School of Computer Science & Engineering
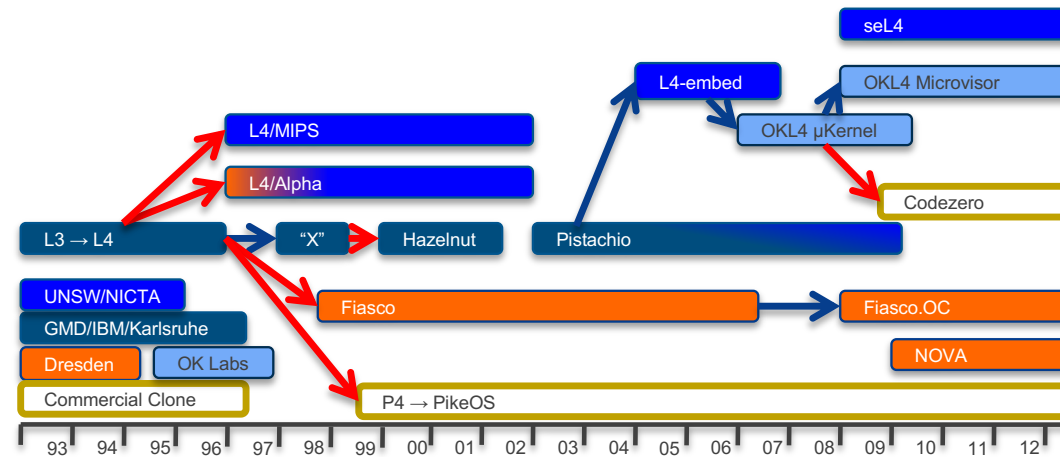
**COMP9242 Advanced Operating Systems**

2020 T2 Week 01a

**Introduction: Microkernels and seL4**

@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:
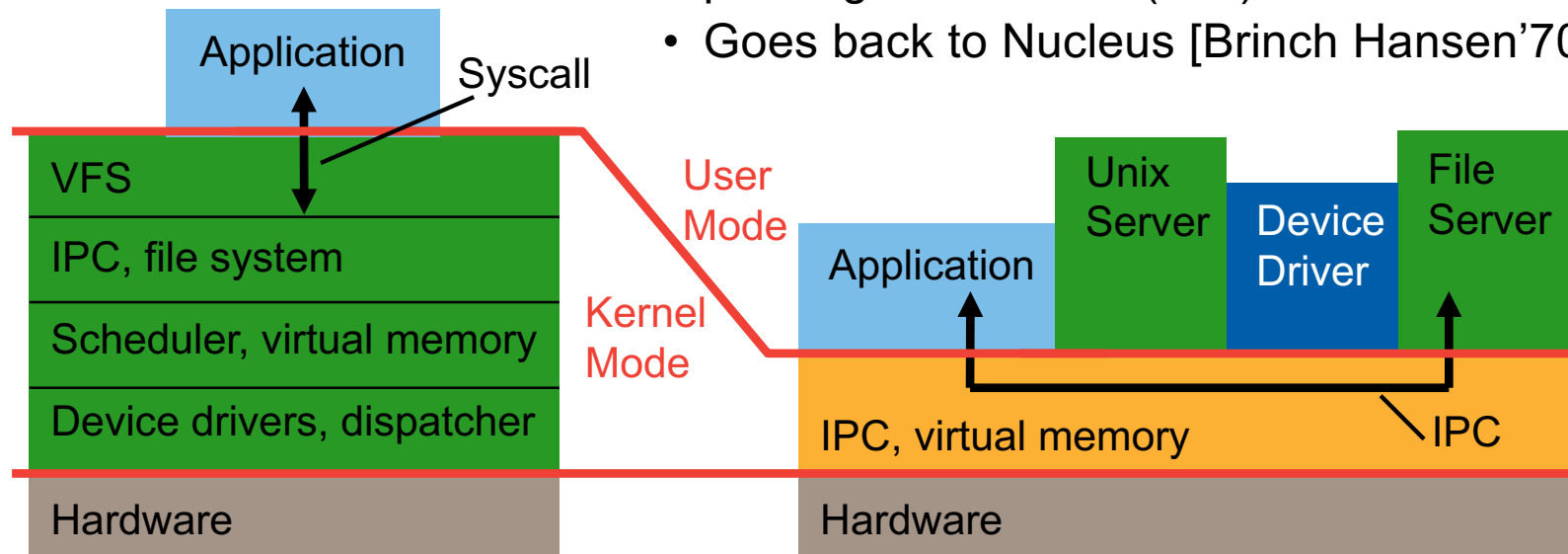
    *"Courtesy of Gernot Heiser, UNSW Sydney"*

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

# Microkernels: Reducing the Trusted Computing Base

IPC performance is critical!

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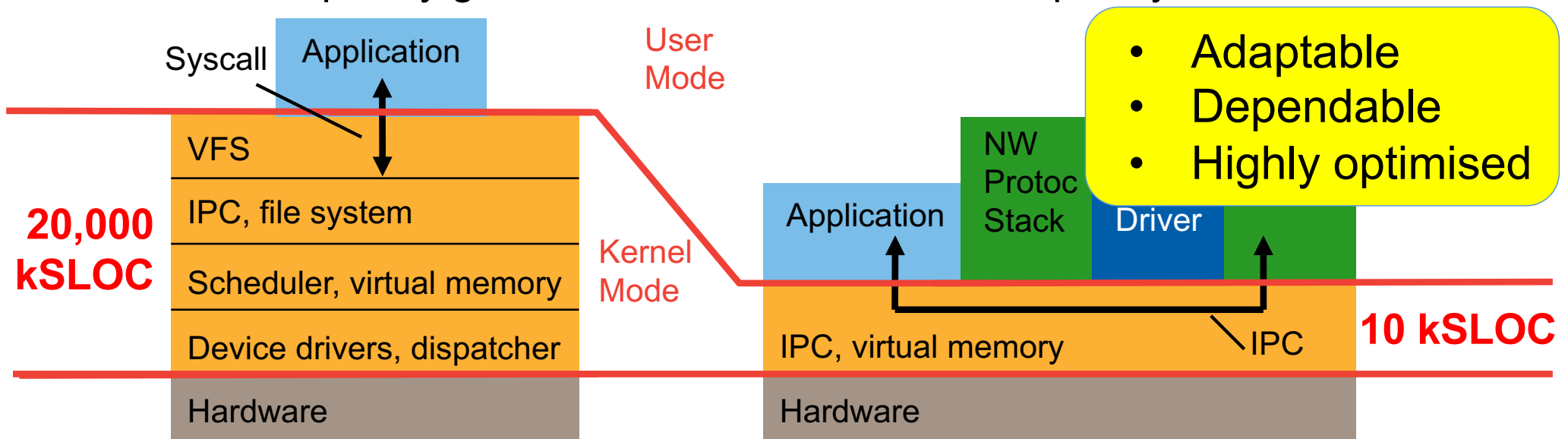  - Goes back to Nucleus [Brinch Hansen'70]

UNSW
SYDNEY

# Monolithic vs Microkernel OS Evolution

## Monolithic OS

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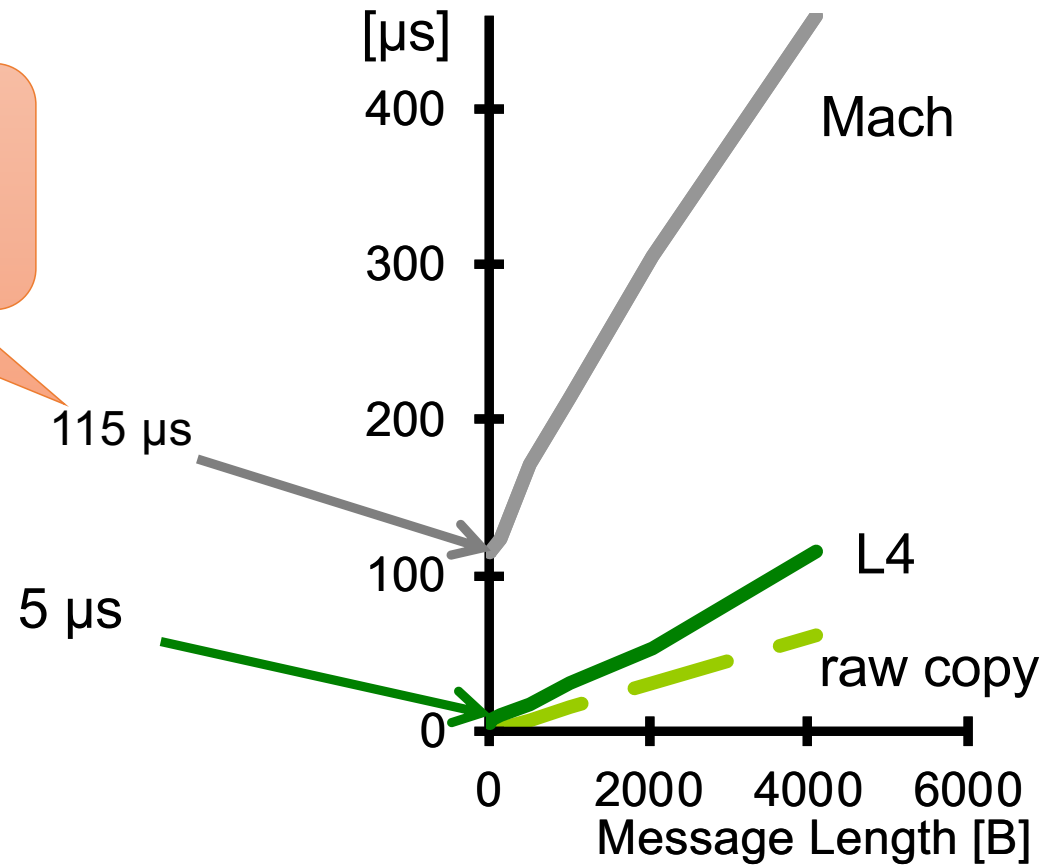- New policies add code kernel
- Kernel complexity grows

## Microkernel OS

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

Syscall

Application

**User Mode**

VFS

IPC, file system

**20,000 kSLOC**

Scheduler, virtual memory

**Kernel Mode**

Device drivers, dispatcher

Hardware

Application

NW Protoc Stack

Driver

- Adaptable
- Dependable
- Highly optimised

IPC, virtual memory

IPC

**10 kSLOC**

Hardware

UNSW SYDNEY

# 1993 "Microkernel": IPC Performance

Culprit:
Cache footprint
[Liedtke'95]

115 µs

5 µs

[µs]

400

300

200

100

0

Mach

i486 @
50 MHz

L4

raw copy

0    2000    4000    6000
Message Length [B]

UNSW
SYDNEY

# 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
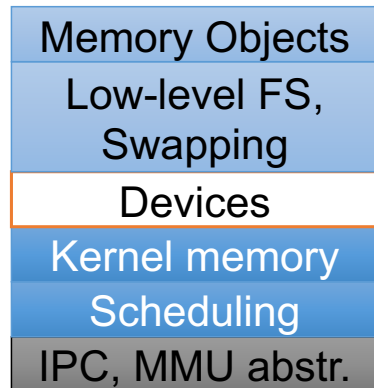  - Kernel design and implementation for high performance

Limited by arch-specific micro-optimisations

Small attack surface, fewer failure modes
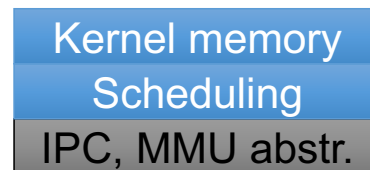
UNSW SYDNEY

# Microkernel Evolution

**First generation**

Mach ['87], QNX, Chorus

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

180 syscalls, 100 kSLOC
100 µs IPC

**Second generation**

L4 ['95], PikeOS, Integrity

| Kernel memory |
| Scheduling |
| IPC, MMU abstr. |

~7 syscalls, ~10 kSLOC
~ 1 µs IPC

**Third generation**

seL4 ['09]

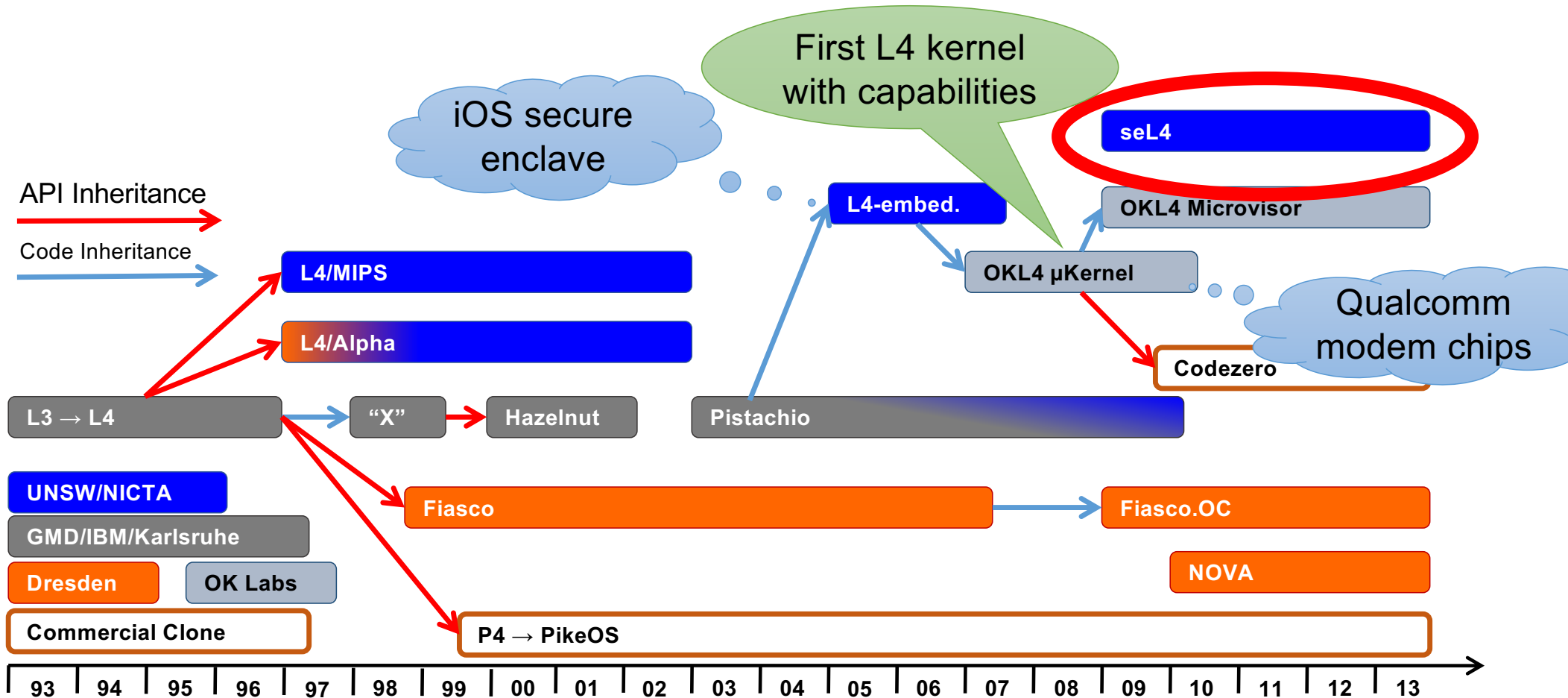| Memory-mangmt library |

| Scheduling |
| IPC, MMU abstr. |

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

UNSW SYDNEY

# L4: 25 Years High Performance Microkernels

# 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 ⇒ covert channels [Shapiro'03]
  - Threads as IPC targets ⇒ insufficient encapsulation
  - Single kernel memory pool ⇒ DoS attacks
  - No delegation of authority ⇒ impacts flexibility, performance
  - Unprincipled management of time

- Addressed by seL4
  - Designed to support safety- and security-critical systems
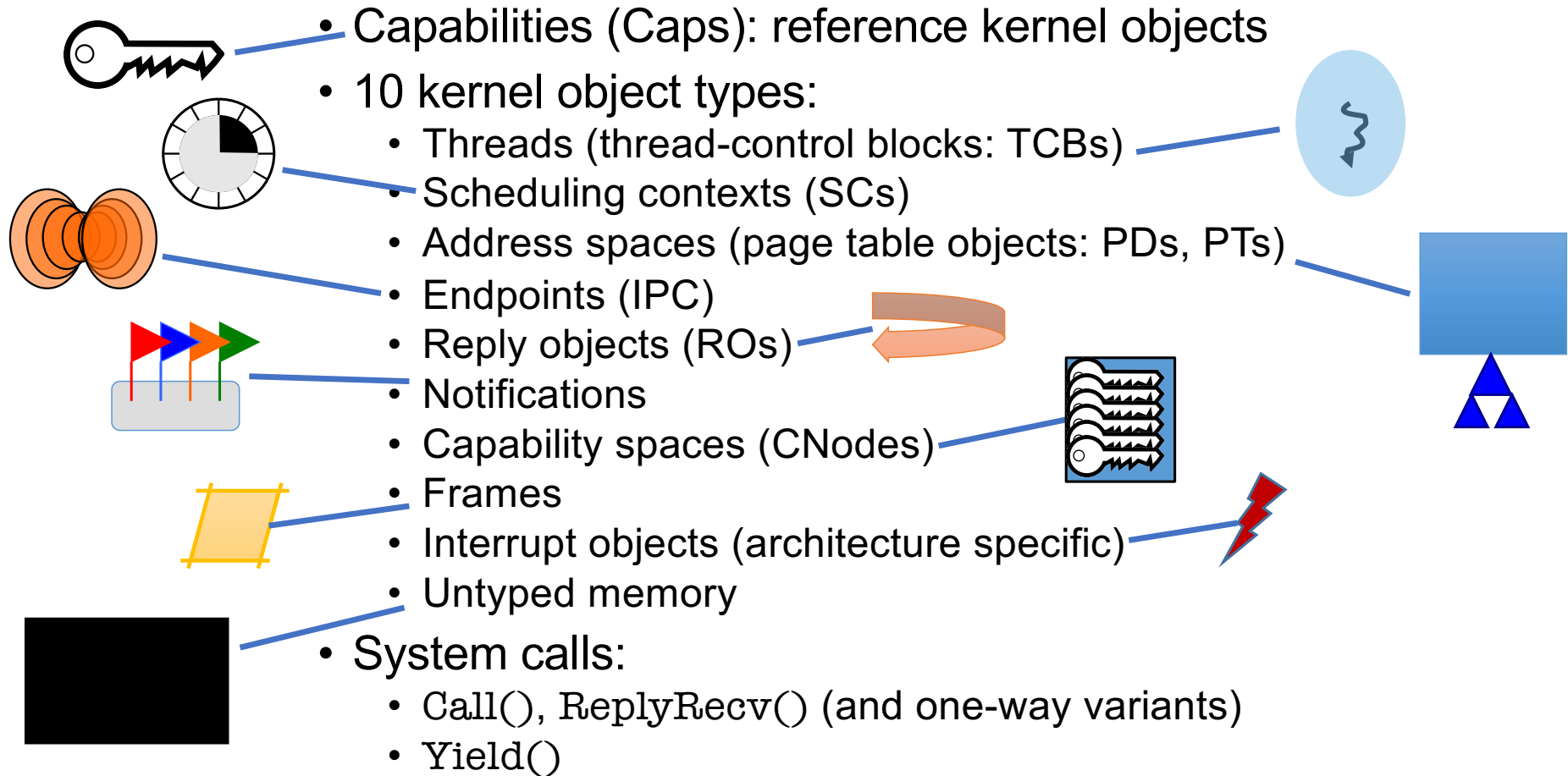  - Principled time management (new MCS configuration)

# The seL4 Microkernel

# Principles

- Single protection mechanism: capabilities
    - Now also for time: MCS configuration [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]

# Concepts in a Slide

- Capabilities (Caps): reference kernel objects
- 10 kernel object types:
  - Threads (thread-control blocks: TCBs)
  - Scheduling contexts (SCs)
  - Address spaces (page table objects: PDs, PTs)
  - Endpoints (IPC)
  - Reply objects (ROs)
  - Notifications
  - Capability spaces (CNodes)
  - Frames
  - Interrupt objects (architecture specific)
  - Untyped memory
- System calls:
  - Call(), ReplyRecv() (and one-way variants)
  - Yield()

UNSW
SYDNEY

# Not a Concept: Hardware Abstraction

**Why?**

- Hardware abstraction violates minimality
- Hardware abstraction introduces policy

**True microkernel:**

- Minimal wrapper of hardware, just enough to safely multiplex
- "CPU driver" [Charles Gray]
- Similarities with Exokernels [Engeler '95]

# What Are (Object) Capabilities?

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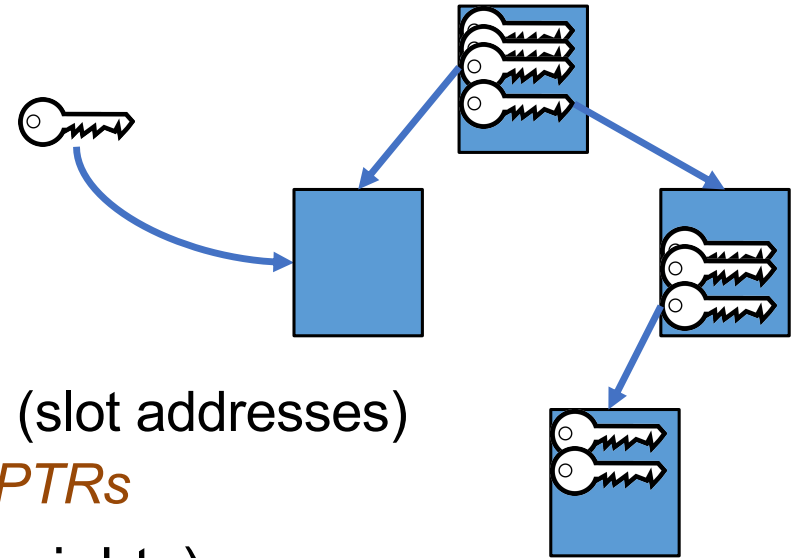Prima-facie evidence of privilege

Object

Eg. thread, address space

Obj reference
Access rights

Eg. read, write, send, execute…

Capabilities provide:
- Fine-grained access control
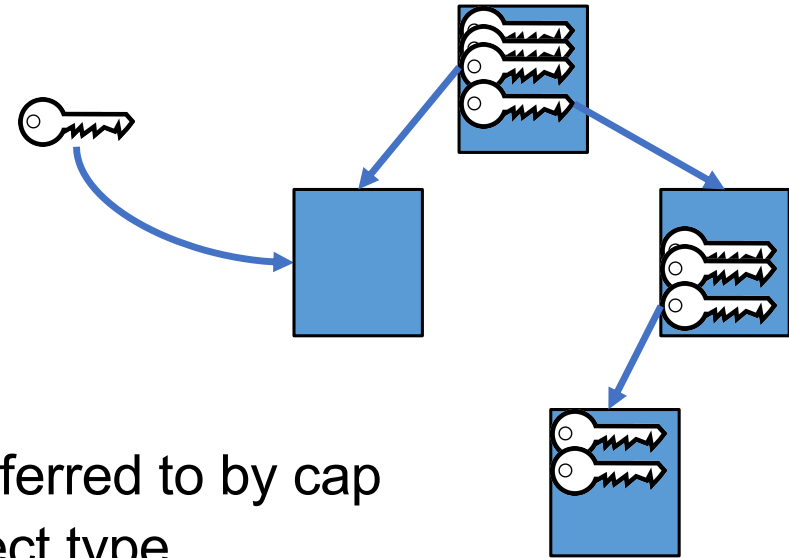- Reasoning about information flow

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

UNSW
SYDNEY

# 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, GrantReply (call), Grant (cap transfer)
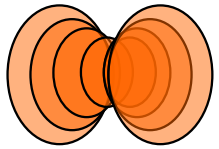
# Capabilities

- 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)
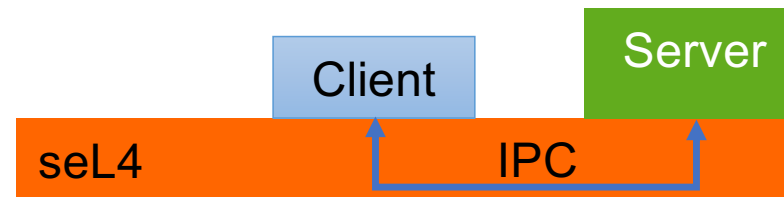  - *Revoke*: delete any derived (eg. copied or minted) caps
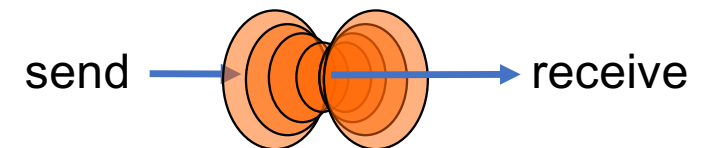
# seL4 Mechanisms

IPC & Notifications

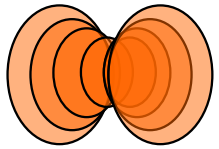# Cross-Address-Space Invocation (IPC)

**Fundamental microkernel operation**

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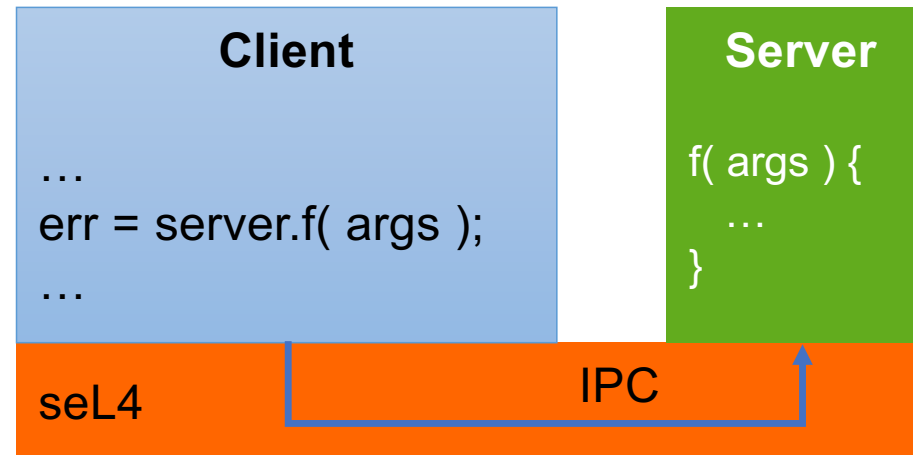- invoked by IPC

- seL4 IPC uses a handshake through *endpoints*:
  - Transfer points without storage capacity
  - Message must be transferred instantly
    - Single-copy user �misc➡ user by kernel

# seL4 IPC: Cross-Domain Invocation

**Client**

…
err = server.f( args );
…

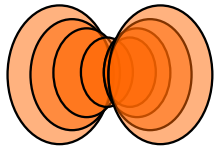**Server**

f( args ) {
…
}

seL4

IPC

seL4 IPC is **not**:
- A mechanism for shipping data
- A synchronisation mechanism
  - side effect, not purpose

seL4 IPC **is**:
- A protected procedure call
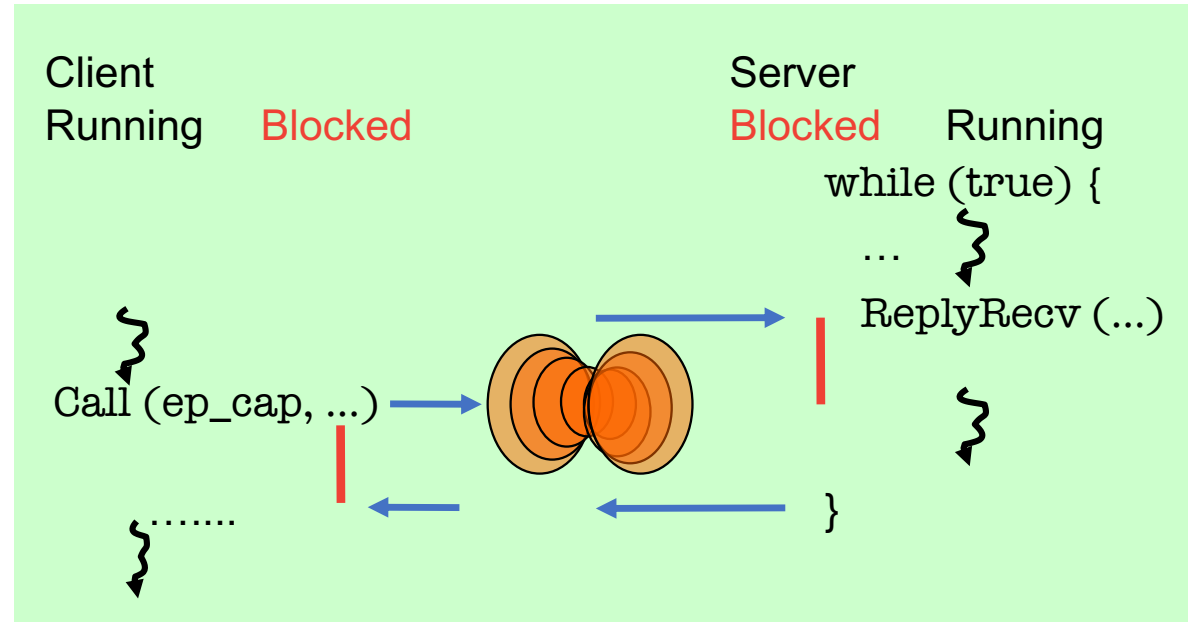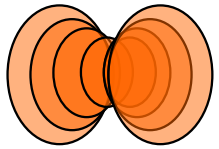- A user-controlled context switch

UNSW
SYDNEY

# IPC: Endpoints

> - Involves 2 threads, but always one blocked
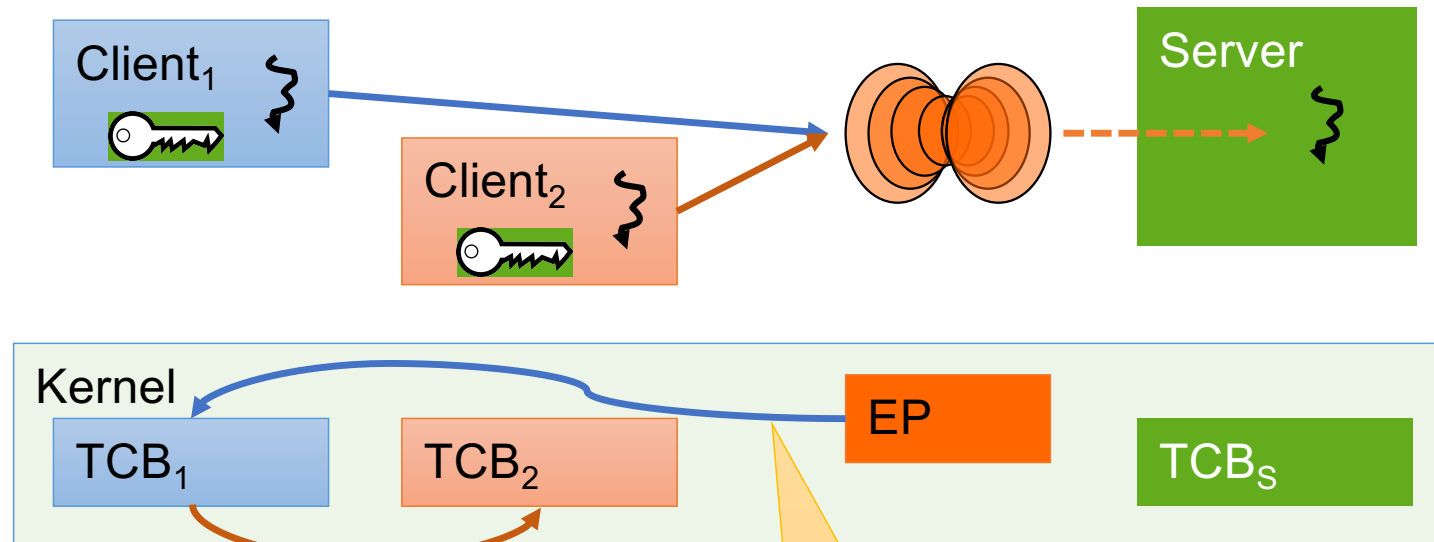> - logically, thread moves between address spaces

- 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")
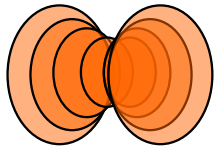    - Should never use anywhere near that much!

```
Client                          Server
Running    Blocked              Blocked    Running
                                         while (true) {
                                           …
                                         ReplyRecv (...)

Call (ep_cap, ...)

        …….                              }
```

UNSW
SYDNEY

# Endpoints are Message Queues

Client$_1$

Client$_2$

Server

Note: Should not normally get queues on a single core, server should have higher priority than clients!

Kernel

TCB$_1$    TCB$_2$    EP    TCB$_S$

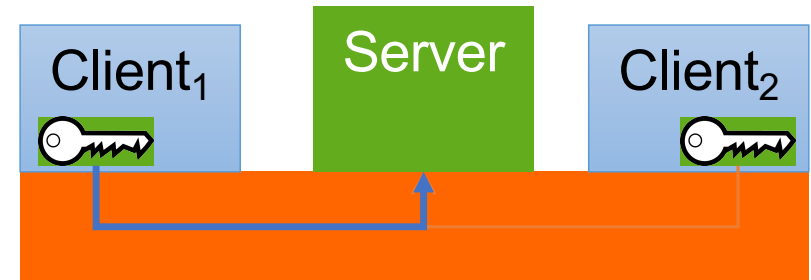Further callers of same direction queue behind

First invocation queues caller

- EP has no sense of direction
- May queue senders or receivers
  - never both at the same time!
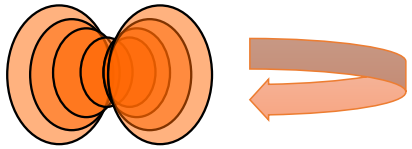- *Communication needs 2 EPs!*

UNSW
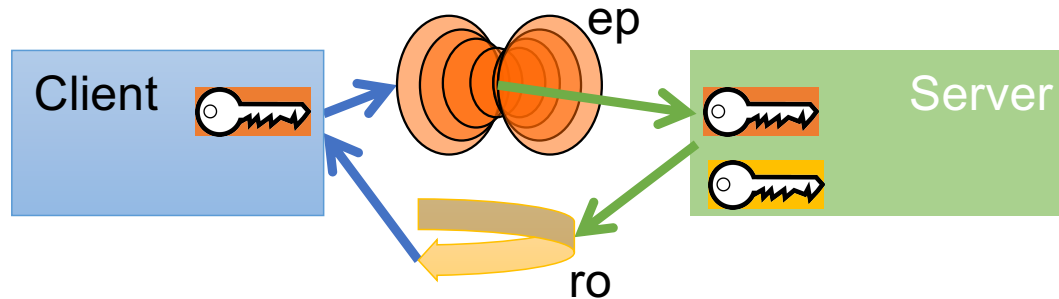SYDNEY

# 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

New MCS kernel semantics!

- seL4 solution: Kernel creates channel in *reply object* (RO)
  - server provides RO in `ReplyRecv()` operation
  - kernel connects RO to client when executing receive phase
  - server invokes RO for send phase (only one send until refreshed)
  - only works when client invokes with `Call()`

UNSW SYDNEY

# Call Semantics



**Client**

Call(ep, args)  →

**Kernel**

*deliver to server*

*block client on RO*  →

*process*

←  *deliver to client*

*process*  ←

**Server**

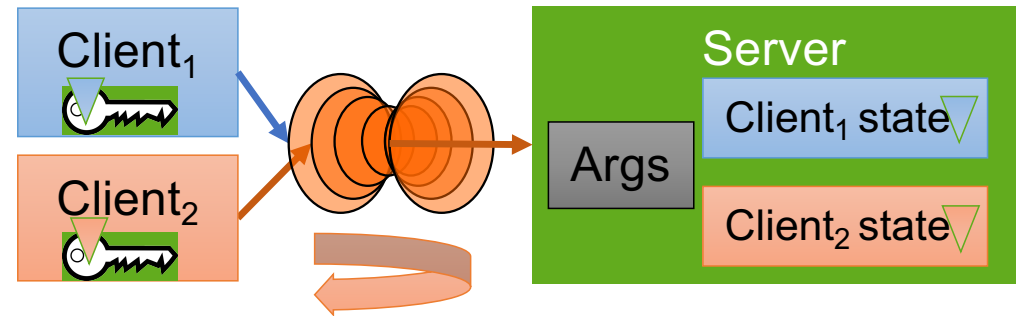ReplyRecv(ro,ep,&args)
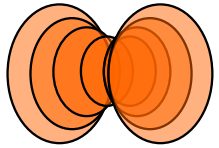
*process*

ReplyRecv(ro,ep,&args)

# Stateful Servers: Identifying Clients
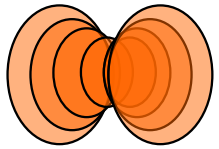
- Server must respond to correct client
  - Ensured by reply cap

- Must associate request
  with correct state



- Could use separate EP per client

  - endpoints are lightweight (16 B)
  - but requires mechanism to wait on a set of EPs (like select)

- Instead, seL4 allows to individually mark ("badge") caps to same EP

  - server provides individually badged (session) caps to clients
    - separate endpoints for opening session, further invocations
  - server tags client state with badge
  - kernel delivers badge to receiver on invocation of badged caps

# IPC Mechanics: Virtual Registers

- Like physical registers, virtual registers are thread state
    - context-switched by kernel
    - implemented as physical registers or thread-local memory
- Message registers
    - contain message transferred in IPC
    - architecture-dependent subset mapped to physical registers
        - 4 on ARM & x64
        - library interface hides details
        - 1$^{st}$ transferred word is special, contains *message tag*
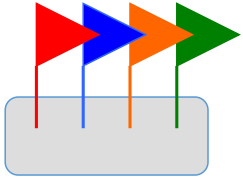    - API MR[0] refers to next word (not the tag!)

# IPC Operations Summary

- `Call (ep_cap, ...)`
  - *Atomic*: guarantees caller is ready to receive reply
  - Sets up server's reply object

- `ReplyRecv (ep_cap, ...)`
  - Invokes RO, waits on EP, re-inits RO

*Not really useful*

- `Recv (ep_cap, ...)`, `Reply(...)`, `Send (ep_cap, ...)`
  - For initialisation and exception handling
  - needs Write, Read permission, respectively

*Need error handling protocol !*

- `NBSend (ep_cap, ...)`
  - Polling send, message lost if receiver not ready
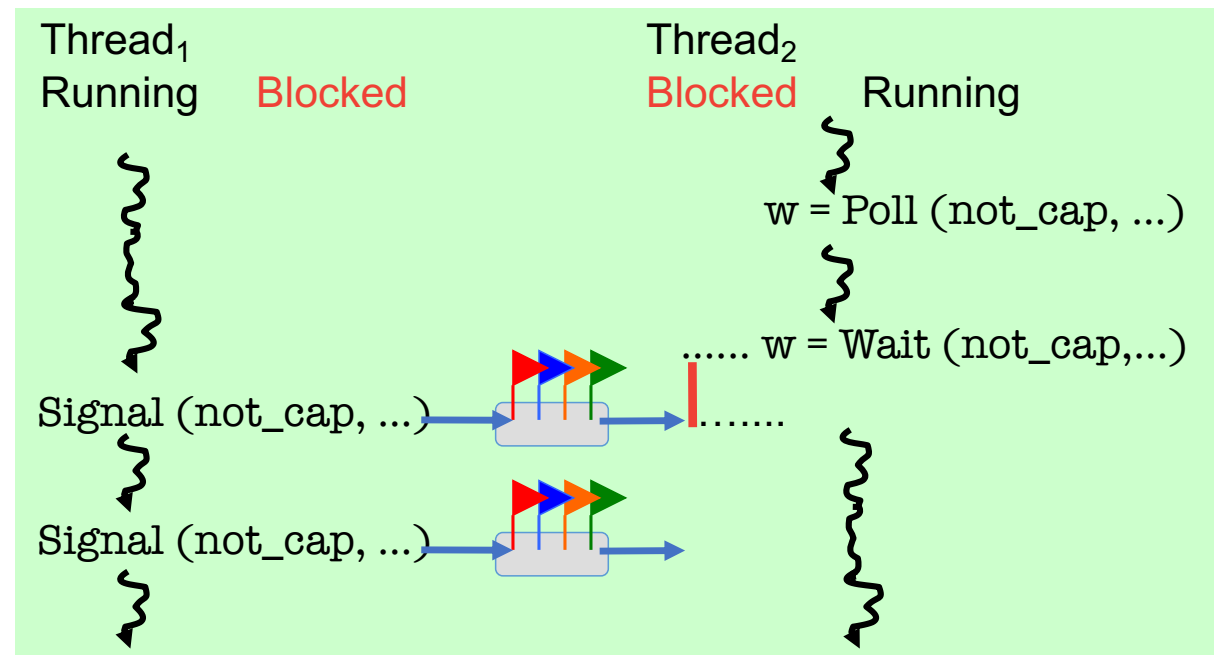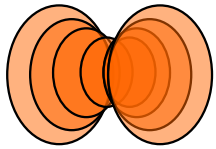
**No failure notification where this reveals info on other entities!**

# Notifications – Synchronisation Objects

- 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

$Thread_1$  Running    Blocked

$Thread_2$  Blocked    Running

w = Poll (not_cap, ...)

...... w = Wait (not_cap,...)

Signal (not_cap, ...)     .......

Signal (not_cap, ...)

UNSW
SYDNEY

# Receiving from EP *and* Notification

**Server with synchronous and asynchronous interface**

Synchronous RPC protocol

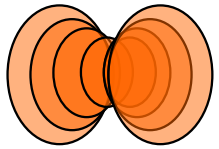Asynchronous completion signals

Client

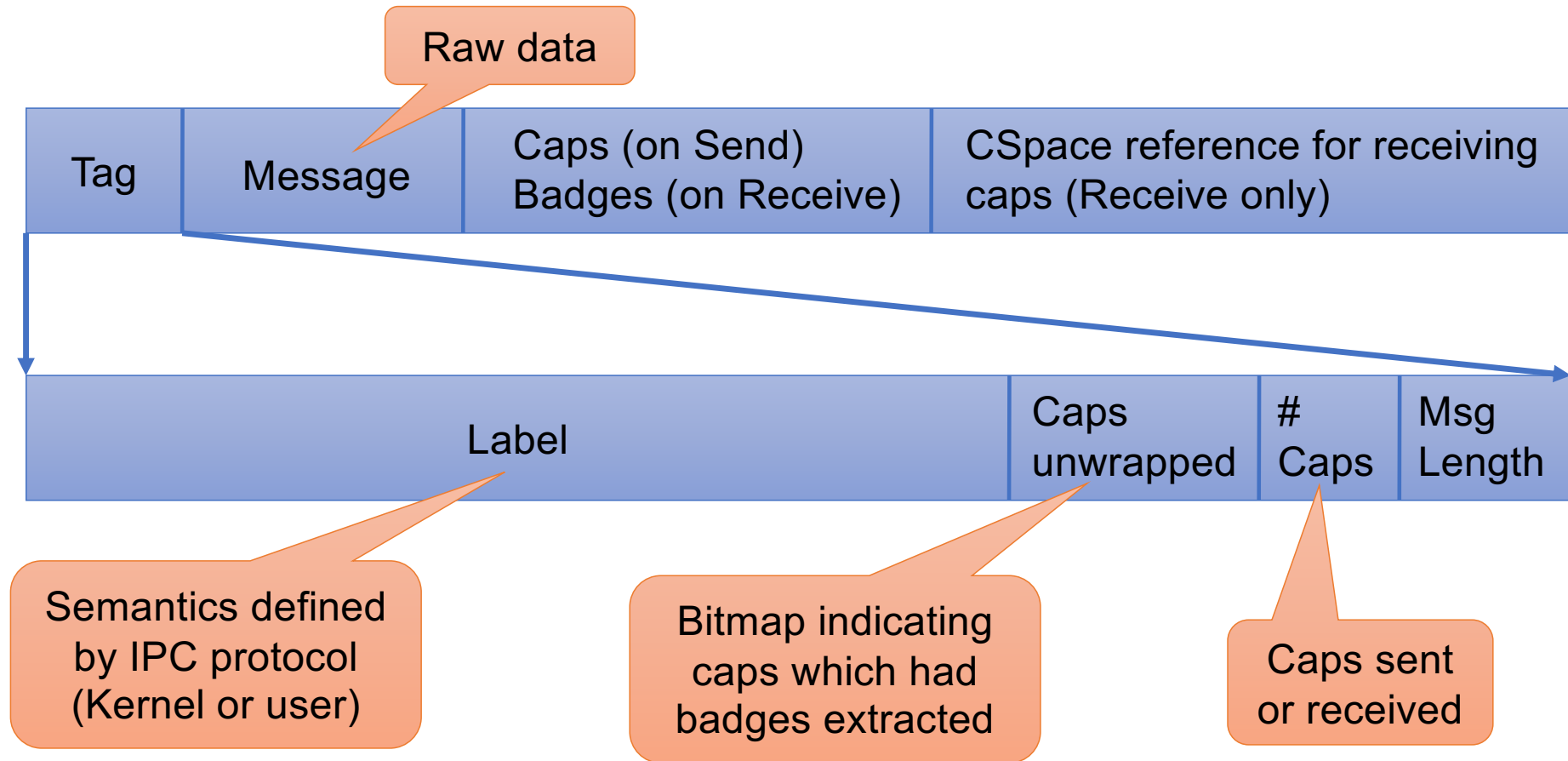Server

Driver

Better: single thread for both interfaces
- Notification "bound" to TCB
- Signal delivered as "IPC" from EP

Separate thread per interface?

Concurrency control, complexity!

# IPC Message Format

Raw data

| Tag | Message | Caps (on Send) Badges (on Receive) | CSpace reference for receiving caps (Receive only) |
|-----|---------|-----------------------------------|---------------------------------------------------|

| Label | Caps unwrapped | # Caps | Msg Length |
|-------|----------------|--------|------------|

Semantics defined by IPC protocol (Kernel or user)

Bitmap indicating caps which had badges extracted

Caps sent or received

UNSW
SYDNEY

# Client-Server IPC Example

**Client**

```
seL4_MessageInfo_t tag = seL4_MessageInfo_new(0, 0, 0, 1);
seL4_SetTag(tag);
seL4_SetMR(0,1);
seL4_Call(server_c, tag);
```

Load into tag register

Set message register #0

**Server**

```
ut_t* reply_ut = ut_alloc(seL4_ReplyBits, &cspace);
seL4_CPtr reply = cspace_alloc_slot(&cspace);
err = cspace_untyped_retype(&cspace, reply_ut->cap, reply,
                            seL4_ReplyObject, seL4_ReplyBits);
seL4_CPtr badged_ep =  cspace_alloc_slot(&cspace);
cspace_mint(&cspace, badged_ep, &cspace, ep, seL4_AllRights, 0xff);
...
seL4_Word badge;
seL4_MessageInfo_t msg = seL4_Recv(ep, &badge, reply);
...
seL4_MessageInfo_t response = seL4_MessageInfo_new(0, 0, 0, 1);
seL4_NBSend(reply,response);
```

Allocate slot & retype to RO

Mint cap with badge 0xff

Wait on EP, receiving badge, setting RO

Reply to sender identified by RO

Should really use ReplyRecv!

UNSW
SYDNEY