

## 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:

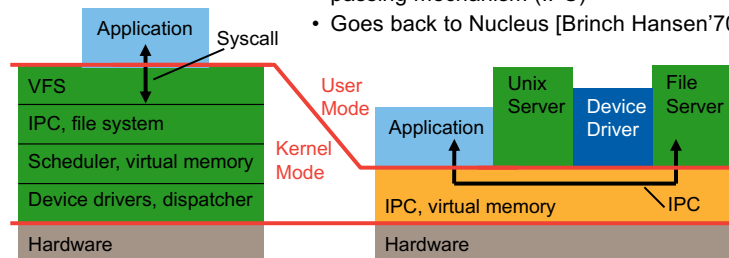
*“Courtesy of Gernot Heiser, UNSW Sydney”*

The complete license text can be found at <http://creativecommons.org/licenses/by/3.0/legalcode>

## Microkernels: Reducing the Trusted Computing Base

- Idea of microkernel:
  - Flexible, minimal platform
  - Mechanisms, not policies
  - OS functionality provided by usermode servers
  - Servers invoked by kernel-provided message-passing mechanism (IPC)
  - Goes back to Nucleus [Brinch Hansen'70]

IPC performance is critical!



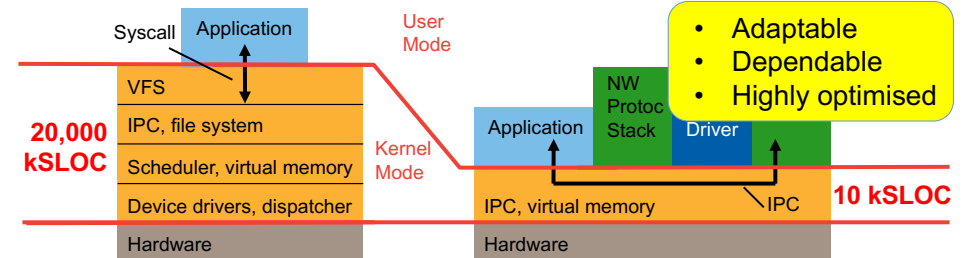
## Monolithic vs Microkernel OS Evolution

### Monolithic OS

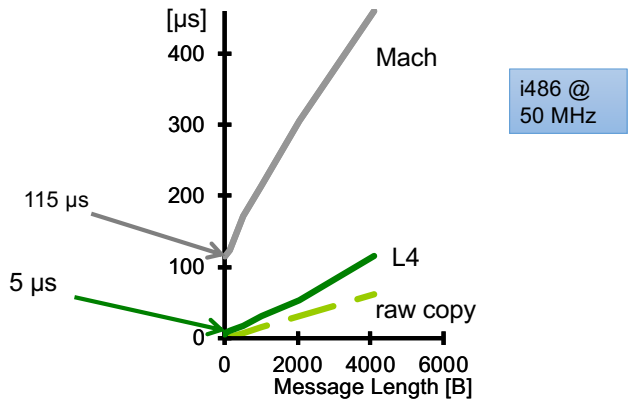
- New features add code kernel
- New policies add code kernel
- Kernel complexity grows

### Microkernel OS

- Features add usermode code
- Policies replace usermode code
- Kernel complexity is stable



# 1993 "Microkernel": IPC Performance



**Culprit: Cache footprint [Liedtke'95]**

i486 @ 50 MHz

# Microkernel Principle: Minimality



*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. [SOSP'95]*

- Advantages of resulting small kernel:
  - Easy to implement, port?
  - Easier to optimise
  - Hopefully enables a minimal *trusted computing base*
  - Easier debug, maybe even *prove* correct?
- Challenges:
  - API design: generality despite small code base
  - Kernel design and implementation for high performance

Limited by arch-specific micro-optimisations

Small attack surface, fewer failure modes

# Microkernel Evolution

**First generation** Mach [’87], QNX, Chorus  
**Second generation** L4 [’95], PikeOS, Integrity  
**Third generation** seL4 [’09]

- Memory Objects
- Low-level FS, Swapping
- Devices
- Kernel memory
- Scheduling
- IPC, MMU abstr.

180 syscalls, 100 kSLOC  
100 μs IPC

- Kernel memory
- Scheduling
- IPC, MMU abstr.

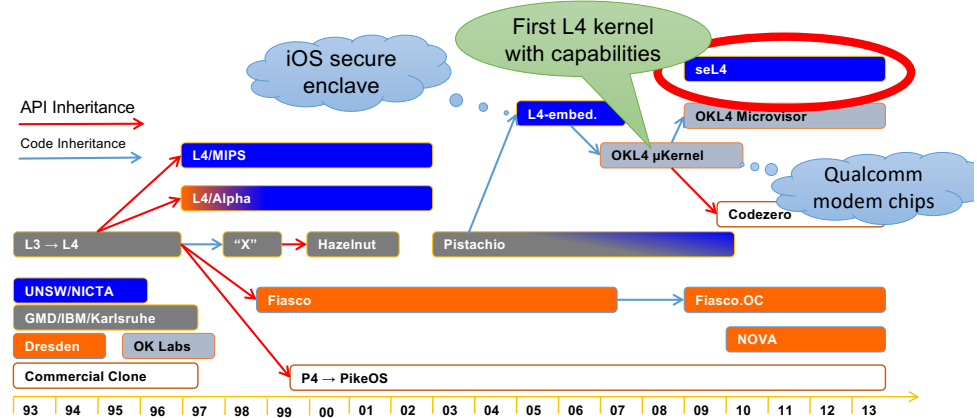
~7 syscalls, ~10 kSLOC  
~ 1 μs IPC

Memory-mangmt library

- Scheduling
- IPC, MMU abstr.

~3 syscalls, ~10 kSLOC  
0.1 μs IPC  
*Capabilities*  
*Design for isolation*

# L4: 25 Years High Performance Microkernels



## Issues With 2G Microkernels

- L4 solved microkernel performance [Härtig et al, SOSP'97]
- Left a number of security issues unsolved
- Problem: ad-hoc approach to protection 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  $\Rightarrow$  limited flexibility, performance issues
  - Unprincipled management of time
- Addressed by seL4
  - Designed to support safety- and security-critical systems
  - Principled time management (MCS branch)

8

COMP9242 2019T2 W01a

© Gernot Heiser 2019 – CC Attribution License



## The seL4 Microkernel

9

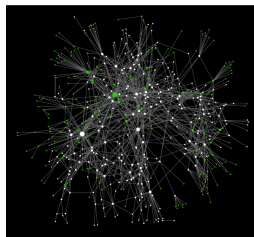
COMP9242 2019T2 W01a

© Gernot Heiser 2019 – CC Attribution License



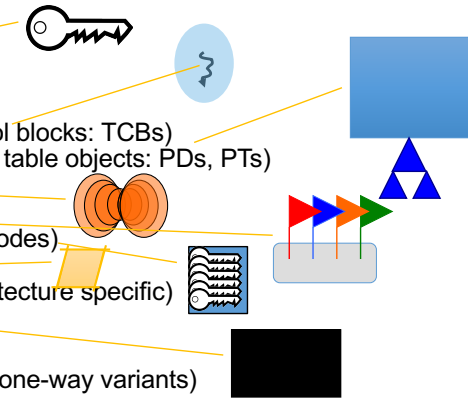
## seL4 Principles

- Single protection mechanism: capabilities
  - Now also for time [Lyons et al, EuroSys'18]
- All resource-management policy at user level
  - Painful to use
  - Need to provide standard memory-management library
    - Results in L4-like programming model
- Suitable for formal verification
  - Proof of implementation correctness
  - Attempted since '70s
  - Finally achieved by L4.verified project at NICTA [Klein et al, SOSP'09]



## seL4 Concepts

- Capabilities (Caps)
  - mediate access
- Kernel objects:
  - Threads (thread-control blocks: TCBs)
  - Address spaces (page table objects: PDs, PTs)
  - Endpoints (IPC)
  - Notifications
  - Capability spaces (CNodes)
  - Frames
  - Interrupt objects (architecture specific)
  - Untyped memory
- System calls:
  - Call, Reply&Wait (and one-way variants)
  - Yield



10

COMP9242 2019T2 W01a

© Gernot Heiser 2019 – CC Attribution License



11

COMP9242 2019T2 W01a

© Gernot Heiser 2019 – CC Attribution License



# seL4 What Are (Object) Capabilities?

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



Obj reference  
Access rights



Object

Eg. thread,  
address space

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

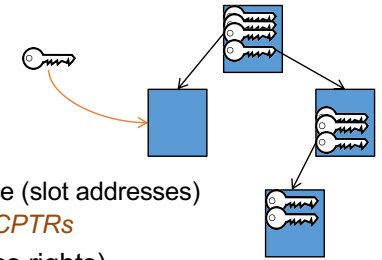
Capabilities provide:

- Fine-grained access control
- Reasoning about information flow

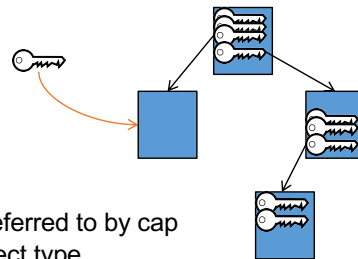
Any system call is invoking a capability:  
`err = cap.method( args );`

# seL4 Capabilities

- Stored in cap space (*Cspace*)
  - Kernel object made up of *CNodes*
  - each an array of cap "slots"
- Inaccessible to userland
  - But referred to by pointers into CSpace (slot addresses)
  - These CSpace addresses are called *CPTRs*
- Caps convey specific privilege (access rights)
  - Read, Write, Execute, Grant (cap transfer)



# Capabilities



- Main operations on caps:
  - Invoke*: perform operation on object referred to by cap
    - Possible operations depend on object type
  - Copy/Mint/Grant*: create copy of cap with *same/lesser* privilege
  - Move/Mutate*: transfer to different address with same/lesser privilege
  - Delete*: invalidate slot (cleans up object if this is the only cap to it)
  - Revoke*: delete any derived (eg. copied or minted) caps

# Cross-Address-Space Invocation (IPC)

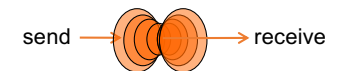
## Fundamental microkernel operation

- Kernel provides no services, only mechanisms
- OS services provided by (protected) user-level server processes
- invoked by IPC

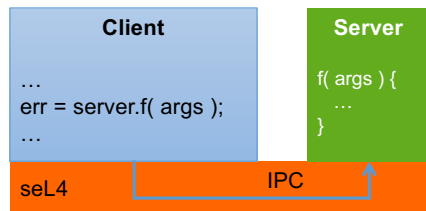


- seL4 IPC uses a handshake through *endpoints*:

- Transfer points without storage capacity
- Message must be transferred instantly
  - Single-copy user → user by kernel



## seL4 IPC: Cross-Domain Invocation



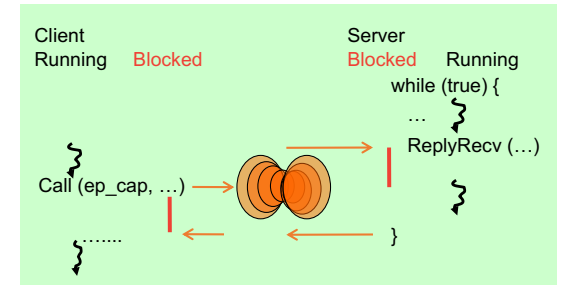
seL4 IPC is **not**:

- A mechanism for shipping data
- A synchronisation mechanism

seL4 IPC is:

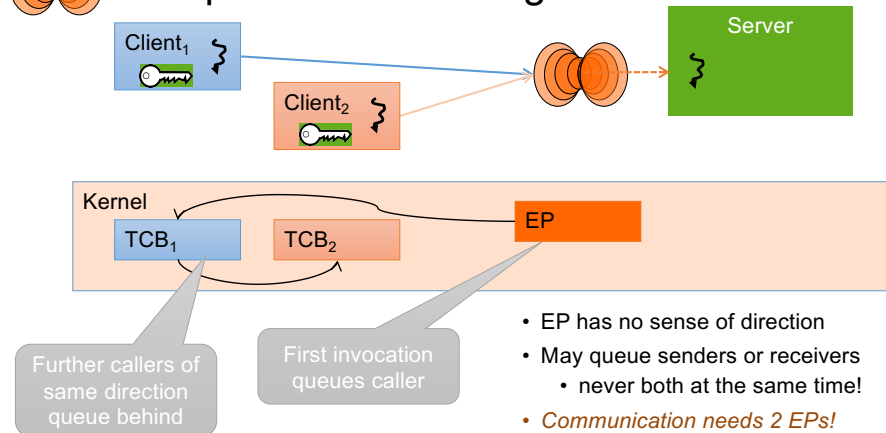
- A protected procedure call
- A user-controlled context switch

## IPC: Endpoints



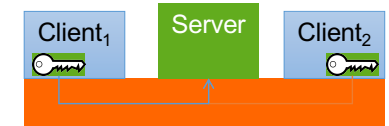
- Threads must rendez-vous
  - One side blocks until the other is ready
  - Implicit synchronisation
- Message copied from sender's to receiver's *message registers*
  - Message is combination of caps and data words
    - Presently max 121 words (484B, incl message “tag”)
    - Should never use anywhere near that much!

## Endpoints are Message Queues



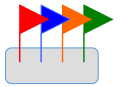
## Server Invocation & Return

- Asymmetric relationship:
  - Server widely accessible, clients not
  - How can server reply back to client (distinguish between them)?
- Client can pass (session) reply cap in first request
  - server needs to maintain session state
  - forces stateful server design
- seL4 solution: Kernel provides single-use *reply cap*
  - only for Call operation
  - allows server to reply to client
  - cannot be copied/minted/re-used but can be moved
  - one-shot (automatically destroyed after first use)



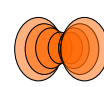
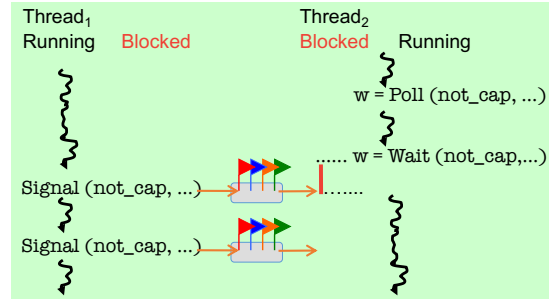
MCS kernel removes the magic



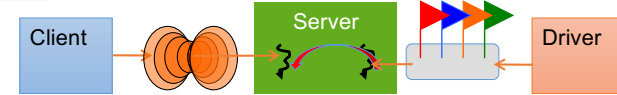


## Notifications

- Logically, a Notification is an array of binary semaphores
  - Multiple signalling, select-like wait
  - Not a message-passing IPC operation!
- Implemented by *data word* in Notification
  - Send OR-s sender's *cap badge* to data word
  - Receiver can poll or wait
    - waiting returns and clears data word
    - polling just returns data word



## Receiving from EP and Notification



- Example: file system
  - synchronous (RPC-style) client protocol
  - asynchronous notifications from driver
- Could have separate threads waiting on endpoints
  - forces multi-threaded server, concurrency control
- Alternative: allow single thread to wait on both events
  - Notification is *bound* to thread with `TCB_BindNotification()`
  - thread waits on Endpoint
  - Notification delivered as if caller had been waiting on Notification

Server with synchronous and asynchronous interface