (32-bit)	2,660	288	0.11
seL4	2013	ARM11	532	188	0.35
seL4	2018	i7 Haswell (64-bit)	3,400	442	0.13
seL4	2018	Cortex A9	1,000	303	0.30

6

COMP9242 2019T2 W05b: Microkernel D&I

© Gemot Heiser 2019 - CC Attribution License



Minimality: Source-Code Size

Name	Architecture	C/C++	asm	total kSLOC
Original	i486	0	6.4	6.4
L4/Alpha	Alpha	0	14.2	14.2
L4/MIPS	MIPS64	6.0	4.5	10.5
Hazelnut	x86	10.0	0.8	10.8
Pistachio	x86	22.4	1.4	23.0
L4-embedded	ARMv5	7.6	1.4	9.0
OKL4 3.0	ARMv6	15.0	0.0	15.0
Fiasco.OC	x86	36.2	1.1	37.6
seL4	ARMv6	9.7	0.5	10.2

7

COMP9242 2019T2 W05b: Microkernel D&I

© Gemot Heiser 2019 - CC Attribution License



What Have We Learnt in 25 Years?

8

COMP9242 2019T2 W05b: Microkernel D&I

© Gemot Heiser 2019 - CC Attribution License



Issues With 2G Microkernels

- L4 solved microkernel performance [Härtig et al, SOSP'97] left a number of issues unsolved
- Problem: ad-hoc approach to security and resource management
 - Global thread name space \Rightarrow covert channels [Shapiro'03]
 - Threads as IPC targets \Rightarrow insufficient encapsulation
 - Single kernel memory pool \Rightarrow DoS attacks
 - No delegation of authority \Rightarrow impacts flexibility, performance
 - Unprincipled management of time

Solved by capabilities

9

COMP9242 2019T2 W05b: Microkernel D&I

© Gemot Heiser 2019 - CC Attribution License

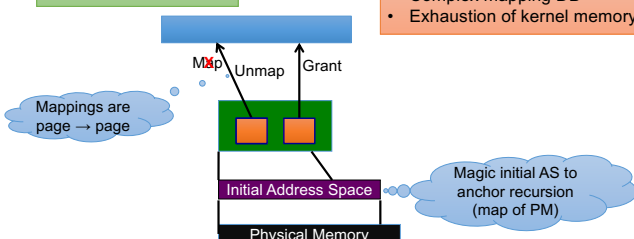


Traditional L4: Recursive Address Spaces

Replaced by magic-free seL4 resource model

Issues:

- Complex mapping DB
- Exhaustion of kernel memory



10

COMP9242 2019T2 W05b: Microkernel D&I

© Gemot Heiser 2019 - CC Attribution License



Issues With 2G Microkernels

- L4 solved microkernel performance [Härtig et al, SOSP'97] left a number of issues unsolved
- Problem: ad-hoc approach to security and resource management
 - Global thread name space \Rightarrow covert channels [Shapiro'03]
 - Threads as IPC targets \Rightarrow insufficient encapsulation
 - Single kernel memory pool \Rightarrow DoS attacks
 - No delegation of authority \Rightarrow impacts flexibility, performance
 - Unprincipled management of time

Solved by seL4 memory management model

11

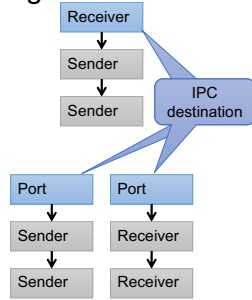
COMP9242 2019T2 W05b: Microkernel D&I

© Gemot Heiser 2019 - CC Attribution License

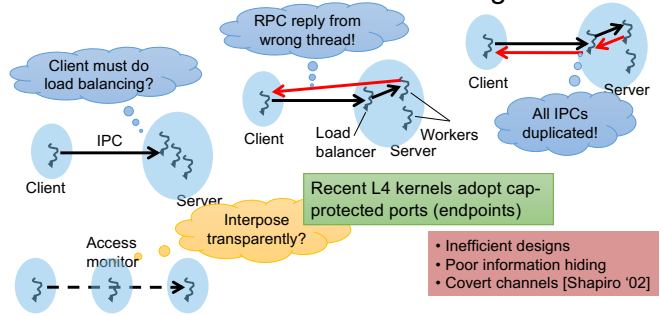


Direct vs Indirect IPC Addressing

- Direct: Queue senders/messages at receiver
 - Need unique thread IDs
- Kernel guarantees identity of sender
 - useful for authentication
- Indirect: Mailbox/port object
 - Just a user-level handle for the kernel-level queue
 - Extra object type – extra weight?
 - Communication partners are anonymous
 - Need separate mechanism for authentication



Other Issues with L4 IPC Addressing



Issues With 2G Microkernels

- L4 solved microkernel performance [Härtig et al, SOS'97] left a number of issues unsolved
- Problem: ad-hoc approach to security and resource management
 - Global thread name space \Rightarrow covert channels [Shapiro'03]
 - Threads as IPC targets \Rightarrow insufficient encapsulation
 - Single kernel memory pool \Rightarrow DoS attacks
 - No delegation of authority \Rightarrow impacts flexibility, performance
 - **Unprincipled management of time**

Solved by caps & endpoints

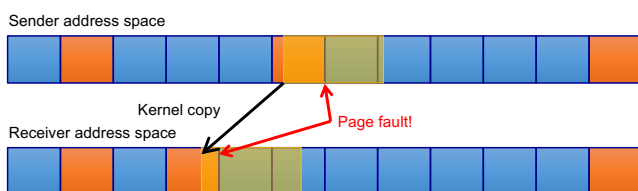
Examine later

Other Design & Implementation Issues

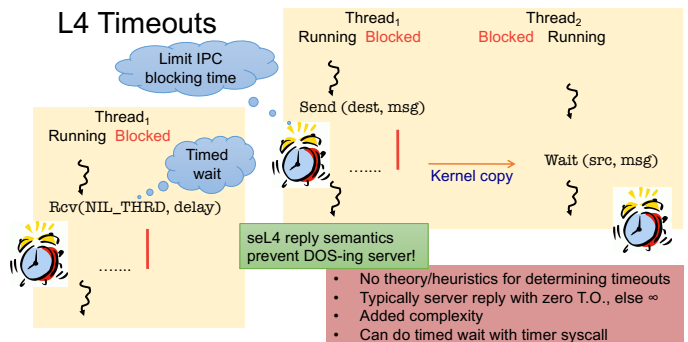
L4 "Long" IPC

Abandoned in seL4

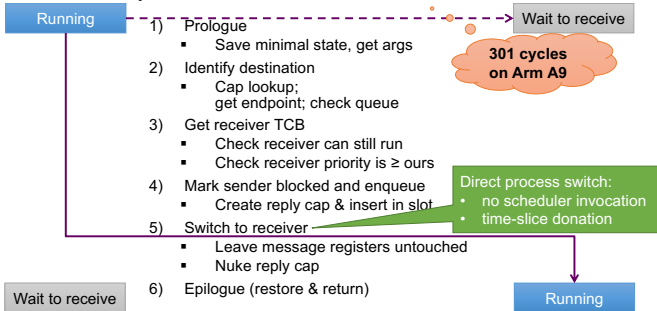
- Not minimal
- Source of kernel complexity:
 - nested exceptions
 - concurrency in kernel
 - must upcall PF handlers during IPC
 - timeouts to prevent DOS attacks



L4 Timeouts



IPC Fastpath: Send Phase of Call



Fastpath Coding Tricks

```
slow = cap_get_capType(en_c) != cap_endpoint_cap ||
!cap_endpoint_cap_get_capCanSend(en_c);
if (slow) enter_slow_path();
```

- Reduces branch-prediction footprint
- Avoids mispredicts, stalls & flushes
- Uses ARM instruction predication
- But: increases slow-path latency (slightly)
 - should be minimal compared to basic slow-path cost

Common case: 0
Common case: 1

How About Real-Time Support?

- Kernel runs with interrupts disabled
 - No concurrency control \Rightarrow simpler kernel
 - Easier reasoning about correctness
 - Better average-case performance

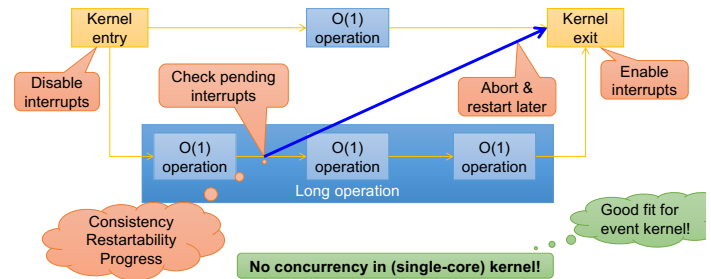
How about long-running system calls?

Most protected-mode RTOSes are fully preemptible



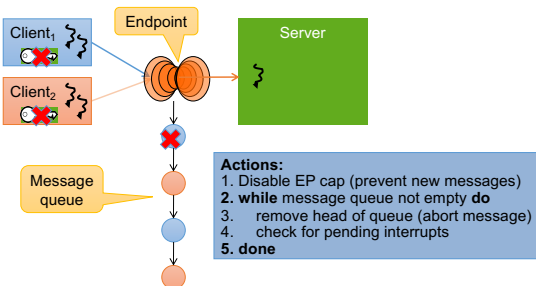
Lots of concurrency in kernel!

sel4 Incremental Consistency Paradigm

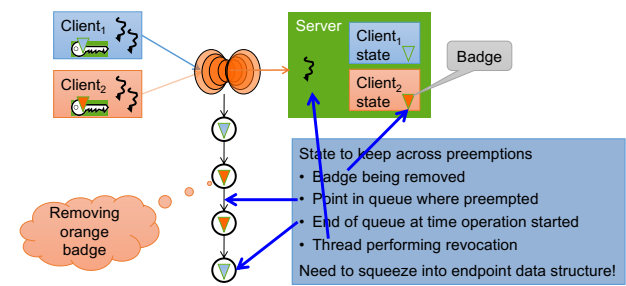


No concurrency in (single-core) kernel!

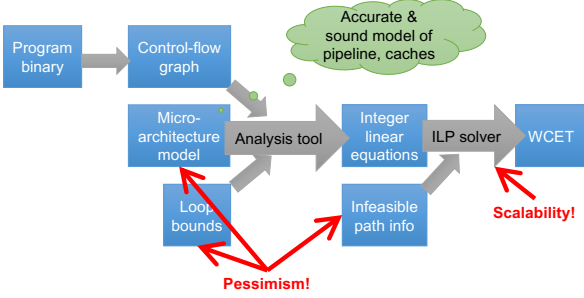
sel4 Example: Destroying IPC Endpoint



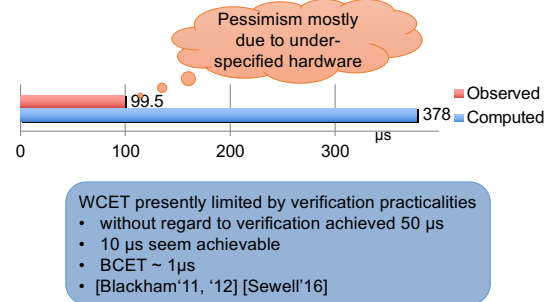
sel4 Difficult Example: Revoking Badge



WCET Analysis



sel4 WCET Analysis on ARM11



L4 Scheduler Optimisation: Lazy Scheduling

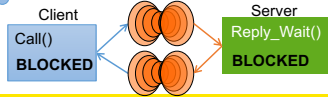
```

thread_t schedule() {
  foreach (prio in priorities) {
    foreach (thread in runQueue[prio]) {
      if (isRunnable(thread))
        return thread;
      else
        schedDequeue(thread);
    }
  }
  return idleThread;
}
  
```

Problem: Unbounded scheduler execution time!

Idea: leave blocked threads in ready queue, scheduler cleans up

- Frequent blocking/unblocking in IPC-based systems
- Many ready-queue manipulations



sel4 Scheduler: Benno Scheduling

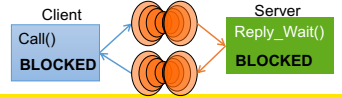
```

thread_t schedule() {
  foreach (prio in priorities) {
    foreach (thread in runQueue[prio]) {
    if (thread=head(runQueue[prio]))
      return thread;
    else
      schedDequeue(thread);
  }
  return idleThread;
}
  
```

Only current thread needs fixing up at preemption time!

Idea: Lazy on unblocking instead on blocking

- Frequent blocking/unblocking in IPC-based systems
- Many ready-queue manipulations



Scheduler Optimisation: Direct Process Switch

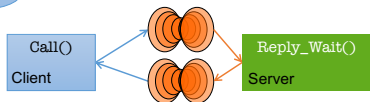
- Sender was running \Rightarrow had highest prio
- If receiver prio \geq sender prio \Rightarrow run receiver

- Arguably, sender should donate back if it's a server replying to a Call()
- Hence, always donate on Reply_Wait()

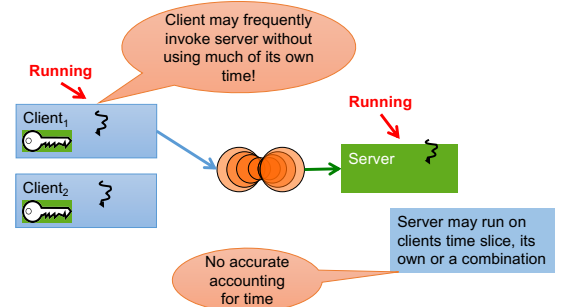
Implication: Time slice donation – receiver runs on sender's time slice

Idea: Don't invoke scheduler if you know who'll be chosen

- Frequent context switches in IPC-based systems
- Many scheduler invocations



Remember: Delegation of Critical Sections



sel4 New Model: Scheduling Contexts

Classical thread attributes

- Priority
- Time slice

New thread attributes

- Priority
- Scheduling context capability

Not runnable if null

Presently being merged into mainline

Limits CPU access!



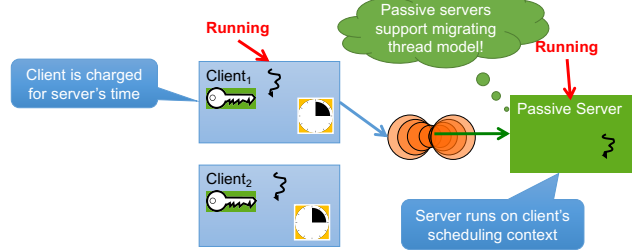
Scheduling context object

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

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

Capability for time

sel4 Delegation with Scheduling Contexts



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

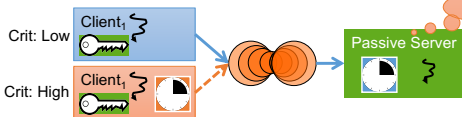
sel4 Mixed-Criticality Support

For supporting **mixed-criticality systems (MCS)**, OS must provide:

- **Temporal isolation**, to force jobs to adhere to declared WCET
- Mechanisms for **safely sharing resources** across criticalities

Solved by scheduling contexts

What if budget expires while shared server executing on Low's scheduling context?



sel4 Timeout Exceptions

Policy-free mechanism for dealing with budget depletion

Possible actions:

- Provide emergency budget to leave critical section
- Cancel operation & roll-back server
- Reduce priority of low-crit client (together with one of the above)
- Implement priority inheritance (if you must...)

Issues With 2G Microkernels

- L4 solved microkernel performance [Härtig et al, SOSP'97] left a number of issues unsolved
- Problem: ad-hoc approach to security and resource management
 - Global thread name space \Rightarrow covert channels [Shapiro'03]
 - Threads as IPC targets \Rightarrow insufficient encapsulation
 - Single kernel memory pool \Rightarrow DoS attacks
 - No delegation of authority \Rightarrow impacts flexibility, performance
- Unprincipled management of time

Solved by scheduling contexts & time-out exceptions

Lessons & Principles

Original L4 Design and Implementation

Implement. Tricks [SOSP'93]

- ~~Process kernel~~
- ~~Virtual TCB array~~ Modified
- ~~Lazy scheduling~~
- ~~Direct process switch~~
- ~~Non-preemptible~~ Retained
- ~~Non-portable~~
- ~~Non-standard calling convention~~
- ~~Assembler~~

Design Decisions [SOSP'95]

- ~~Synchronous IPC~~
- ~~Rich message structure, arbitrary out-of-line messages~~
- ~~Zero-copy register messages~~
- ~~User-mode page-fault handlers~~
- ~~Threads as IPC destinations~~
- ~~IPC timeouts~~
- ~~Hierarchical IPC control~~
- ~~User-mode device drivers~~
- ~~Process hierarchy~~
- ~~Recursive address space construction~~

Reflecting on Changes

Original L4 design had two major shortcomings:

- Insufficient/impractical resource control
 - Poor/non-existent control over kernel memory use
 - Inflexible & costly process hierarchies (policy!)
 - Arbitrary limits on number of address spaces and threads (policy!)
 - Poor information hiding (IPC addressed to threads)
 - Insufficient mechanisms for authority delegation
- Over-optimised IPC abstraction, mangles:
 - Communication
 - Synchronisation
 - Memory management – sending mappings
 - Scheduling – time-slice donation

sel4 Design Principles

- Fully delegatable access control
 - All resource management is subject to user-defined policies
 - Applies to kernel resources too!
 - Performance on par with best-performing L4 kernels
 - Prerequisite for real-world deployment!
 - Suitability for real-time use
 - Important for safety-critical systems
 - Suitable for *formal verification*
 - Requires small size, avoid complex constructs
- Largely in line with traditional L4 approach!

A Thirty-Year Dream!

Specification and Verification of the UCLA Unix¹ Security Kernel

Benno J. Wolk, Richard A. Kemmerer, and Gerald J. Popek
University of California, Los Angeles

This paper reports the specification and verification of the security kernel of the UCLA Unix system. The kernel is a multi-processor system, and its verification is a complex task. The paper describes the development of the kernel, the specification of the security kernel, and the verification of the kernel. The kernel is a multi-processor system, and its verification is a complex task. The paper describes the development of the kernel, the specification of the security kernel, and the verification of the kernel.

Our research seeks to develop means by which an operating system can be shown data secure, meaning that direct access to data must be possible only if the recorded protection policy permits it. The two major components

Communications of the ACM February 1980 Volume 23 Number 2