**TECHNISCHE UNIVERSITÄT DRESDEN**

**Faculty of Computer Science** Institute of Systems Architecture, Operating Systems Group

# MICROKERNEL-BASED OPERATING SYSTEMS

# INTRODUCTION

**CARSTEN WEINHOLD**

- Provide deeper understanding of OS mechanisms

- Illustrate alternative OS design concepts

- Promote OS research at TU Dresden

- Make you all enthusiastic about OS development in general and microkernels in particular

- Lecture every week
  - Online, as downloadable videos
  - Uploaded at the beginning of the week

- Subscribe to mailing list:
  - https://os.inf.tu-dresden.de/mailman/listinfo/mos2020

- Exercises (roughly) bi-weekly Tuesday, 2:50 PM

- Practical exercises:

  - Computer lab or online (announced on mailing list)

- Paper reading exercises:

  - Read a paper beforehand

  - Sum it up and prepare 3 questions

  - We expect you to actively participate in discussion

- First exercise: in two weeks

  - To be announced on website and mailing list
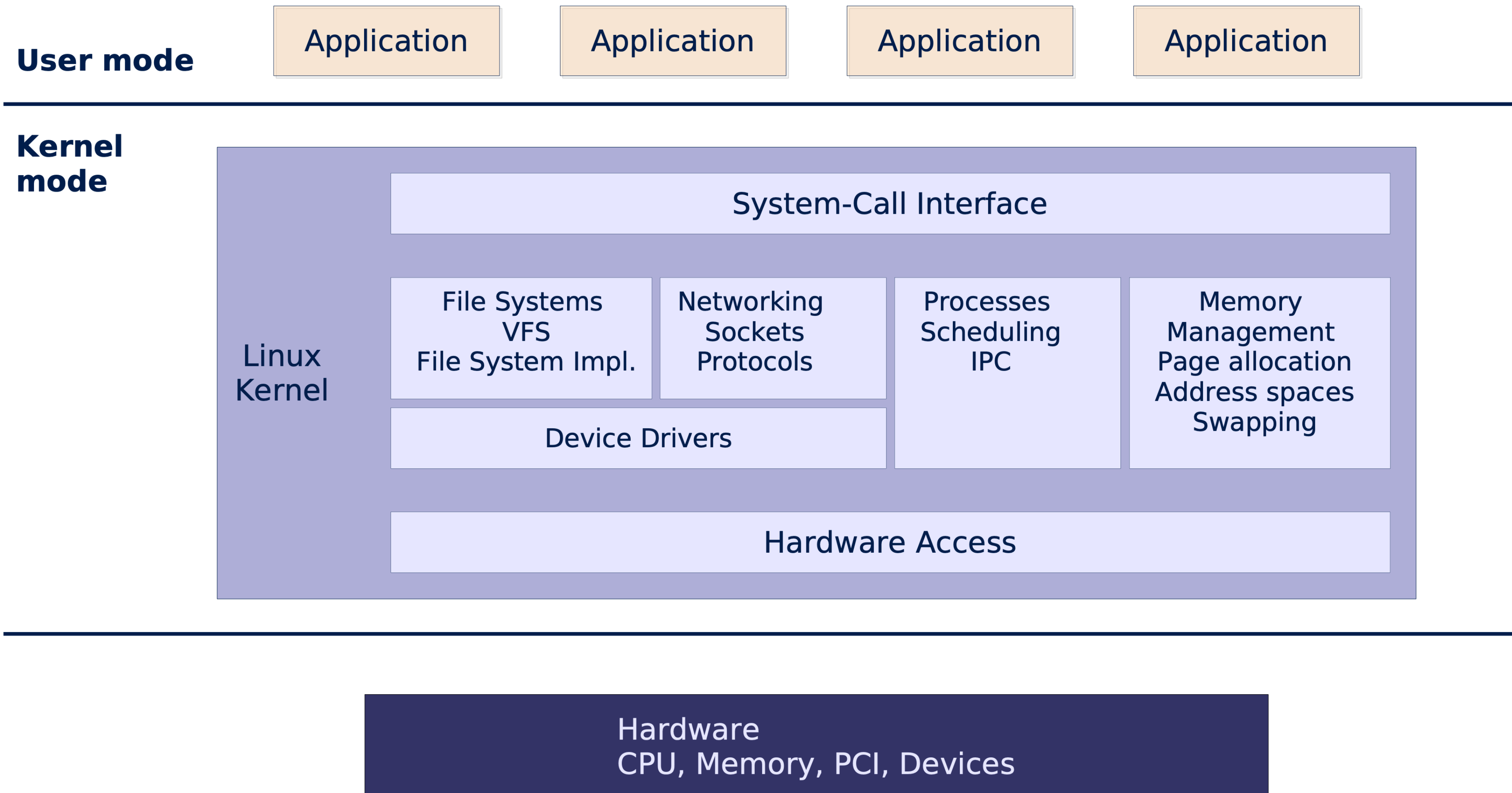
# More Hands-On: Complex Lab

- Complex lab Microkernel-based Operating Systems
  - Build several components of an OS

- This term, in parallel to lecture
  - Starts in week of November 23, 2020
  - Watch for announcement on complex lab website and on MOS lecture mailing list

# MICROKERNEL-BASED OPERATING SYSTEMS

## INTRODUCTION

- Manage the available resources
  - Hardware (CPU, memory, …)
  - Software (file systems, networking stack, …)

- Provide easier-to-use interface to access resources
  - Unix: read/write data from/to sockets instead of fiddling with TCP/IP packets on your own

- Perform privileged / HW-specific operations
  - x86: ring 0 vs. ring 3
  - Device drivers

- Provide separation and collaboration
  - Isolate users / processes from each other
  - Allow cooperation if needed (e.g., sending messages between processes)

# Monolithic Kernels: Linux

**User mode**

| Application | Application | Application | Application |

**Kernel mode**

Linux Kernel

System-Call Interface

| File Systems VFS File System Impl. | Networking Sockets Protocols | Processes Scheduling IPC | Memory Management Page allocation Address spaces Swapping |

Device Drivers

Hardware Access

Hardware
CPU, Memory, PCI, Devices

# Monolithic Kernels: Problems

- Security issues
  - All components run in privileged mode
  - Direct access to all kernel-level data
  - Module loading → easy living for rootkits

- Resilience issues
  - Faulty drivers can crash the whole system
  - 75% of today's OS kernels are drivers

- Software-level issues
  - Complexity is hard to manage
  - Custom OS for hardware with scarce resources?

- Minimal OS kernel
  - less error prone
  - small *Trusted Computing Base*
  - suitable for verification

- System services in user-level *servers*
  - flexible and extensible

- Protection between individual components
  - More resilient – crashing component does not (necessarily...) crash the whole system
  - More secure – inter-component protection

# The Microkernel Vision

**User mode**

| Application | Application | Application | Application |
|---|---|---|---|

| File Systems VFS File System Impl. | Networking Sockets Protocols | Memory Management Page allocation Swapping |
|---|---|---|

| Device Drivers | |

**Kernel mode**

| System-Call Interface | Address Spaces Threads Scheduling IPC |
|---|---|
| Hardware Access | |

Microkernel

Hardware
CPU, Memory, PCI, Devices

- OS personalities

- Customizability
  - Servers may be configured to suit the target system (small embedded systems, desktop PCs, SMP systems, …)
  - Remove unneeded servers

- Enforce reasonable system design
  - Well-defined interfaces between components
  - No access to components besides these interfaces
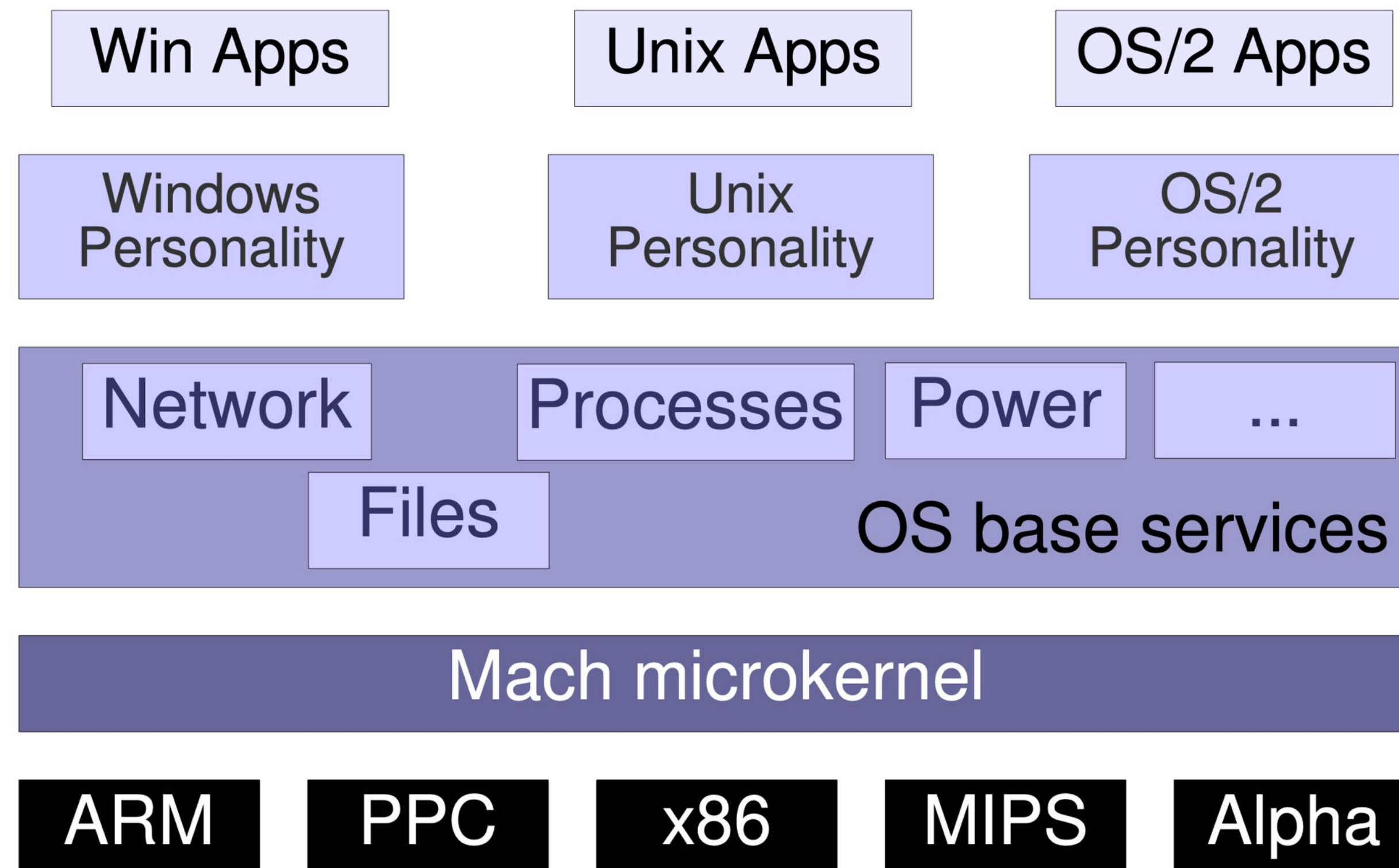  - Improved maintainability

- Mach – developed at CMU, 1985 - 1994
  - Rick Rashid (former head of MS Research)
  - Avie Tevanian (former Apple CTO)
  - Brian Bershad (professor @ U. of Washington)
  - ...

- Foundation for several real systems
  - Single Server Unix (BSD4.3 on Mach)
  - MkLinux (OSF)
  - IBM Workplace OS
  - NeXT OS → Mac OS X

- Simple, extensible *communication kernel*
  - "Everything is a pipe." – *ports* as secure communication channels
- Multiprocessor support
- Message passing by mapping
- Multi-server OS personality
- POSIX-compatibility
- Shortcomings
  - performance
  - drivers still in the kernel

**TECHNISCHE UNIVERSITÄT DRESDEN**

- Main goals:
  - multiple OS personalities
  - run on multiple HW architectures

| Win Apps | Unix Apps | OS/2 Apps |
|---|---|---|
| Windows Personality | Unix Personality | OS/2 Personality |

Network | Processes | Power | ...

Files

OS base services

Mach microkernel

| ARM | PPC | x86 | MIPS | Alpha |

- Never finished (but spent 1 billion $)
- Failure causes:
  - Underestimated difficulties in creating OS personalities
  - Management errors, forced divisions to adopt new system without having a system
  - "Second System Effect": too many fancy features
  - Too slow
- Conclusion: Microkernel worked, but system atop the microkernel did not

- OS personalities did not work
- Flexibility – but monolithic kernels became flexible, too (Linux kernel modules)
- Better design – but monolithic kernels also improved (restricted symbol access, layered architectures)
- Maintainability – still very complex
- Performance matters a lot

# Mircokernels: Proven Advantages

- Subsystem protection / isolation
- Code size
  - Microkernel-based OS
    - Fiasco kernel:                                        ~ 34,000 LoC
    - "HelloWorld" (+boot loader +root task): ~ 10,000 LoC
  - Linux kernel (3.0.4., x86 architecture):
    - Kernel:           ~ 2.5 million LoC
    - +drivers:        ~ 5.4 million LoC
  - (*generated using David A. Wheeler's 'SLOCCount'*)
- Customizability
  - Tailored memory management / scheduling / … algorithms
  - Adaptable to embedded / real-time / secure / … systems

- We need fast and efficient kernels
  - covered in the "Microkernel construction" lecture in the summer term
- We need fast and efficient OS services
  - Memory and resource management
  - Synchronization
  - Device Drivers
  - File systems
  - Communication interfaces
  - Subject of this lecture

- Minix @ FU Amsterdam (Andrew Tanenbaum)

- Singularity @ MS Research

- EROS/CoyotOS @ Johns Hopkins University

- The L4 Microkernel Family
  - Originally developed by Jochen Liedtke at IBM and GMD
  - $2^{nd}$ generation microkernel
  - Several kernel ABI versions

- Jochen Liedtke:
  *"A microkernel does no real work."*
  - Kernel only provides inevitable mechanisms.
  - Kernel does not enforce policies.

- But what is inevitable?
  - Abstractions
    - Threads
    - Address spaces (tasks)
  - Mechanisms
    - Communication
    - Resource mapping
    - (Scheduling)

# TAKING A CLOSER LOOK:
## L4/FIASCO.OC AS A CASE STUDY

- "Everything is an object"
  - – Task       Address spaces
  - – Thread     Activities, scheduling
  - – IPC Gate   Communication, resource mapping
  - – IRQ        Communication
  - – Factory    Create other objects, enforce resource quotas

- One system call: **invoke_object()**
  - – Parameters passed in UTCB
  - – Types of parameters depend on type of object

- Kernel-provided objects
  - Threads
  - Tasks
  - IRQs
  - ...

- Generic communication object: IPC gate
  - Send message from sender to receiver
  - Used to implement **new objects** in **user-level** applications

- Everything above kernel built using user-level objects that provide a service
  - Networking stack
  - File system
  - ...

- Kernel provides
  - Object creation/management
  - Object interaction: Inter-Process Communication (IPC)

- To call an object, we need an address:
  - Telephone number
  - Postal address
  - IP address
  - ...

Client      Service 1

call(service1.ID)

Kernel

- Simple idea, right?
- ID is wrong? Kernel returns ENOTEXIST
- But not so fast! This scheme is insecure:
  - Client could simply "guess" IDs brute-force.
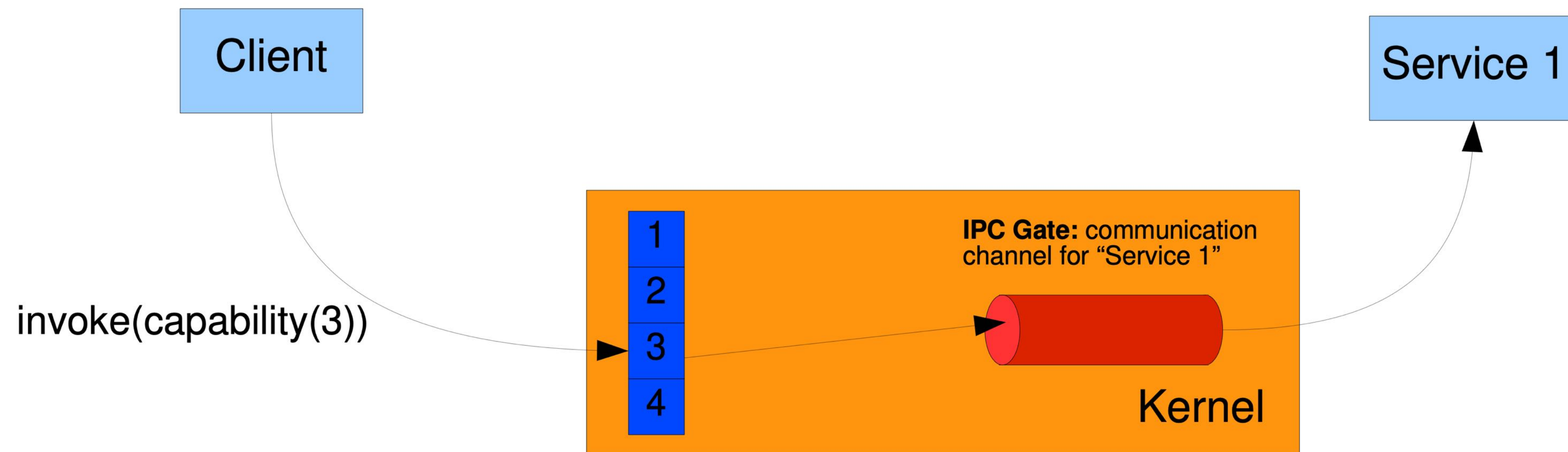  - Existence/non-existence can be used as a covert channel

- Global object IDs are
  - insecure (forgery, covert channels).
  - inconvenient (programmer needs to know about partitioning in advance)

- Solution in Fiasco.OC
  - Task-local *capability space* as an indirection
  - *Object capability* <u>required</u> to invoke object

- Per-task name space
  - Maps names to object capabilities.
  - Configured by task's creator

- Capability:
  - Reference to an object
  - Protected by the Fiasco.OC kernel
    - Kernel knows all capability-object mappings.
    - Managed as a per-process capability table.
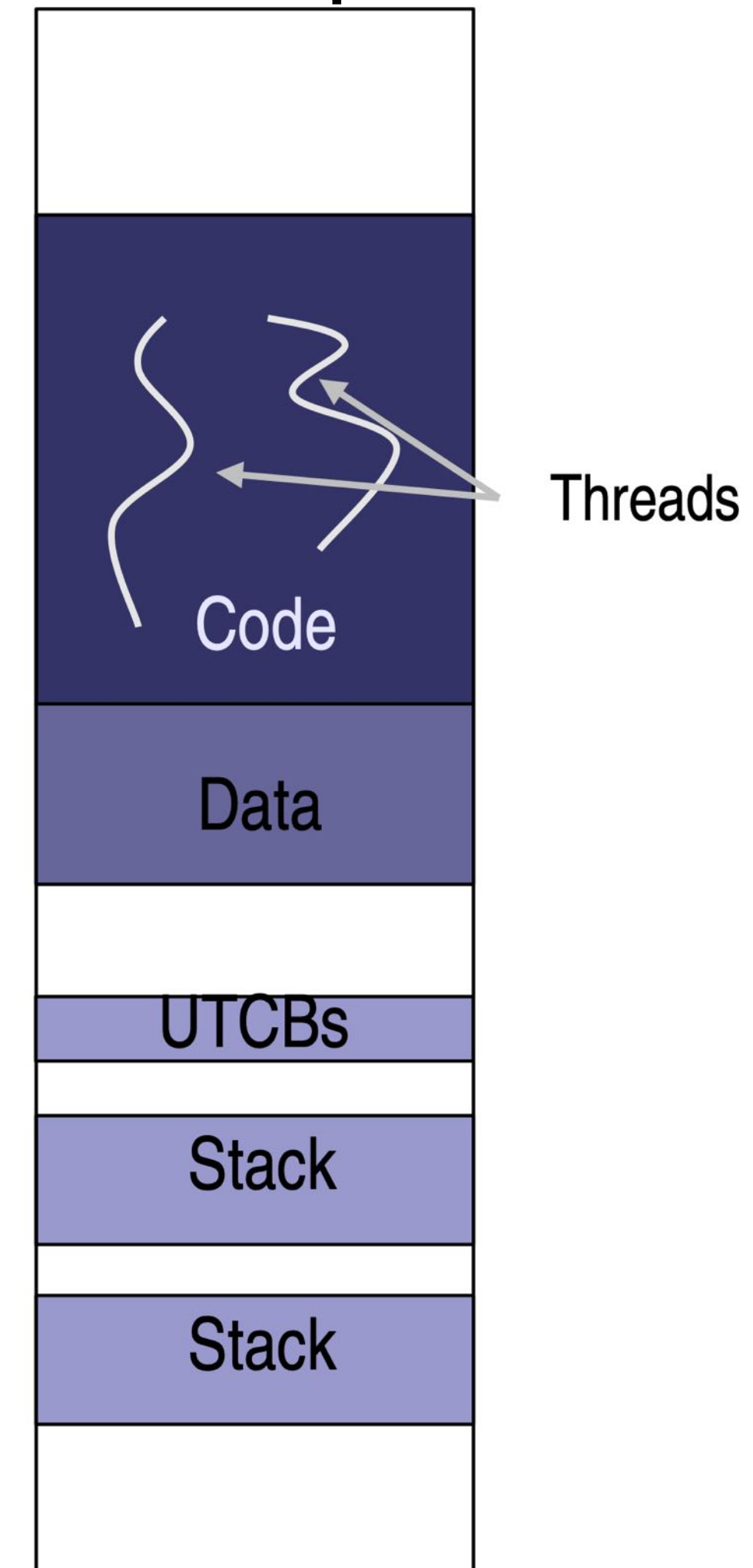    - User processes only use indexes into this table.

Client

Service 1

invoke(capability(3))

1
2
3
4

IPC Gate: communication
channel for "Service 1"

Kernel

- Kernel object for communication: *IPC gate*

- Inter-process communication (IPC)
    - Between threads
    - Synchronous

- Communication using IPC gate:
    - Sender thread puts message into its UTCB
    - Sender invokes IPC gate, blocks sender until receiver ready (i.e., waits for message)
    - Kernel copies message to receiver thread's UTCB
    - Both continue, knowing that message has been transferred/received

# MORE L4 CONCEPTS

**Address Space**

- Thread
  - Unit of Execution
  - Implemented as kernel object

- Properties managed by the kernel:
  - Instruction Pointer (EIP)
  - Stack (ESP)
  - Registers
  - User-level TCB

- User-level applications need to
  - allocate stack memory
  - provide memory for application binary
  - find entry point
  - ...

Threads

Code

Data

UTCBs

Stack

Stack

- Kernel object: IRQ
- Used for hardware and software interrupts
- Provides asynchronous signaling
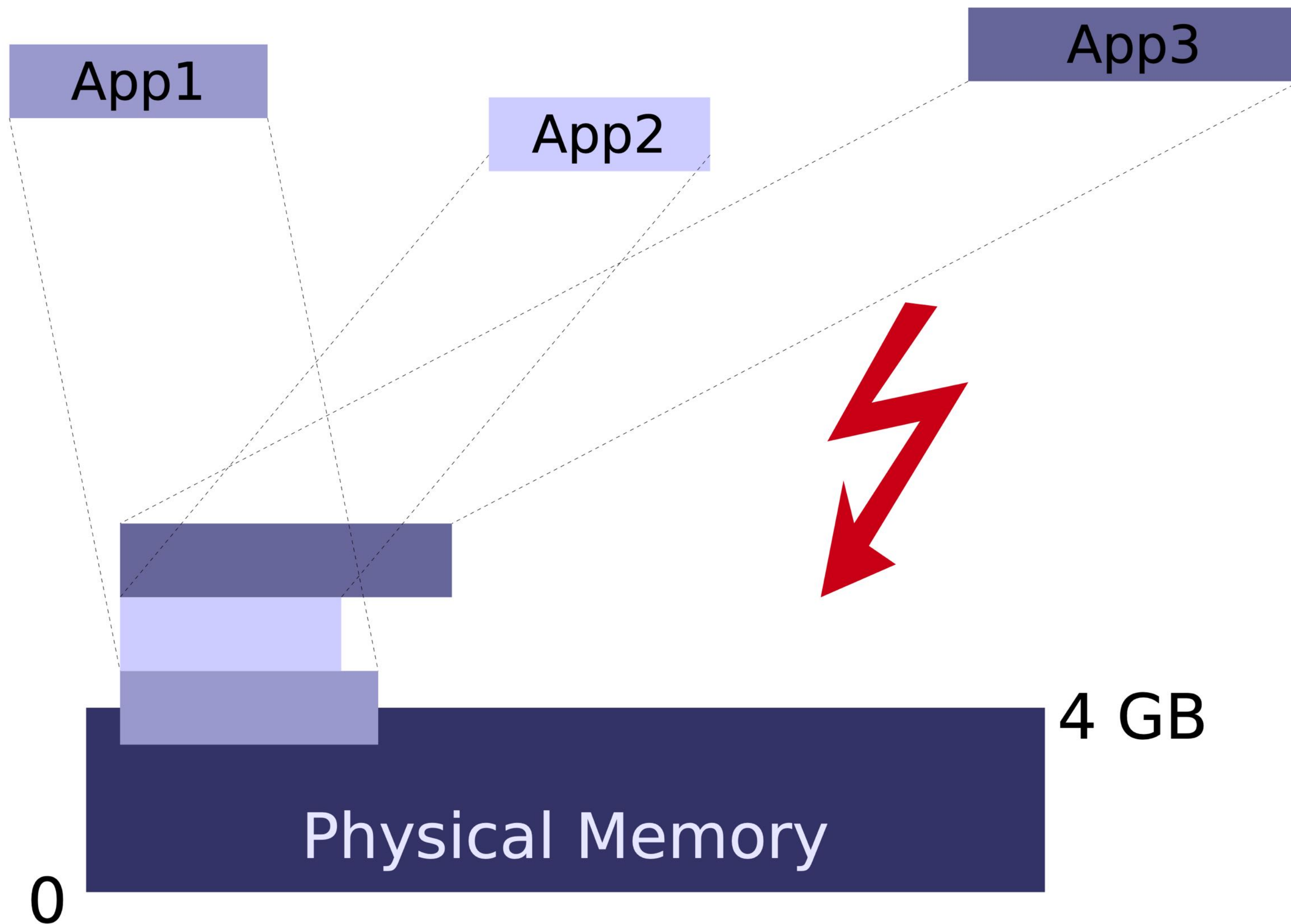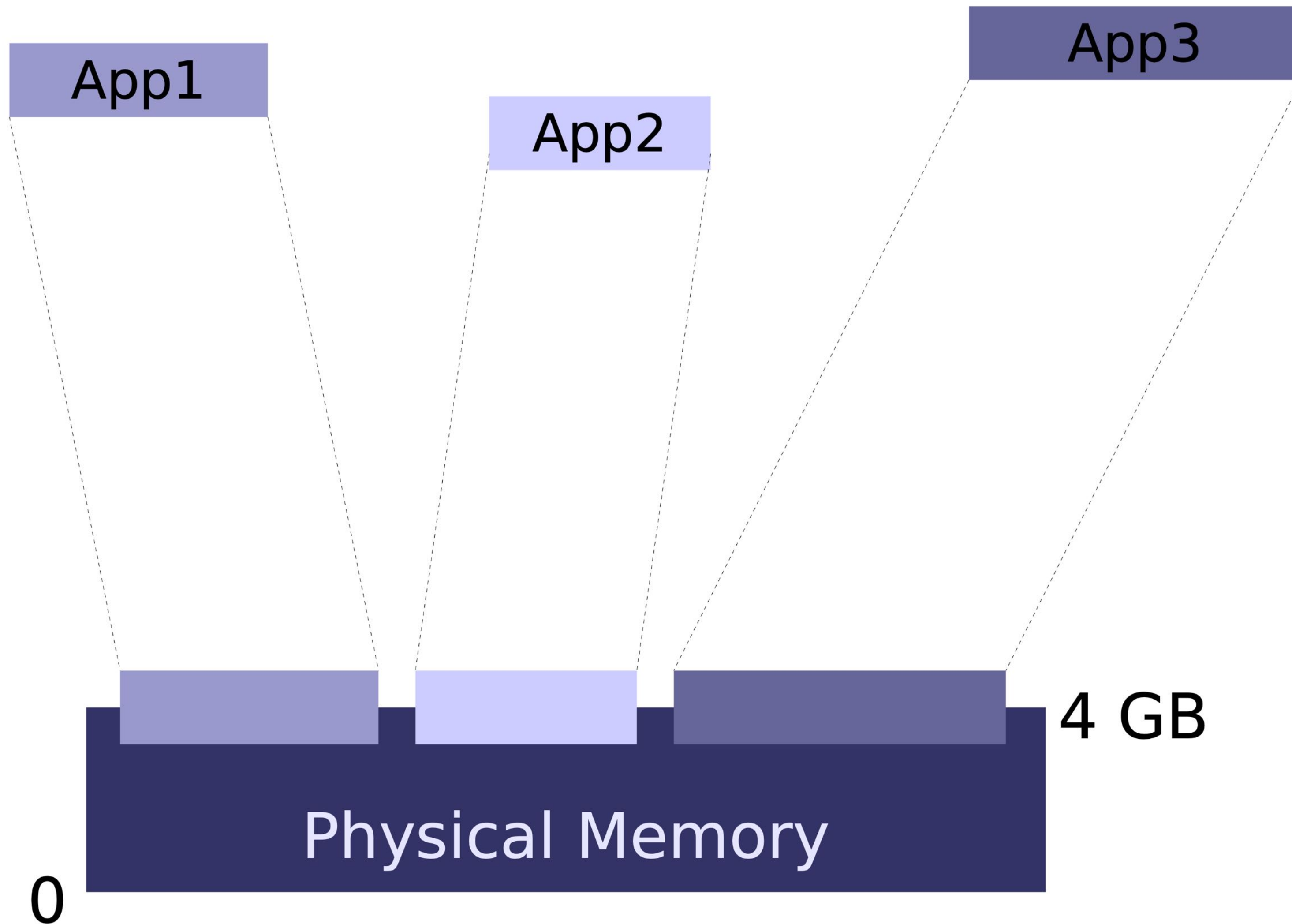  - invoke_object(irq_cap, WAIT)
  - invoke_object(irq_cap, TRIGGER)

IRQ → Kernel → invoke_object (irq_cap, ...) → User-space device driver

Image source:
https://commons.wikimedia.org/File:Ethernet_pci_card.jpg

34

RAM

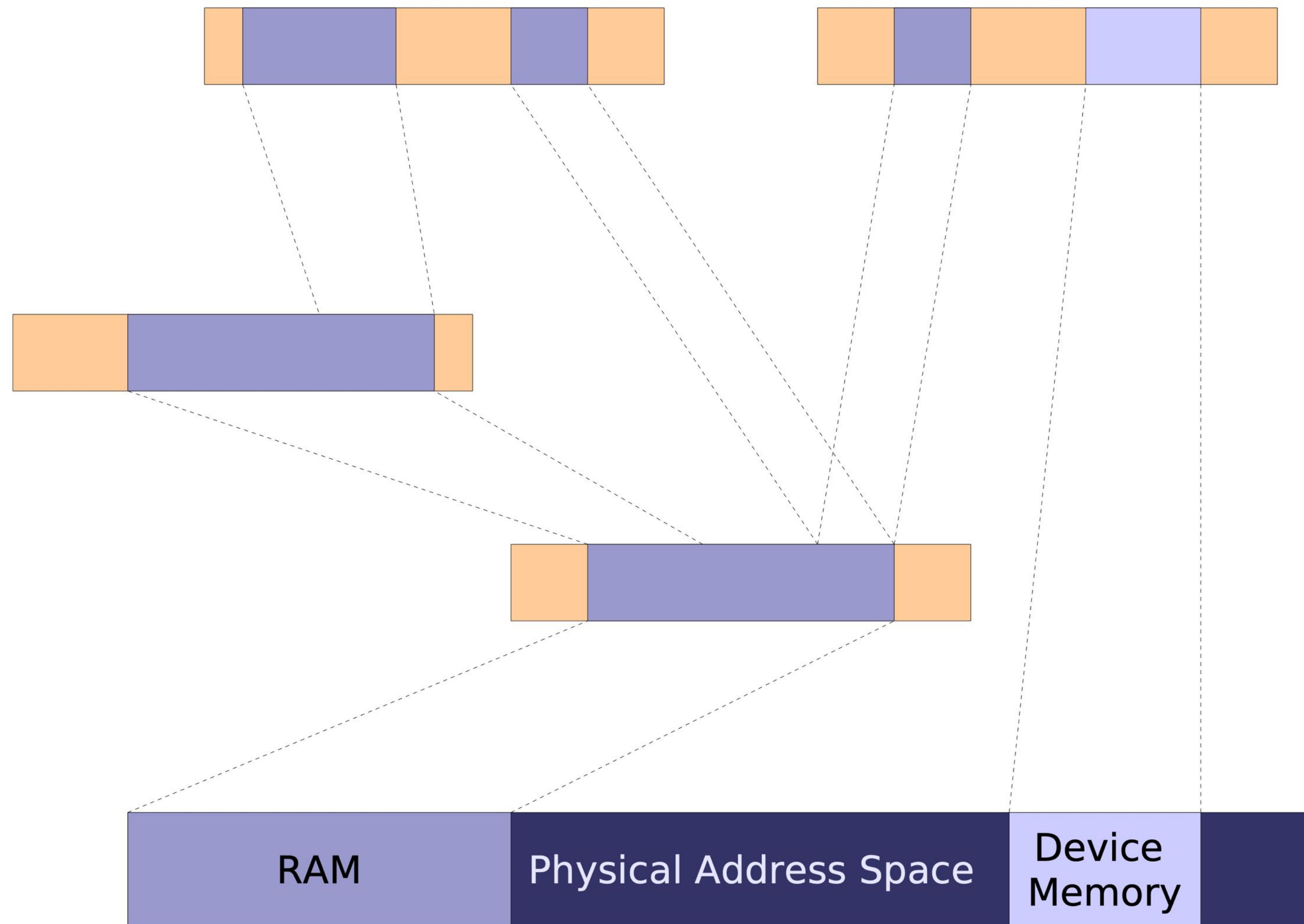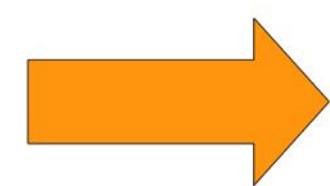Physical Address Space

Device Memory

- If a thread has access to a capability, it can map this capability to another thread

- Mapping / not mapping of capabilities used for implementing access control

- Abstraction for mapping: *flexpage*

- Flexpages describe mapping
  - location and size of resource
  - receiver's rights (read-only, mappable)
  - type (memory, I/O, communication capability)

- Summary of object types
  - Task
  - Thread
  - IPC Gate
  - IRQ
  - Factory

- Each task gets initial set of capabilities for some of these objects at startup

# WHAT CAN BE BUILT ON L4?

TECHNISCHE
UNIVERSITÄT
DRESDEN

- Fiasco.OC is <u>not</u> a full operating system!
  - No device drivers (except UART + timer)
  - No file system / network stack / …

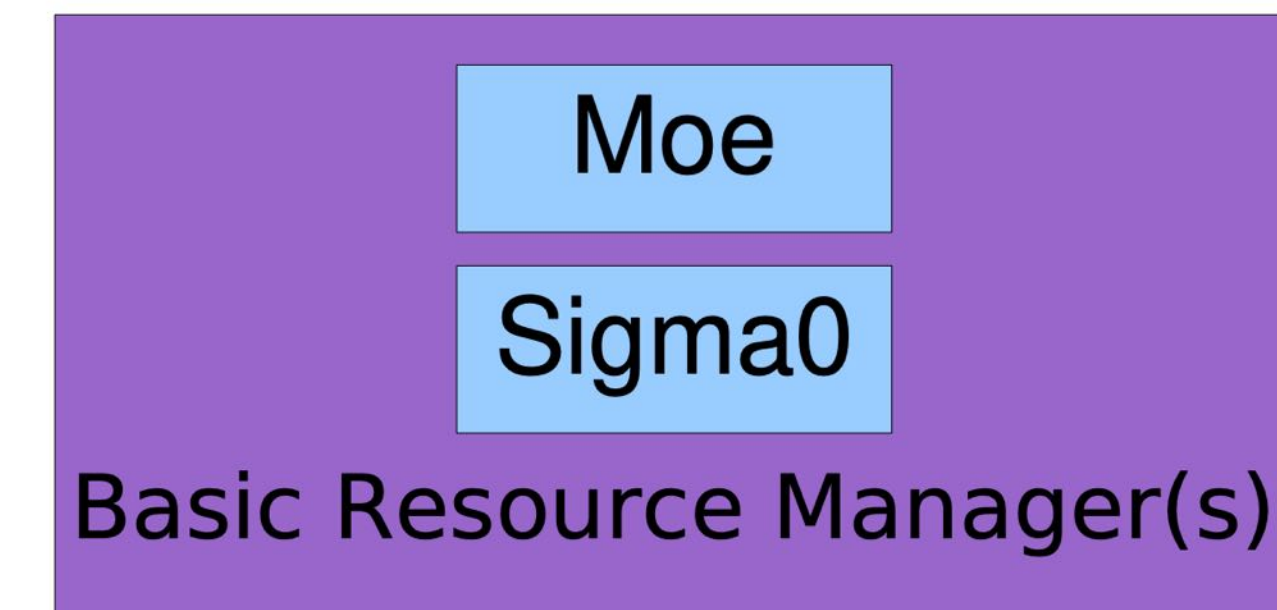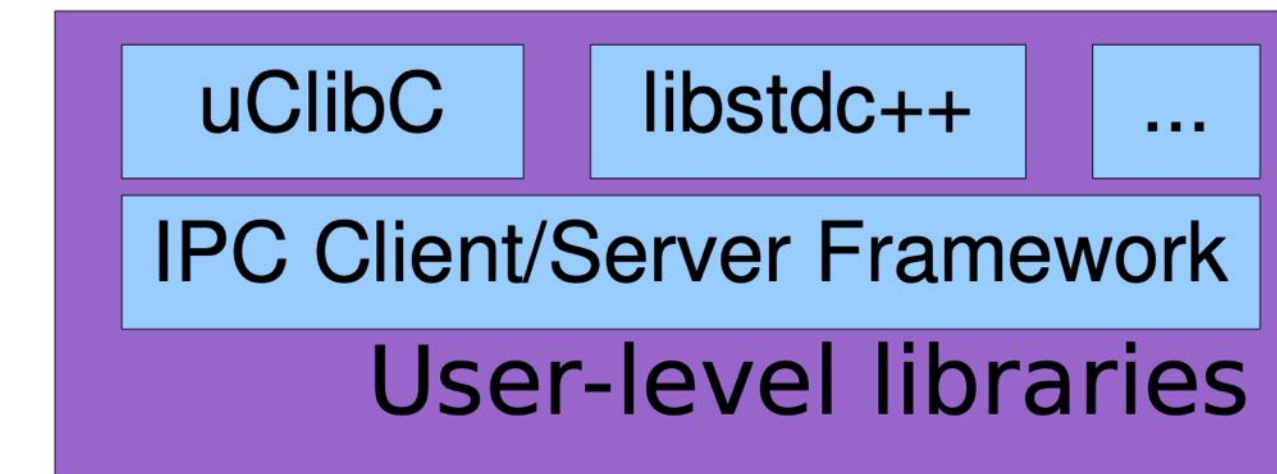- A microkernel-based OS needs to add these services as user-level components
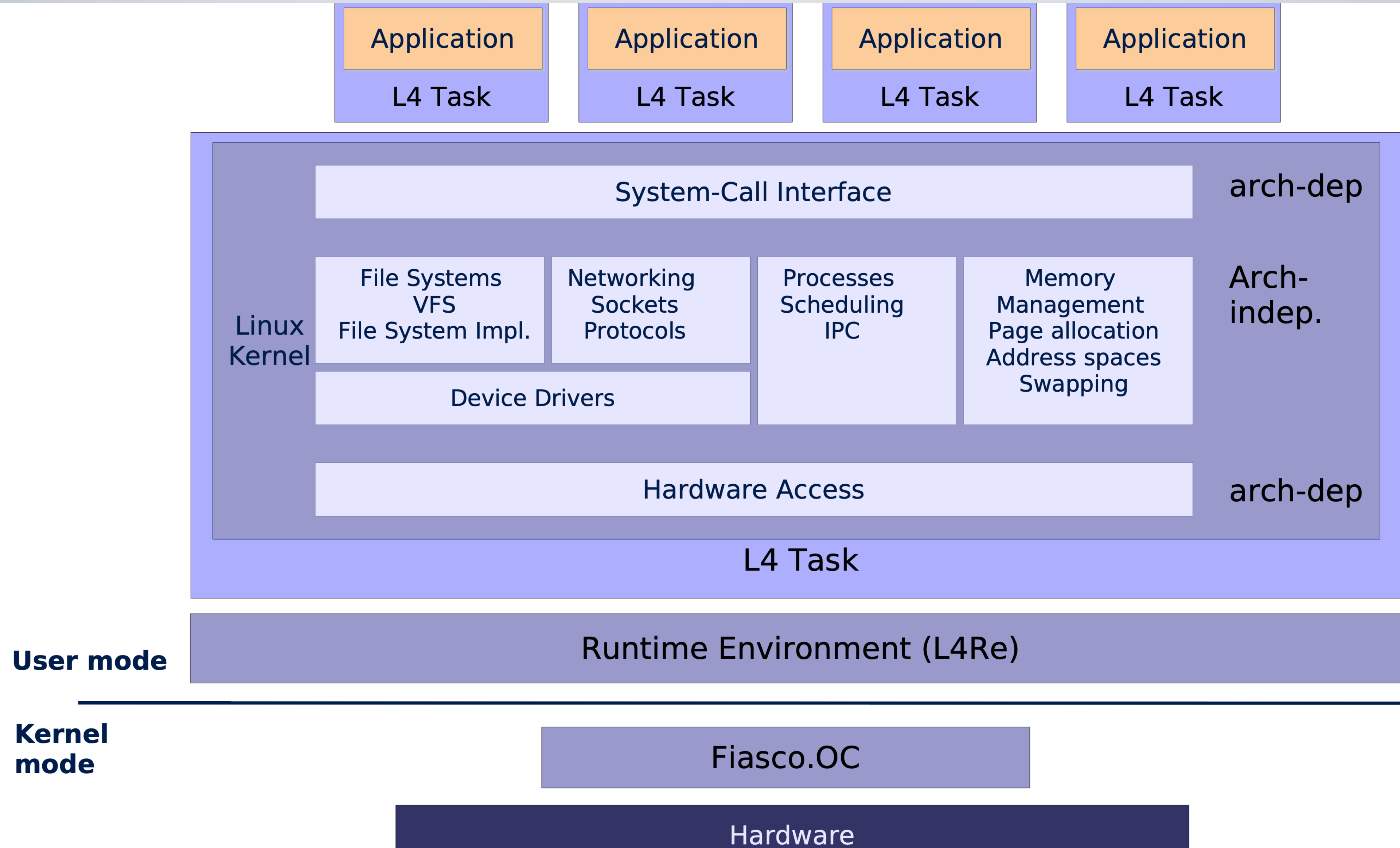
➡ L4Re – L4 Runtime Environment

**L4Re**

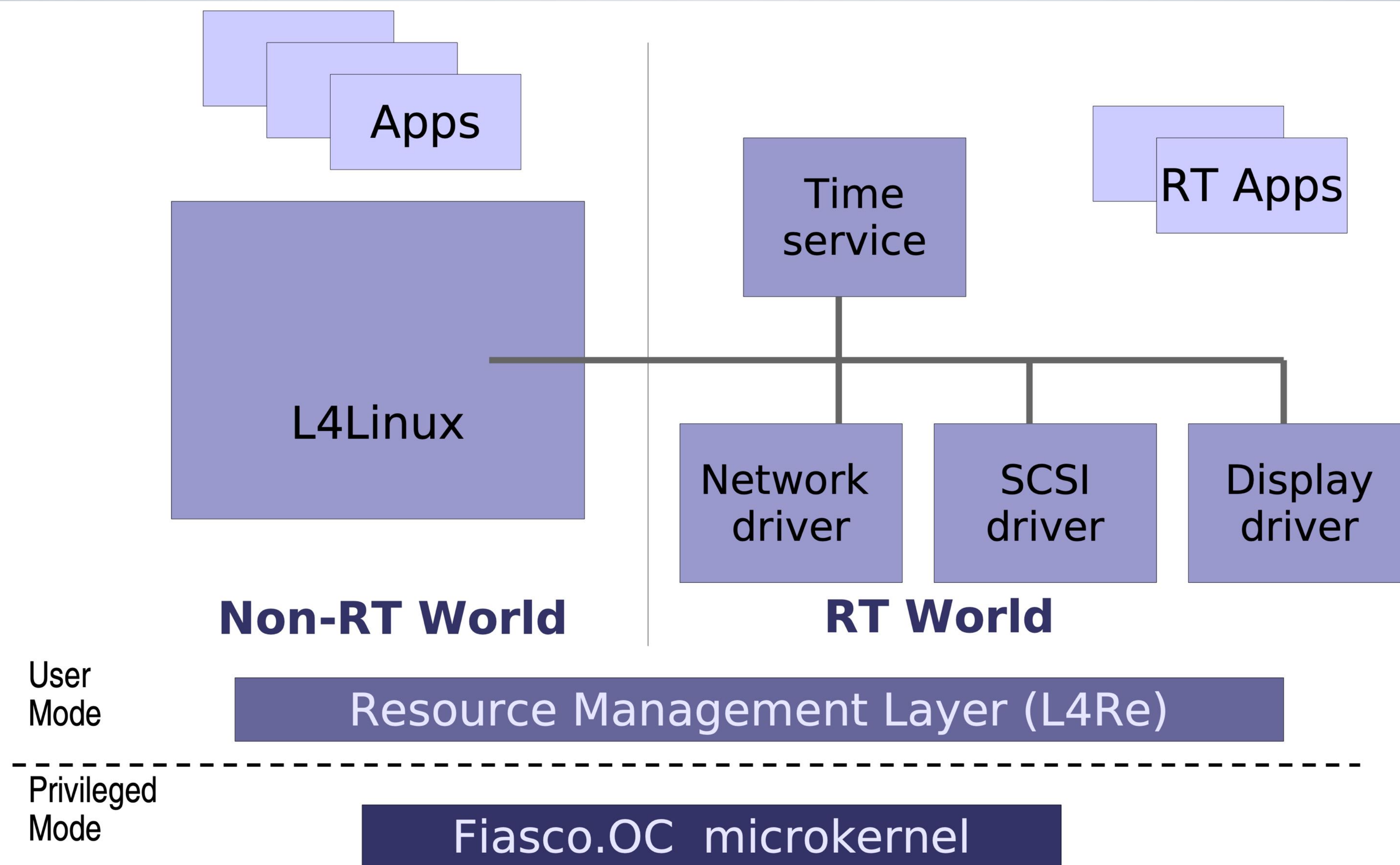| uClibC | libstdc++ | … |

IPC Client/Server Framework

User-level libraries

Ned

Init-style task loader

Moe

Sigma0

Basic Resource Manager(s)

**User mode**

**Kernel mode**

Fiasco.OC

TECHNISCHE
UNIVERSITÄT
DRESDEN

| Application | Application | Application | Application |
|:---:|:---:|:---:|:---:|
| L4 Task | L4 Task | L4 Task | L4 Task |

**Linux Kernel**

System-Call Interface — arch-dep

| File Systems VFS File System Impl. | Networking Sockets Protocols | Processes Scheduling IPC | Memory Management Page allocation Address spaces Swapping |
|:---:|:---:|:---:|:---:|

Device Drivers — Arch-indep.

Hardware Access — arch-dep
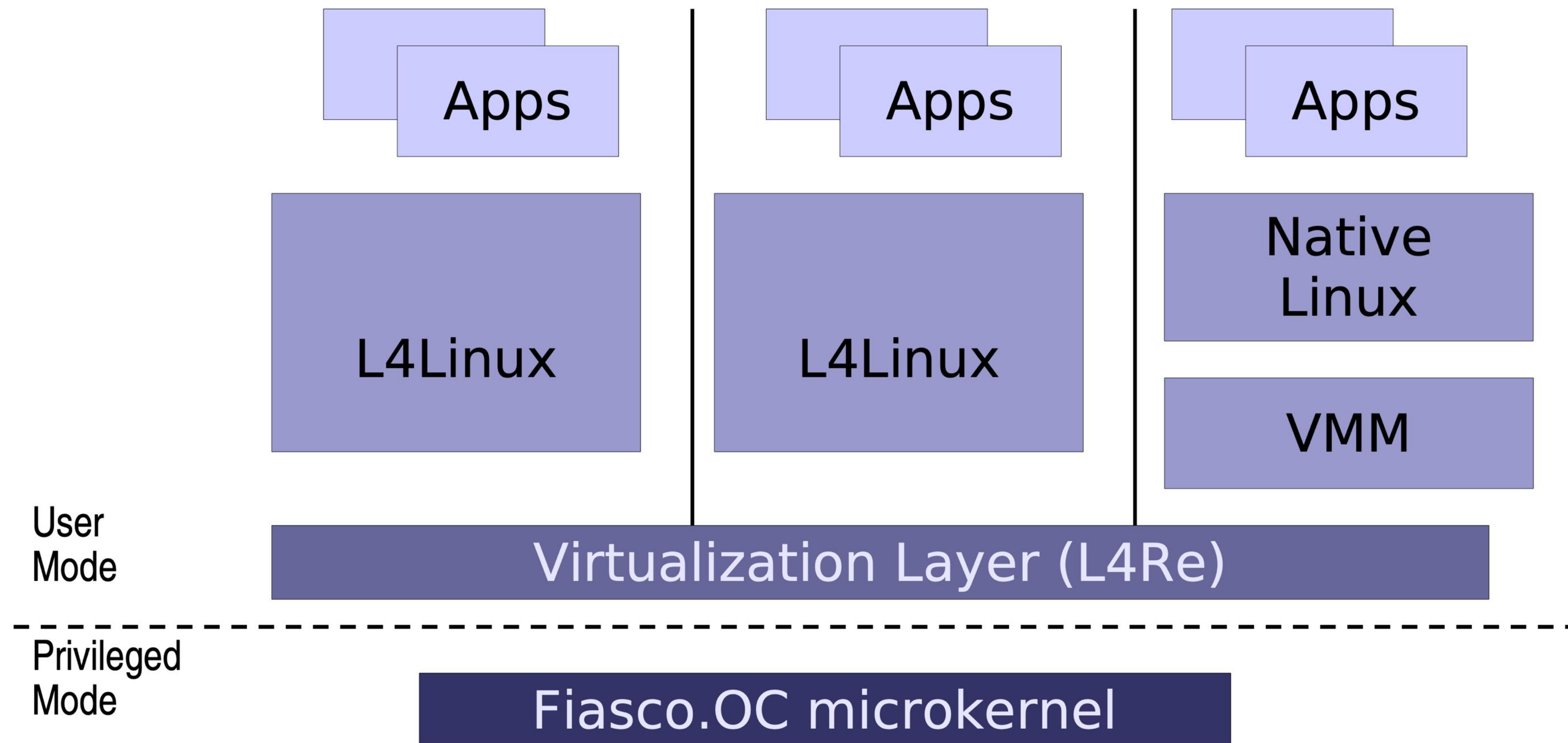
L4 Task

**User mode**

Runtime Environment (L4Re)

**Kernel mode**

Fiasco.OC

Hardware

42

- Isolate not only processes, but also complete Operating Systems (compartments)
- "Server consolidation"



| Apps | Apps | Apps |
|------|------|------|
| L4Linux | L4Linux | Native Linux |
| | | VMM |

User Mode — Virtualization Layer (L4Re)

Privileged Mode — Fiasco.OC microkernel

44

- Genode := C++-based OS framework developed here in Dresden
- Aim: hierarchical system in order to
  - Support resource partitioning
  - Layer security policies on top of each other

- **Basic mechanisms and concepts**
  - Memory management
  - Tasks, Threads, Synchronization
  - Communication
- **Building real systems**
  - What are resources and how to manage them?
  - How to build a secure system?
  - How to build a real-time system?
  - How to reuse existing code (Linux, standard system libraries, device drivers)?
  - How to improve robustness and safety?