



μ -Kernel Construction (1)

Overview, Motivation, Problems



Staff

- Lecturer: Raphael Neider
- PhD student

- Meeting Times
 - Tue, 14:00-15:30h
 - Bldg. 50.34, Room 161



Background

- L4Ka project (<http://l4ka.org>)
- State-of-the-art in microkernel construction
 - L4Ka::Pistachio
 - Microkernel
 - IDL4
 - Interface compiler
 - Virtual machines
 - Target application

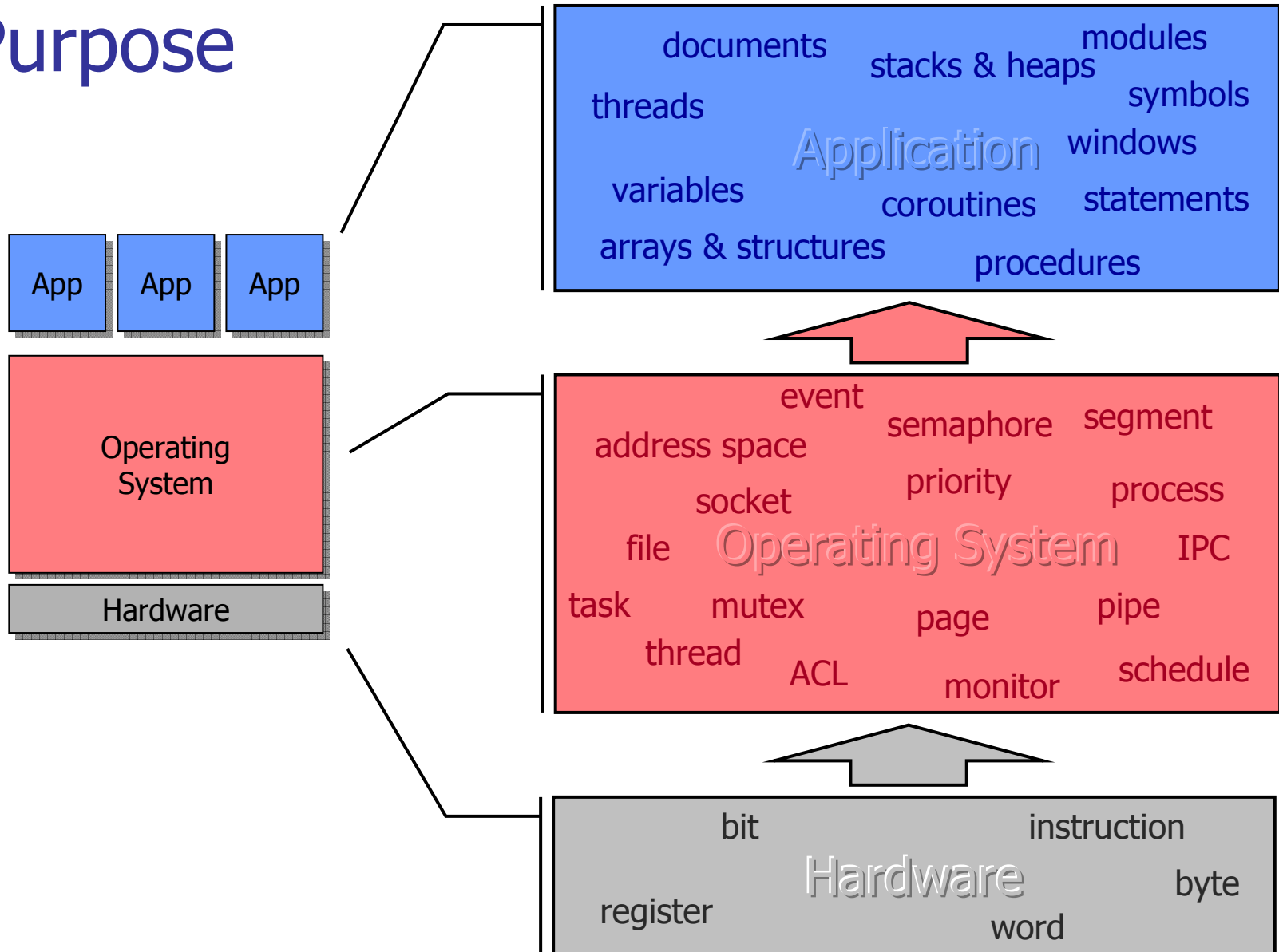


Purpose of Operating Systems

- Abstract from the hardware
 - Interrupts, exceptions
- Provide common services
 - Address spaces and protection
 - Threads and concurrency control
 - Persistency of data
- Bridge semantic gap
 - Application demands vs. hardware provides



Purpose





Operating System Designs

μ-kernel with object interfaces

- ✓ Standard interface
- ✓ User-defined interfaces
- ✓ Runs subsystems from different vendors
- ✓ Good flexibility
- ✓ Good minimalism
- ✓ Good performance
- ✗ Difficult to use
- ✗ Different paradigm

application-specific

- ✓ Ultimate flexibility
- ✓ Ultimate minimalism
- ✓ Ultimate performance
- ✓ Normal programming paradigm
- ✗ Proprietary and incompatible solutions

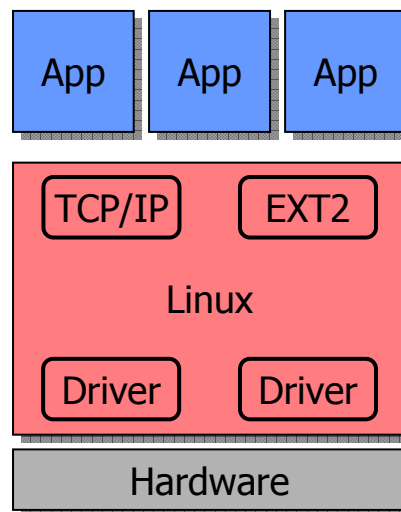
monolithic

- ✓ Standard interface
- ✓ Runs programs from different vendors
- ✗ Compromised interface
- ✗ Poor performance
- ✗ Inflexible



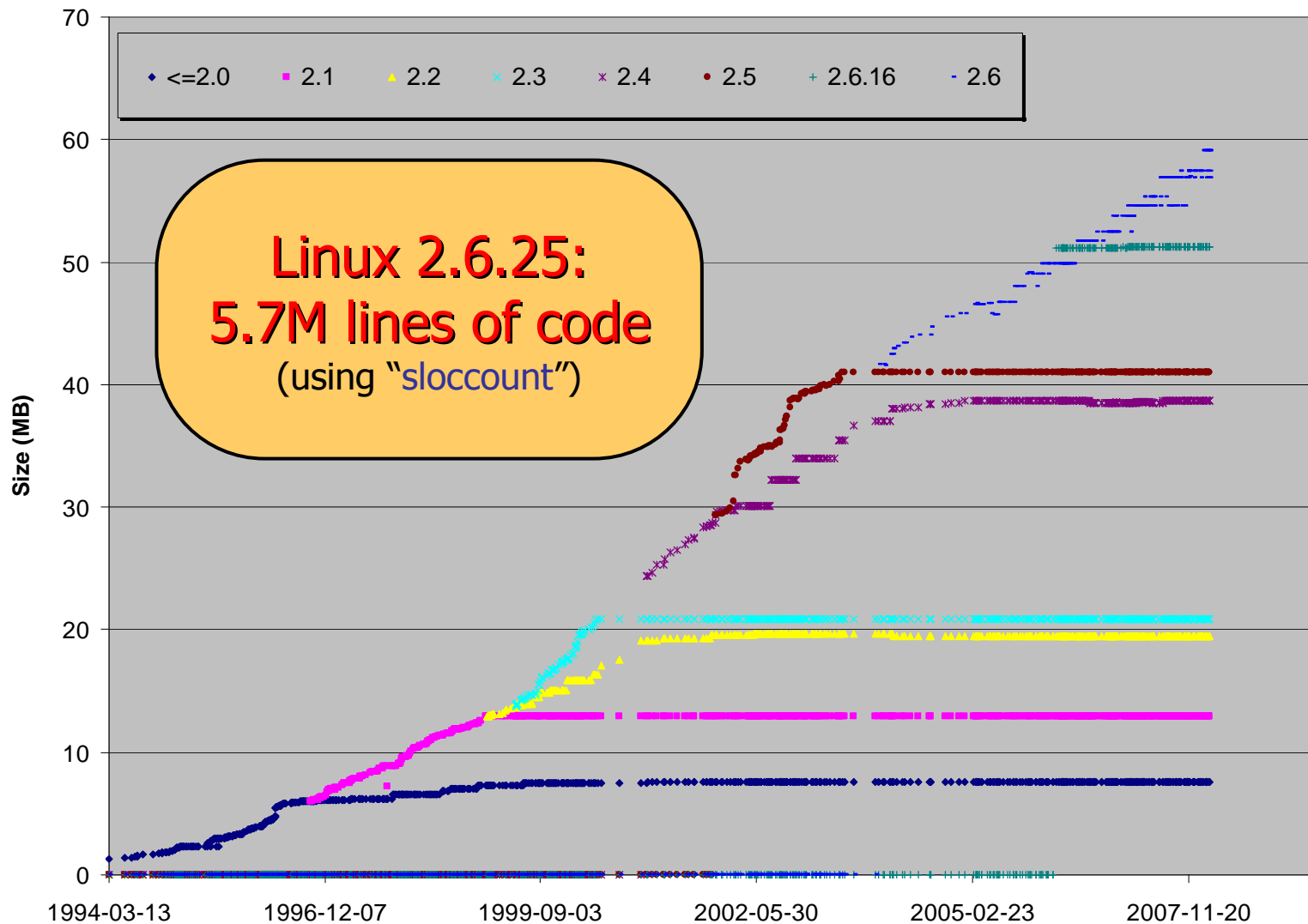
Monolithic Kernels – Advantages

- Kernel has access to everything, potentially
 - All optimizations are possible
 - All techniques/mechanisms/concepts can be implemented
- Kernel extended by adding more code





Linux Kernel Evolution (.tar.gz)





Approaches to Tackling Complexity

- Monolithic approaches
 - Layered Kernels
 - Modular Kernels
 - Object Oriented Kernels
- Alternatives
 - Extensible Kernels
 - **Microkernels**



History

- Monolithic kernels

- 1st-generation μ -kernels

- Mach *CMU, OSF* **External Pager**

- Chorus *Inria, Chorus*

- Amoeba *Vrije Universiteit*

- (L3) *GMD* **User-Level Driver**

- 2nd-generation μ -kernels

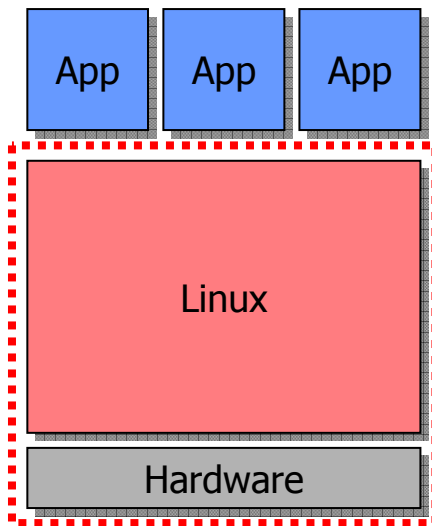
- (Spin) *U Washington*

- Exokernel *MIT*

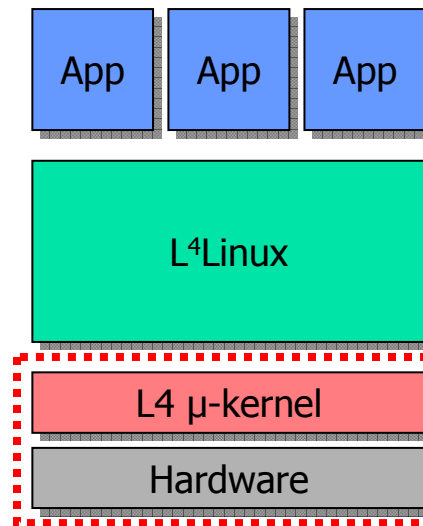
- L4 *GMD / IBM / UKa* **Recursive Address Spaces**



μ -Kernel Based Systems



Monolithic System

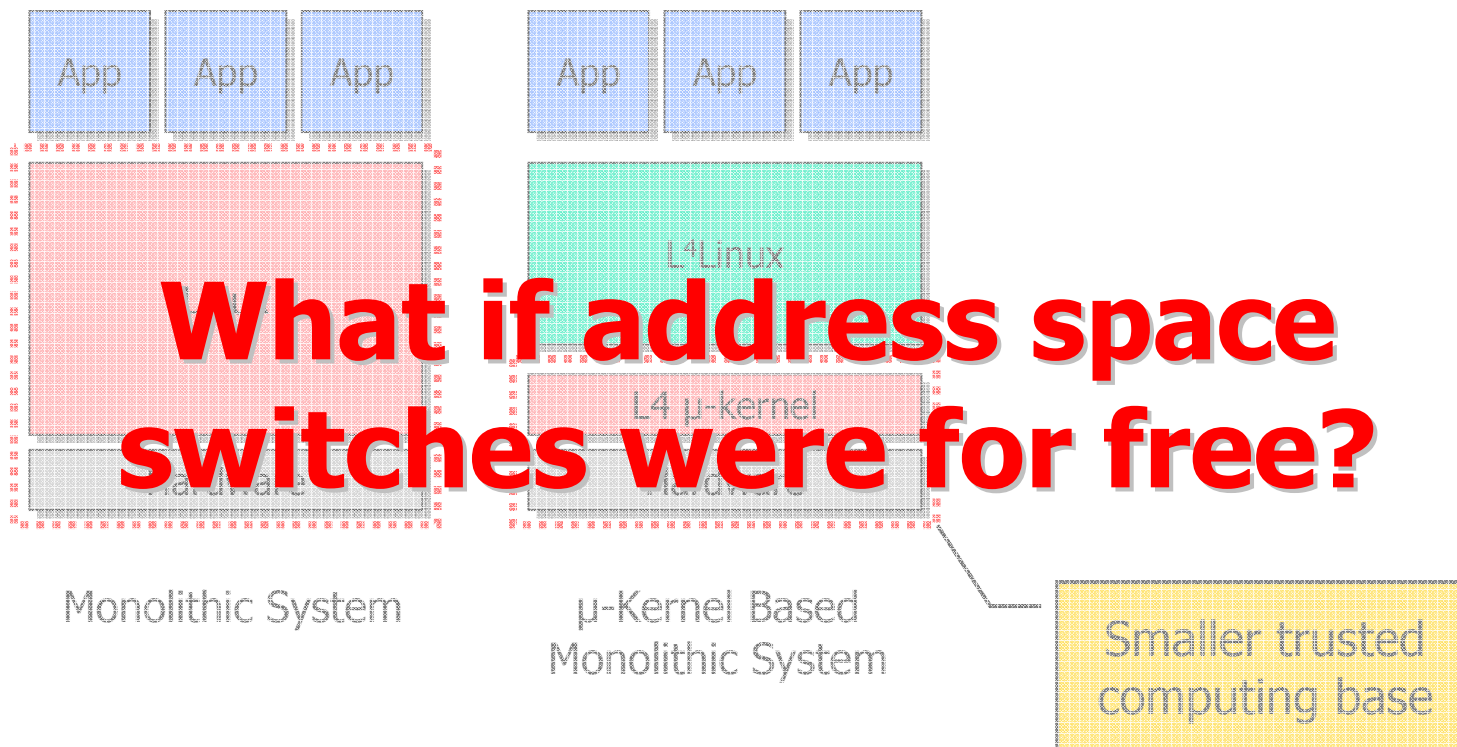


μ -Kernel Based Monolithic System

Smaller trusted computing base

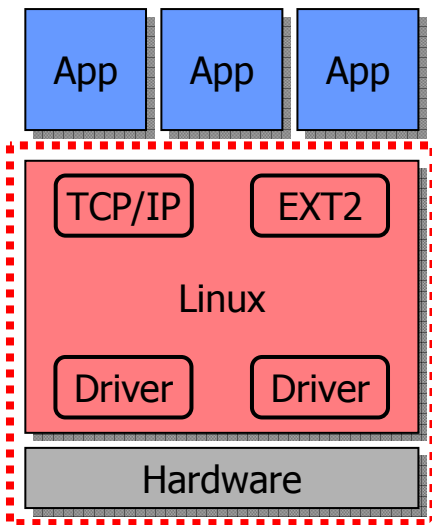


μ -Kernel Based Systems

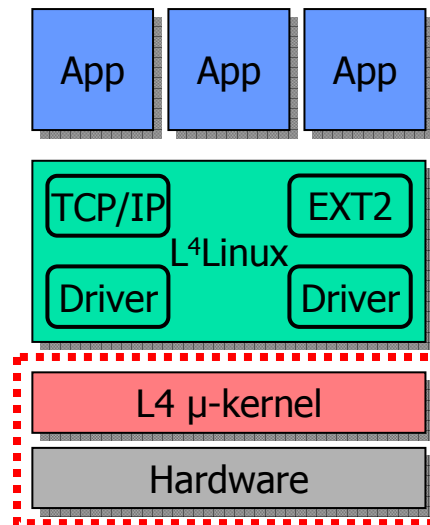




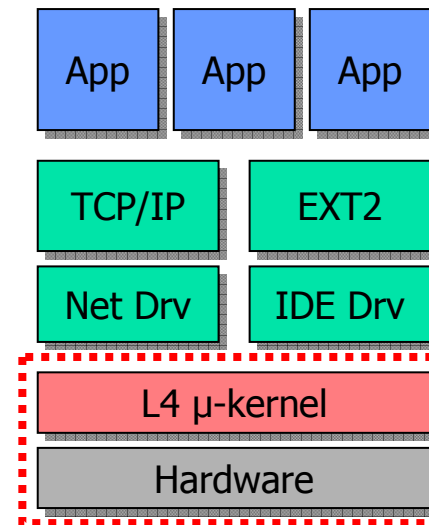
μ -Kernel Based Systems



Monolithic System

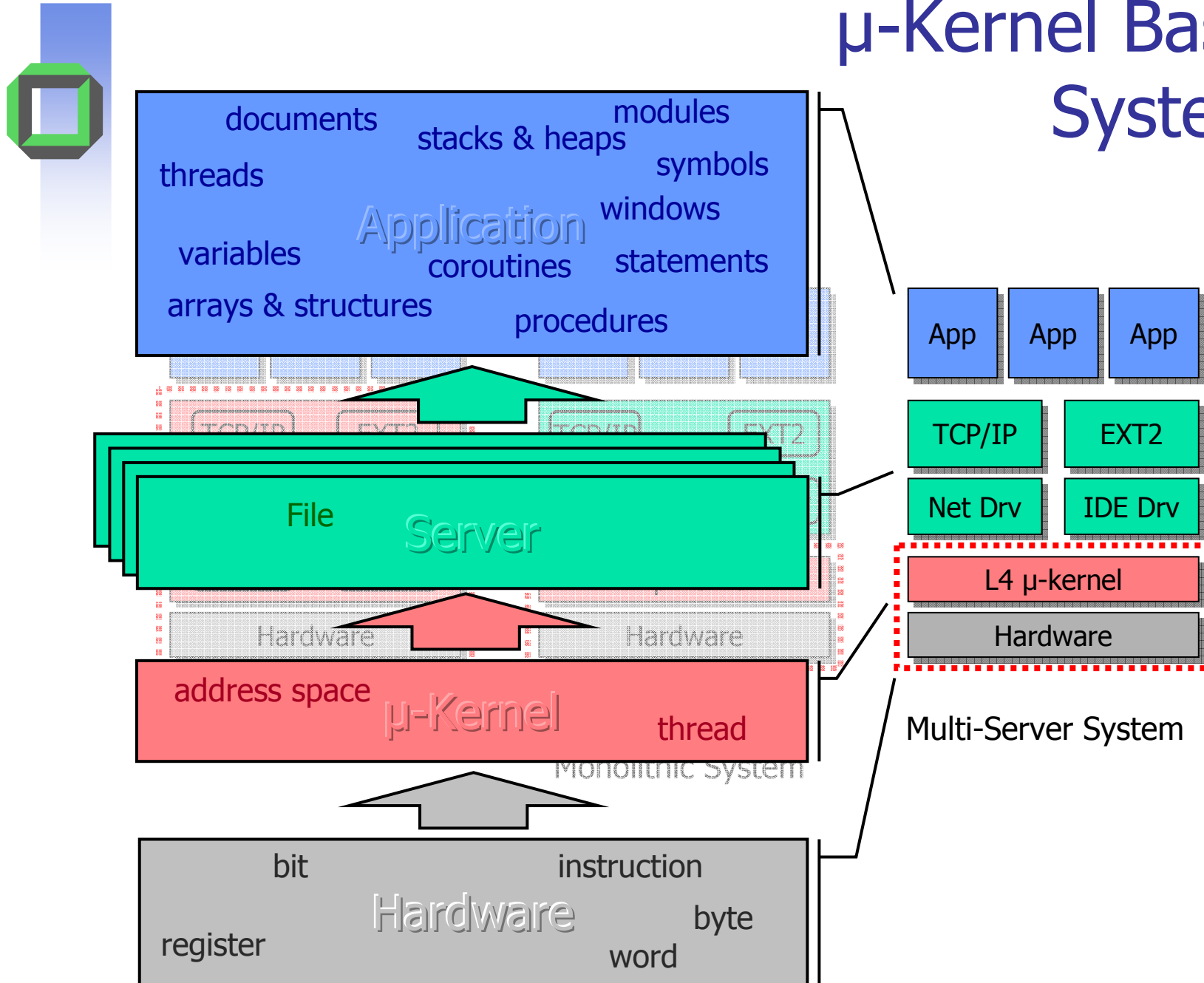


μ -Kernel Based Monolithic System



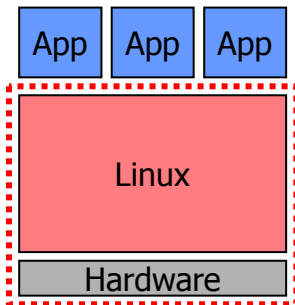
Multi-Server System

μ-Kernel Based Systems

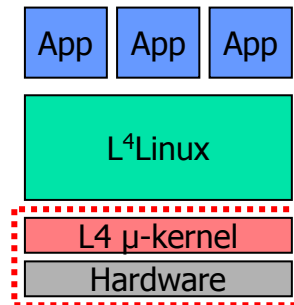




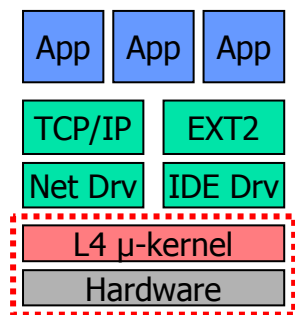
Microkernel Based Systems



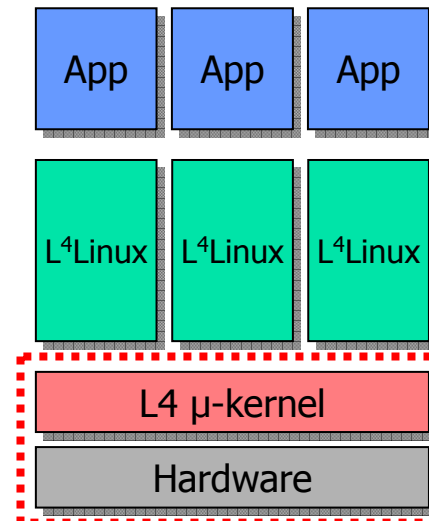
Monolithic System



μ-Kernel Based Monolithic System



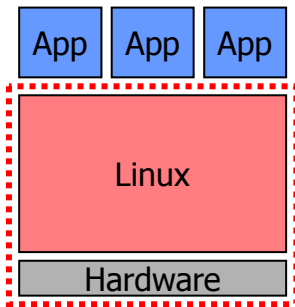
Multi-Server System



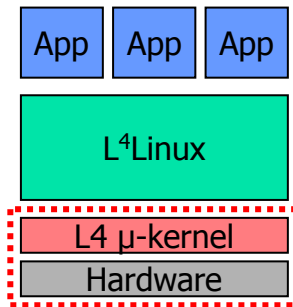
μ-Kernel Based Server Consolidation



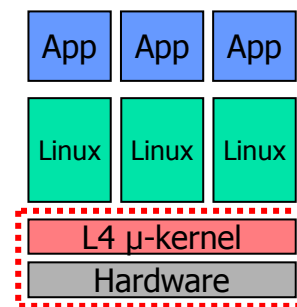
Microkernel Based Systems



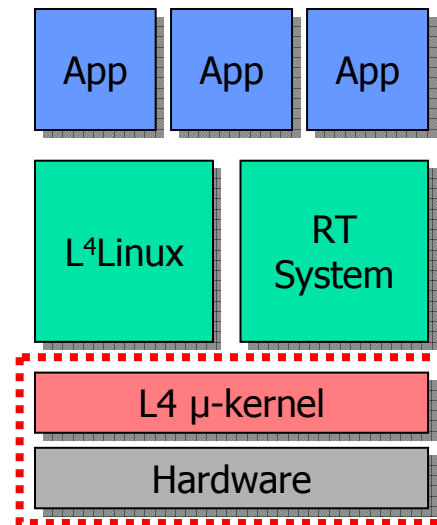
Monolithic System



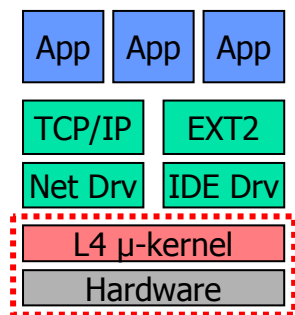
μ-Kernel Based Monolithic System



μ-Kernel Based Server Consolidation



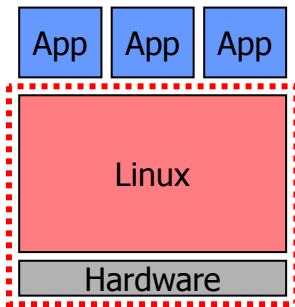
Coupling with Real-Time Systems



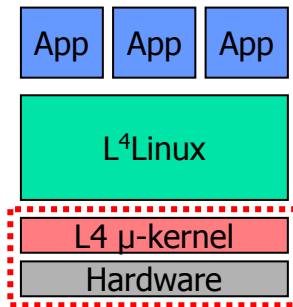
Multi-Server System



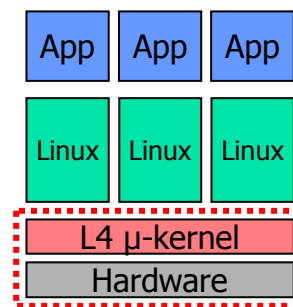
Microkernel Based Systems



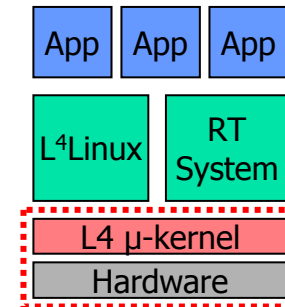
Monolithic System



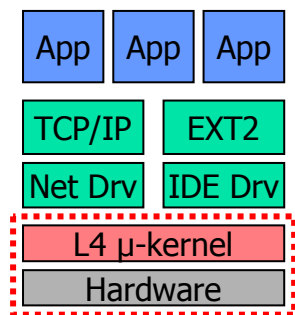
μ-Kernel Based Monolithic System



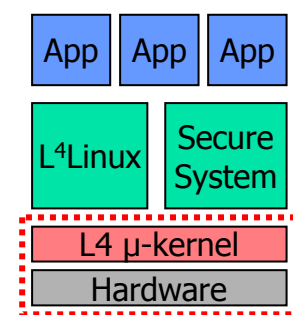
μ-Kernel Based Server Consolidation



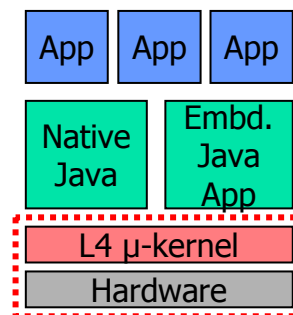
Coupling with Real-Time Systems



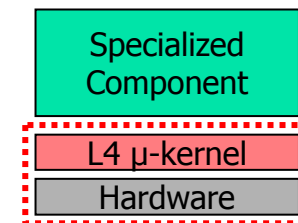
Multi-Server System



Coupling with Secure Systems



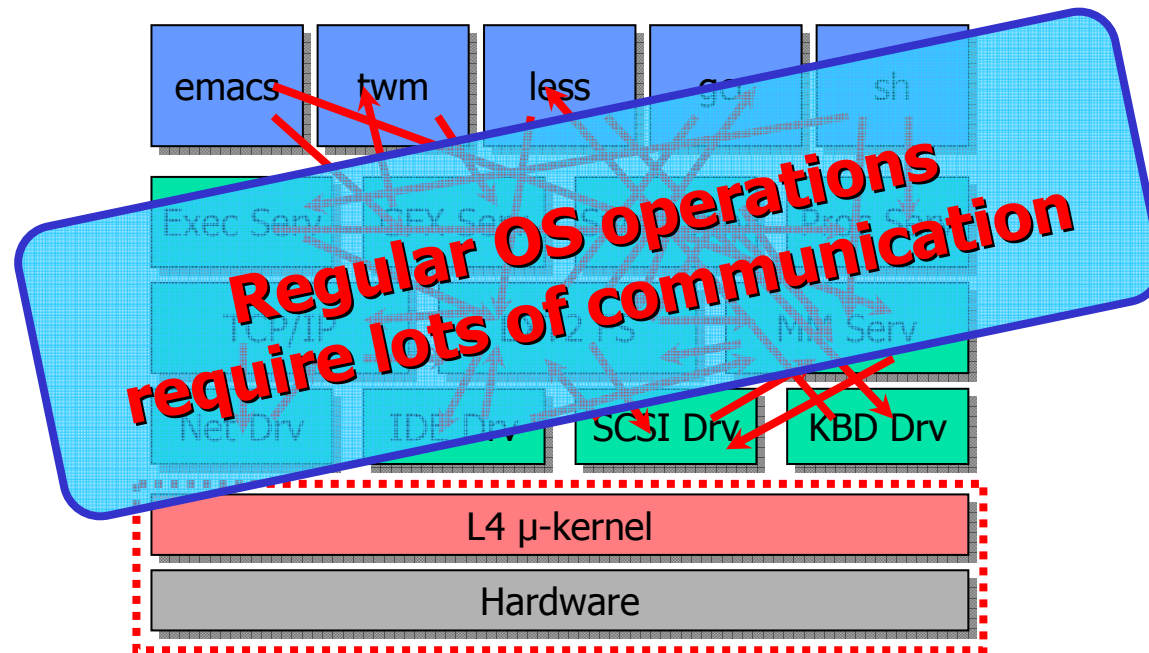
Thin Clients



Specialized Systems



Microkernel Based Systems: The Challenge





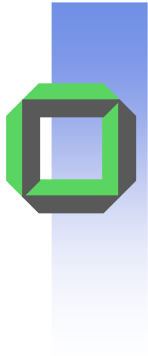
The Great Promise

The Big Disaster

- Coexistence of different
 - APIs
 - File systems
 - OS personalities
- Flexibility
- Extensibility
- Simplicity
- Maintainability
- Security
- Safety

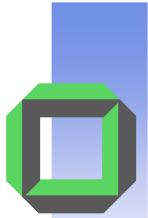
- ***SLOW***
- ***INFLEXIBLE***
- ***LARGE***

**IBM WorkPlace OS:
~2,000,000,000 US\$**

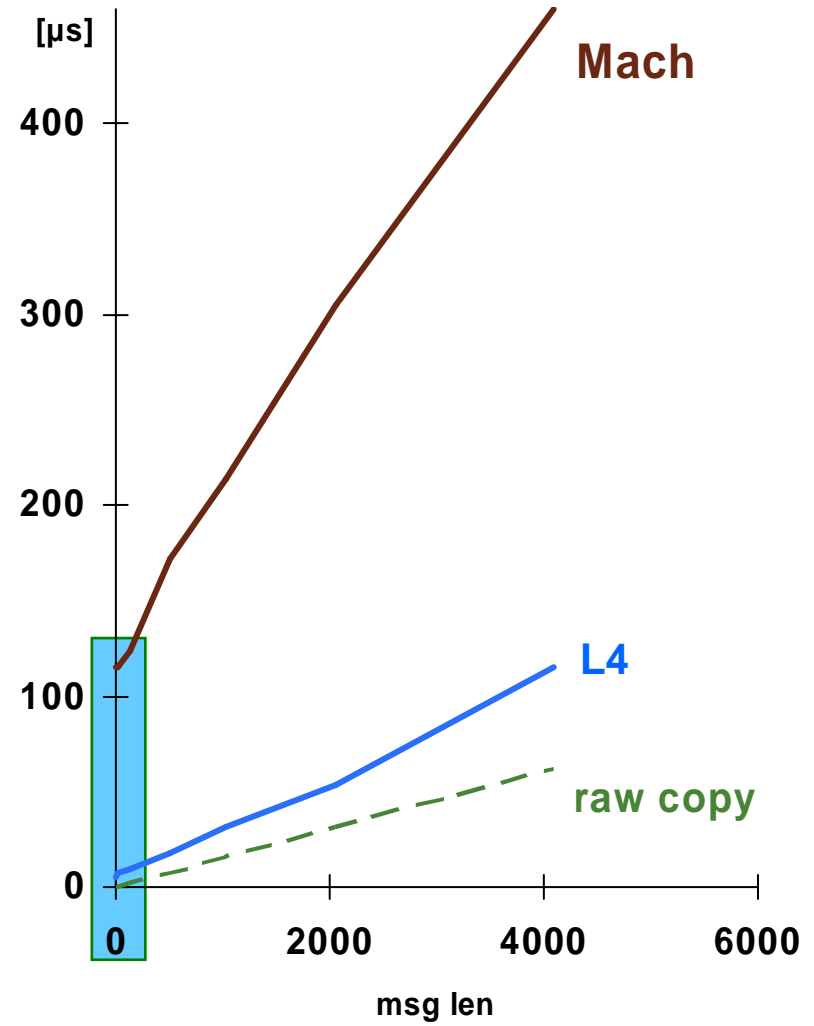
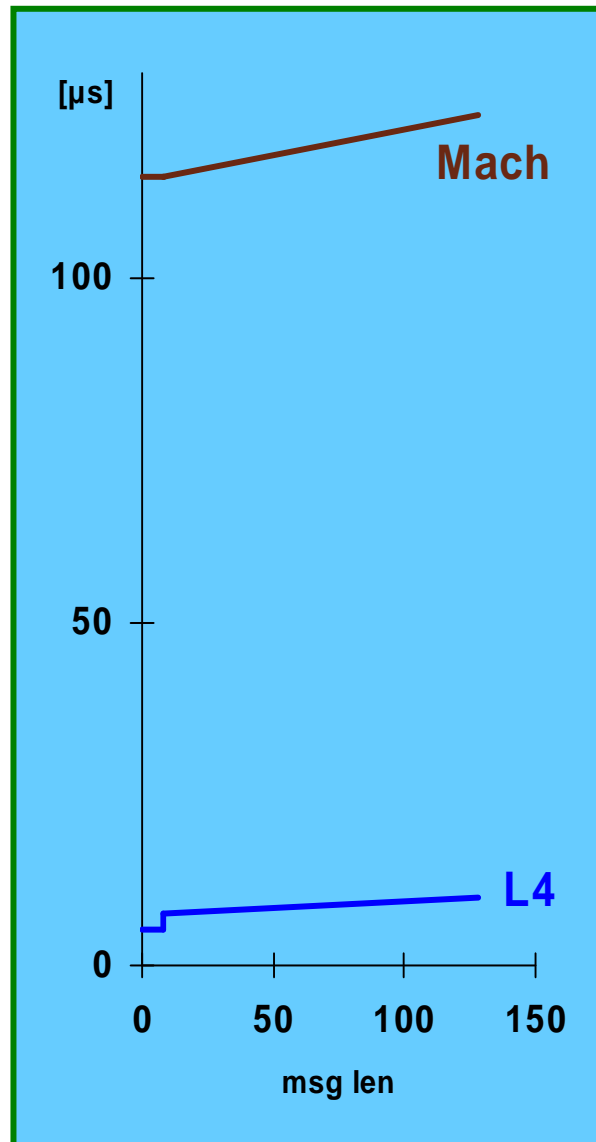


The 100- μ s Disaster

25 MHz 386 → 50 MHz 486 → 90 MHz Pentium → 133 MHz Alpha

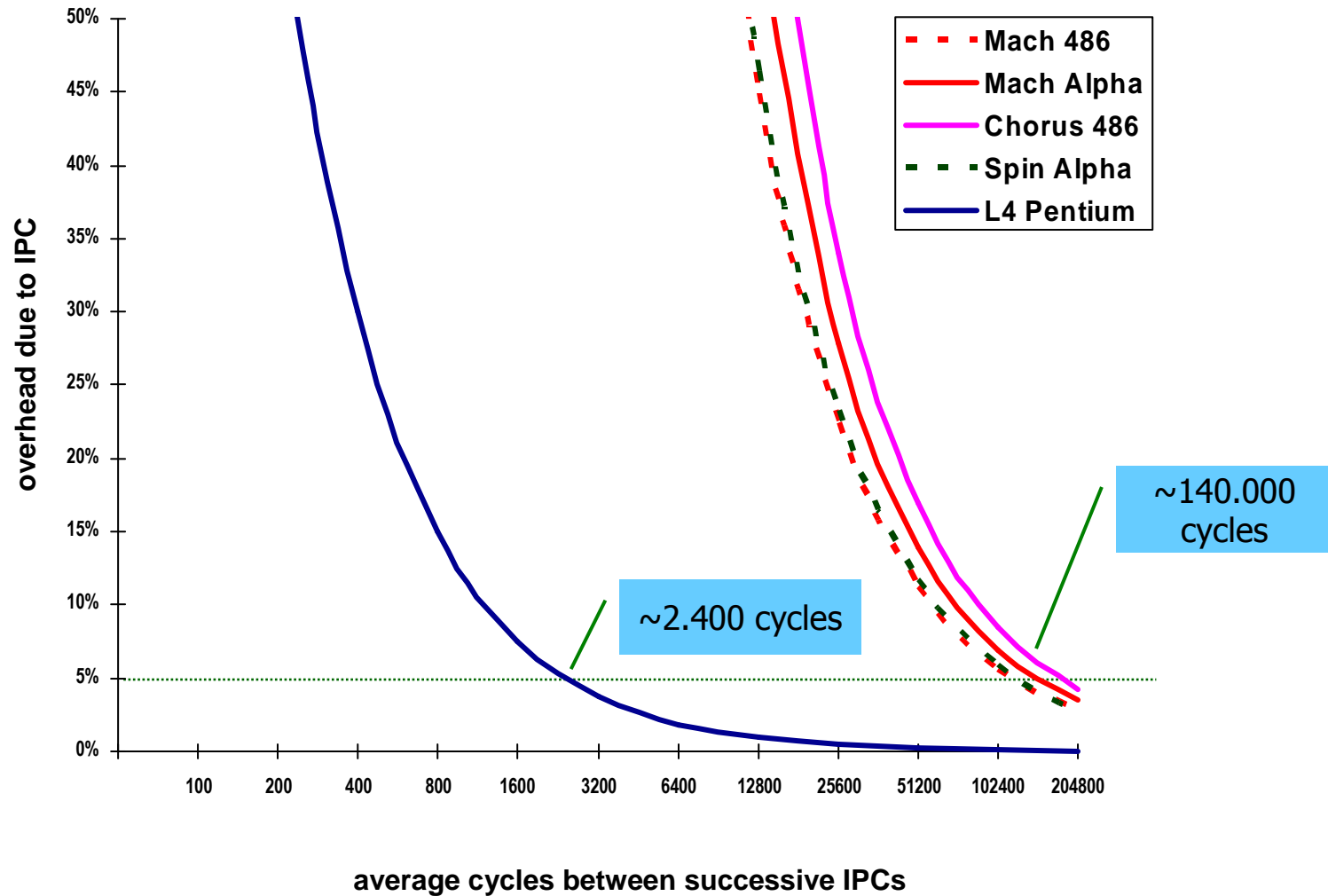


IPC Costs (486, 50 MHz)



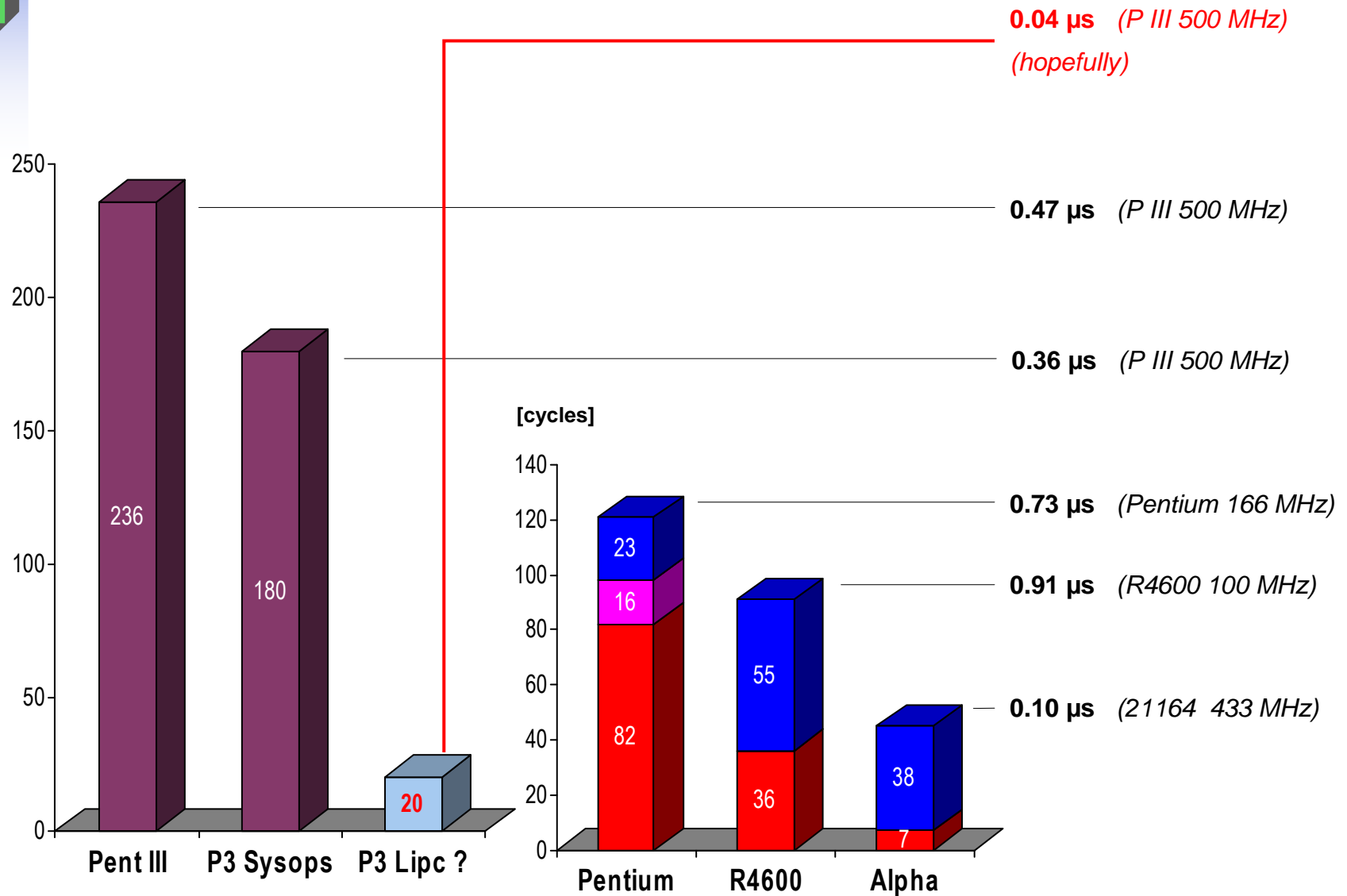


Overhead due to IPC



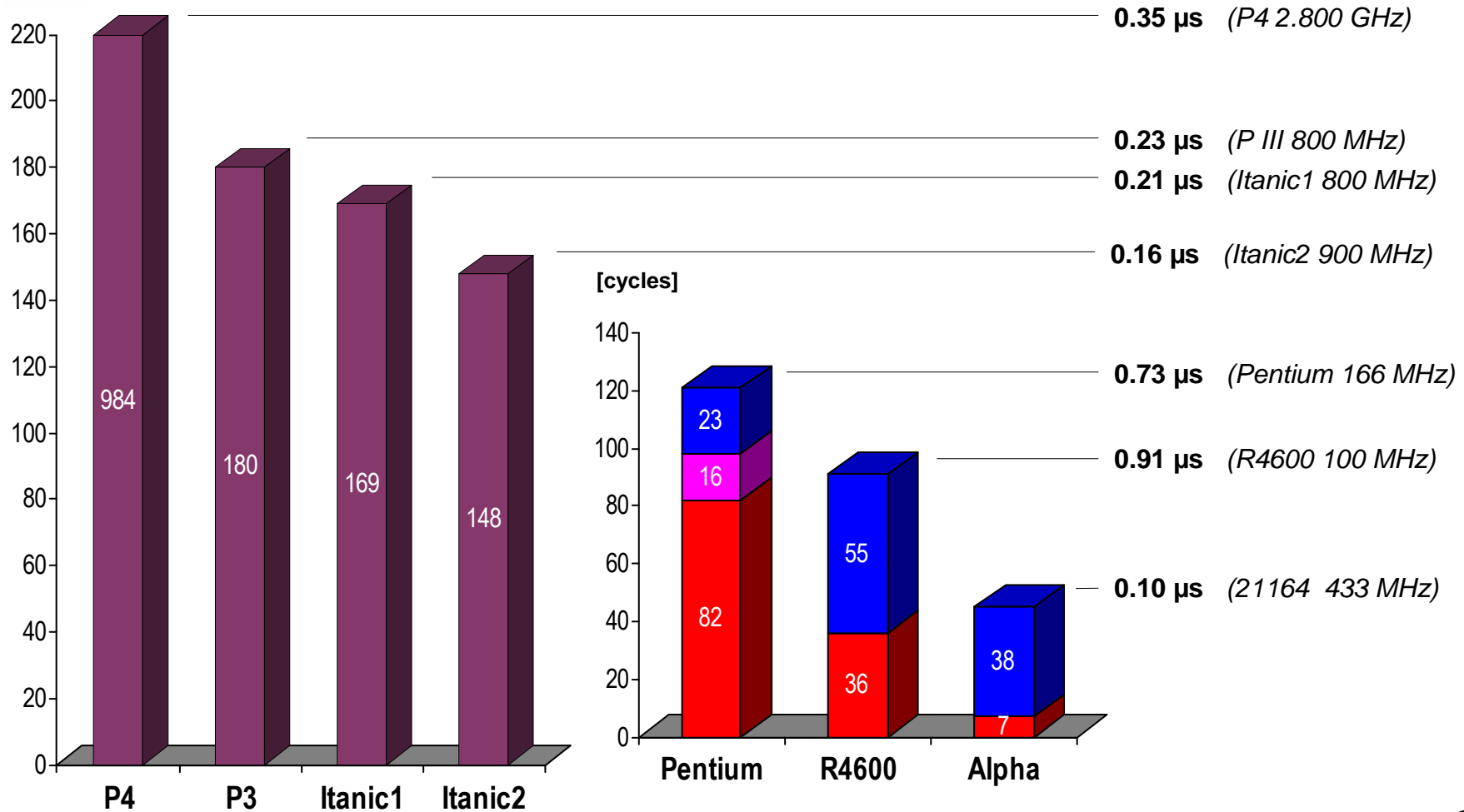


L4 IPC: '95—'00



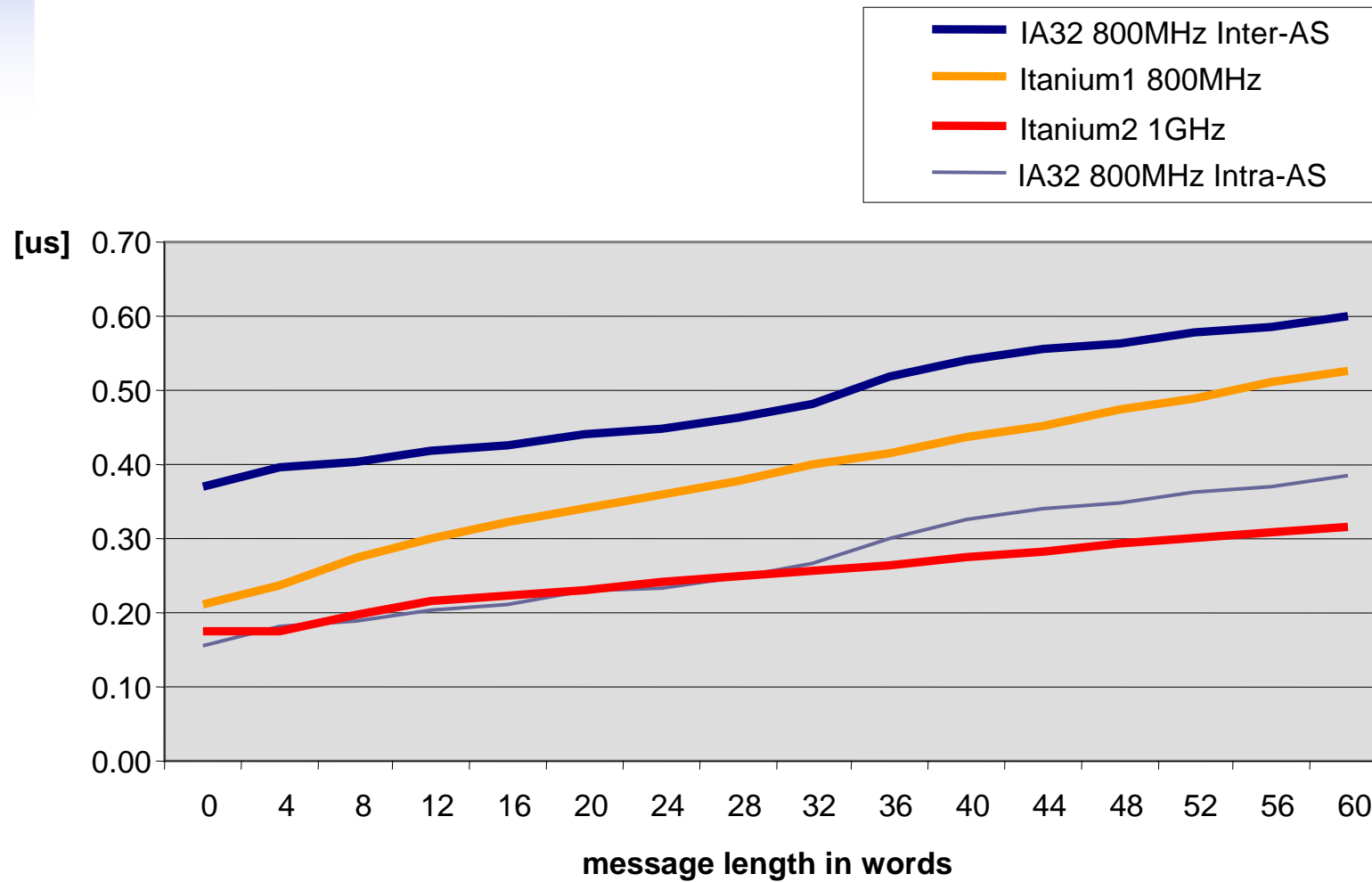


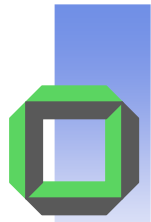
L4 IPC Performance





L4Ka::Pistachio (today)



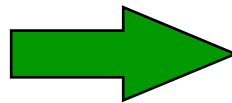


Thesis

A μ -kernel does the job if

- properly designed and
- carefully implemented.

- ◆ **Minimality**
- ◆ **Architectural Integration**
- ◆ **Elegance**



- ◆ **Efficiency**
- ◆ **Flexibility**



Thesis

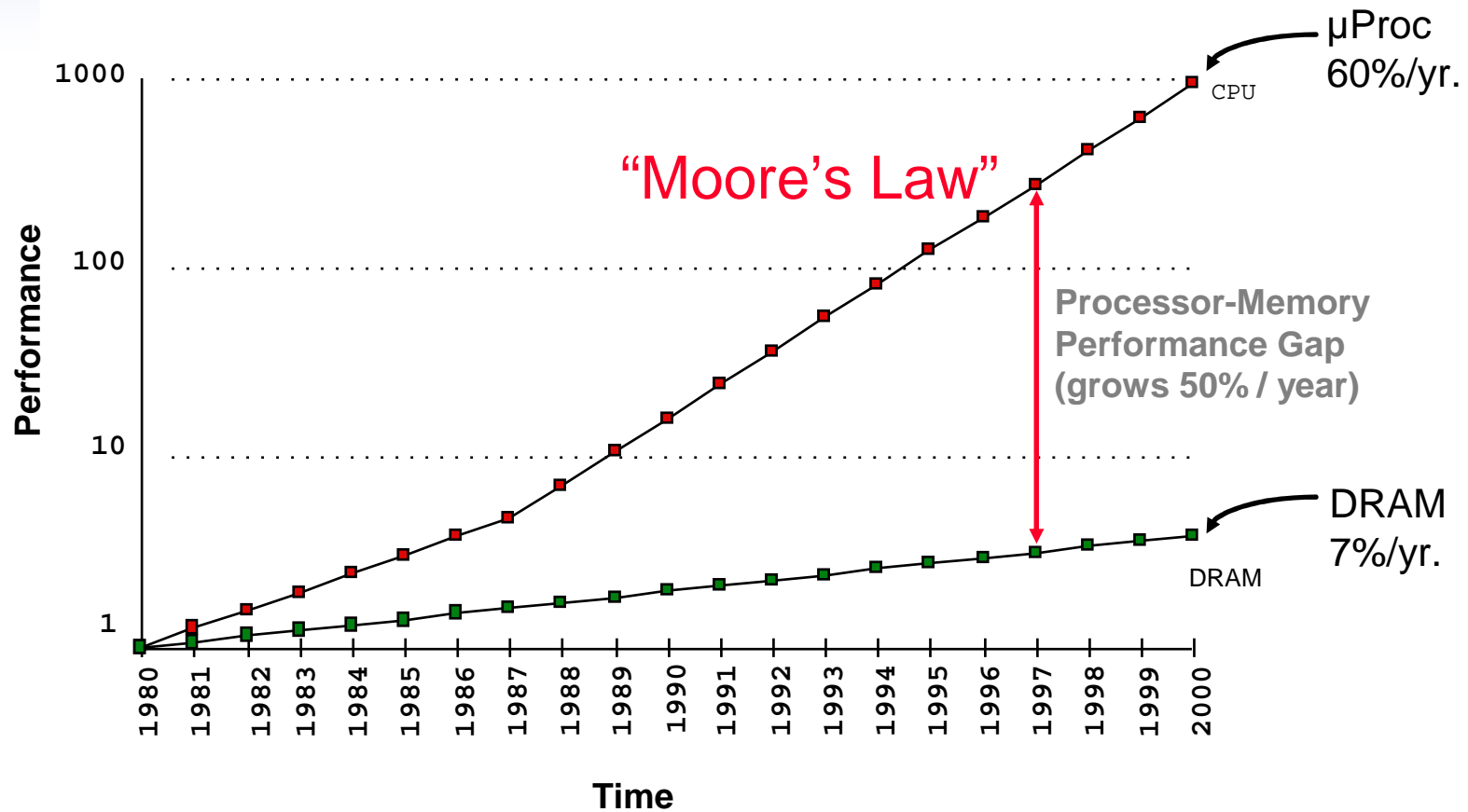
A μ -kernel does the job if

- properly designed and
- carefully implemented.

**When analyzing IPC performance,
cycles are not the only thing to consider!**



Processor-DRAM Gap (Latency)



Slide originally from
Dave Patterson, Parcon 98



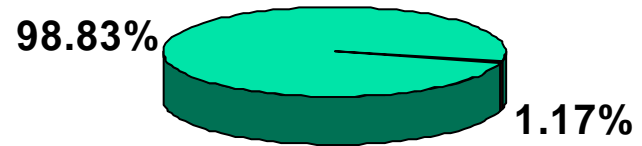
Today's Situation: Microprocessor

- Microprocessor-DRAM performance gap
 - Time of a full cache miss in instructions executed
 - 1st Alpha (7000): $340 \text{ ns} / 5.0 \text{ ns} = 68 \text{ clks} \times 2 \text{ instr.} = 136 \text{ instr.}$
 - 2nd Alpha (8400): $266 \text{ ns} / 3.3 \text{ ns} = 80 \text{ clks} \times 4 \text{ instr.} = 320 \text{ instr.}$
 - 3rd Alpha (t.b.d.): $180 \text{ ns} / 1.7 \text{ ns} = 108 \text{ clks} \times 6 \text{ instr.} = 648 \text{ instr.}$
 - $1/2X \text{ latency} \times 3X \text{ clock rate} \times 3X \text{ instr/clock} \Rightarrow \approx 4.5X$
- Minimize kernel-induced cache misses!

Slide originally from
Dave Patterson, Parcon 98

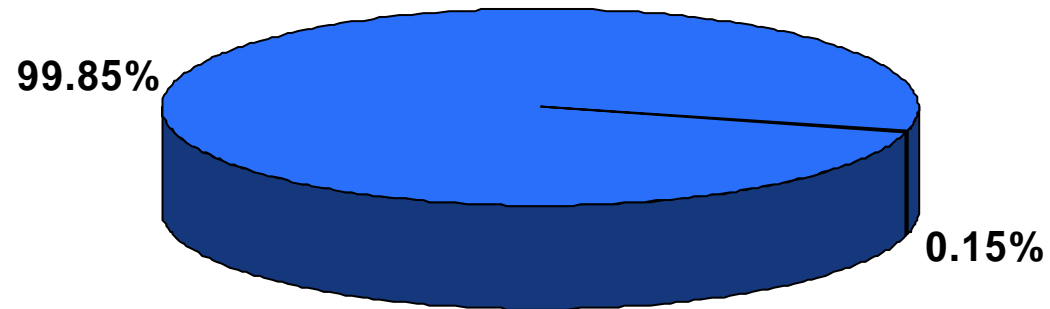


Cache Working Sets



L1 cache

- 1024 cache lines (16K + 16K)
- 12 lines used for IPC



L2 cache

- 8192 cache lines (256K)
- 12 lines used for IPC



Multi-Processor Architectures

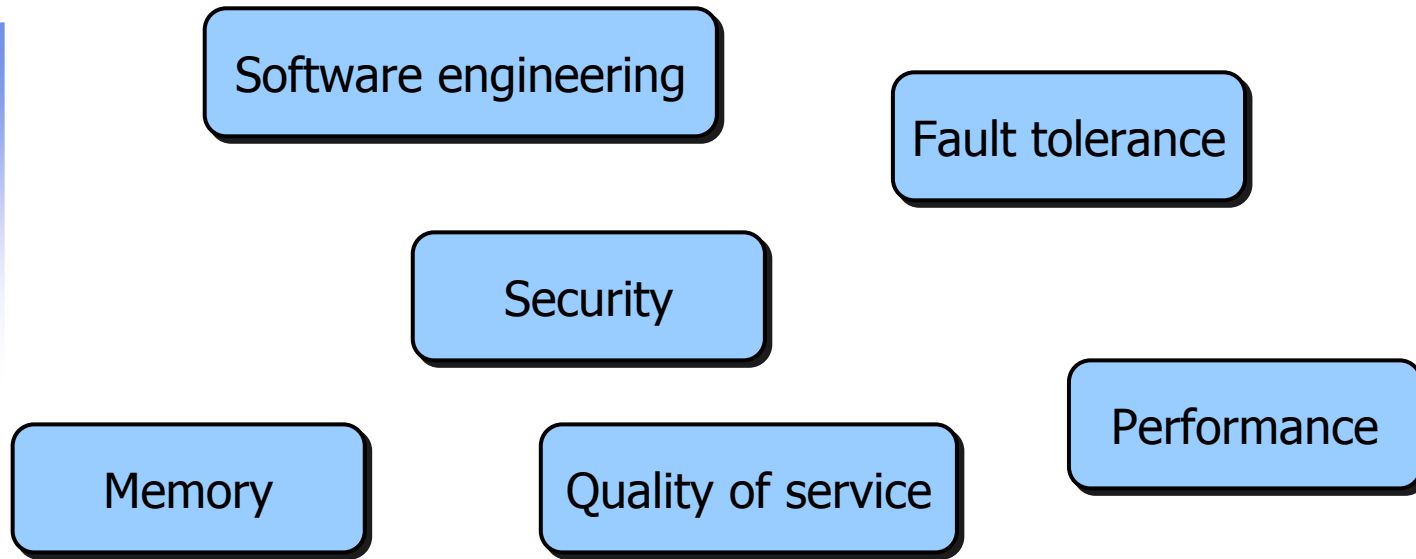
- Synchronization
 - Bus locks
 - Inter-processor interrupts
- NUMA behavior
 - AMD Opteron
 - Simultaneous multithreading (HyperThreading)



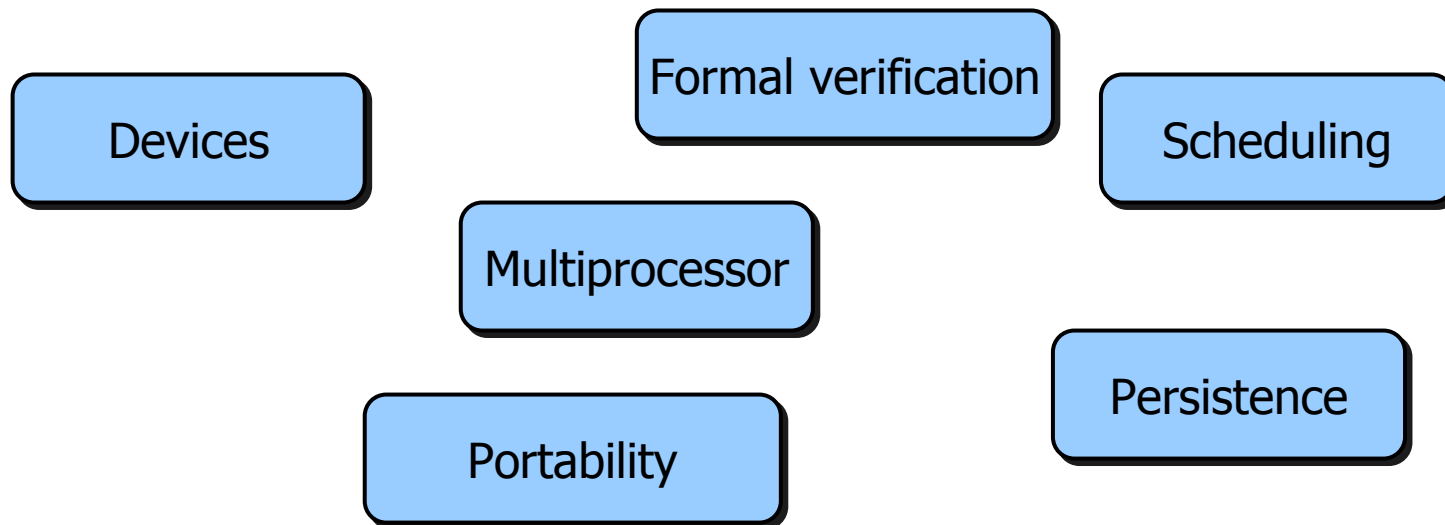
A Perfect Microkernel?

Proving **minimality**, **necessarity** and **completeness** would be nice but is impossible, since there is no agreed-upon metric and all is Turing-equivalent.

Jochen Liedtke, *On μ -Kernel Construction*, SOSPP '95



Small Kernel \neq Small Problem





μ-Kernel Design

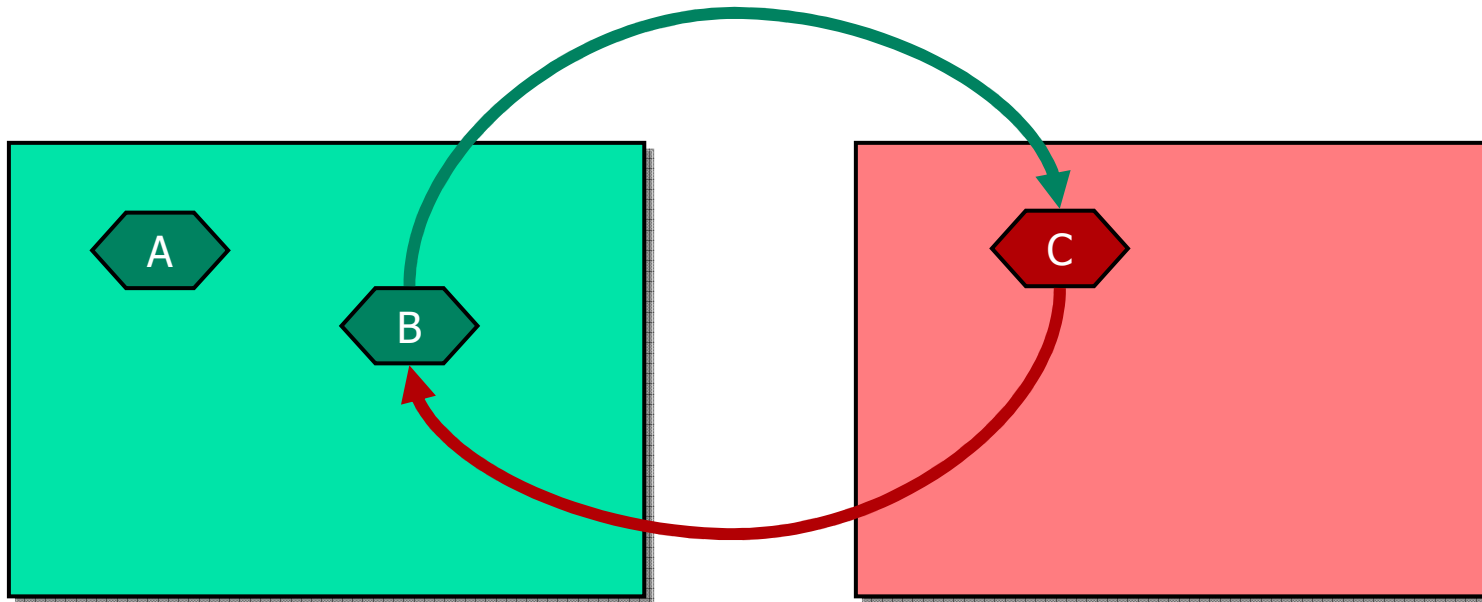
- A μ-kernel does no real work
 - μ-kernel services are only required to overcome μ-kernel constraints
- Therefore, μ-kernels have to be infinitely fast!

Minimality is the key!

- ◆ **Threads** *IPC*
- ◆ **Address Spaces** *Mapping*



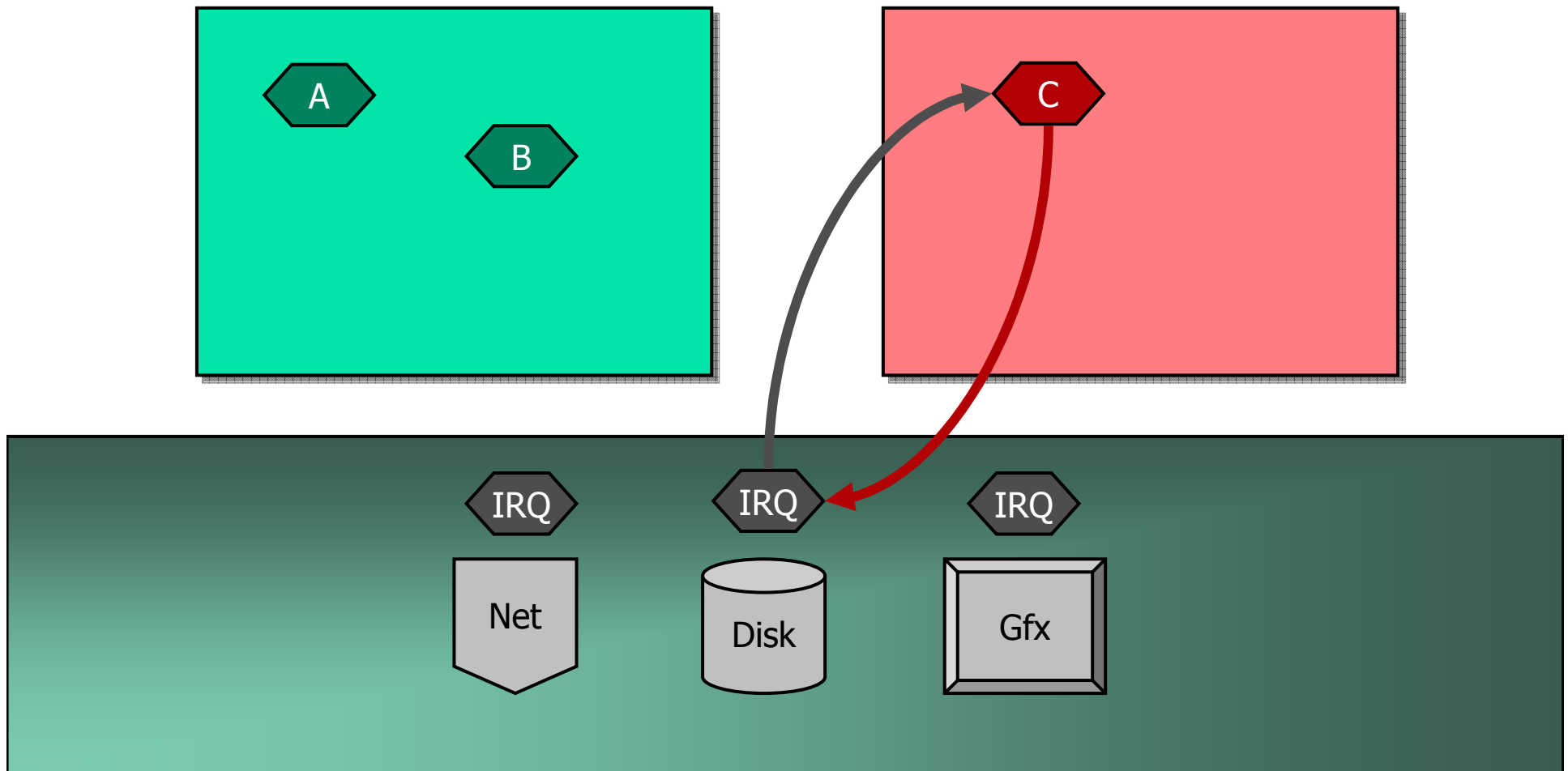
Threads – IPC



- Enable controlled communication across address space boundaries

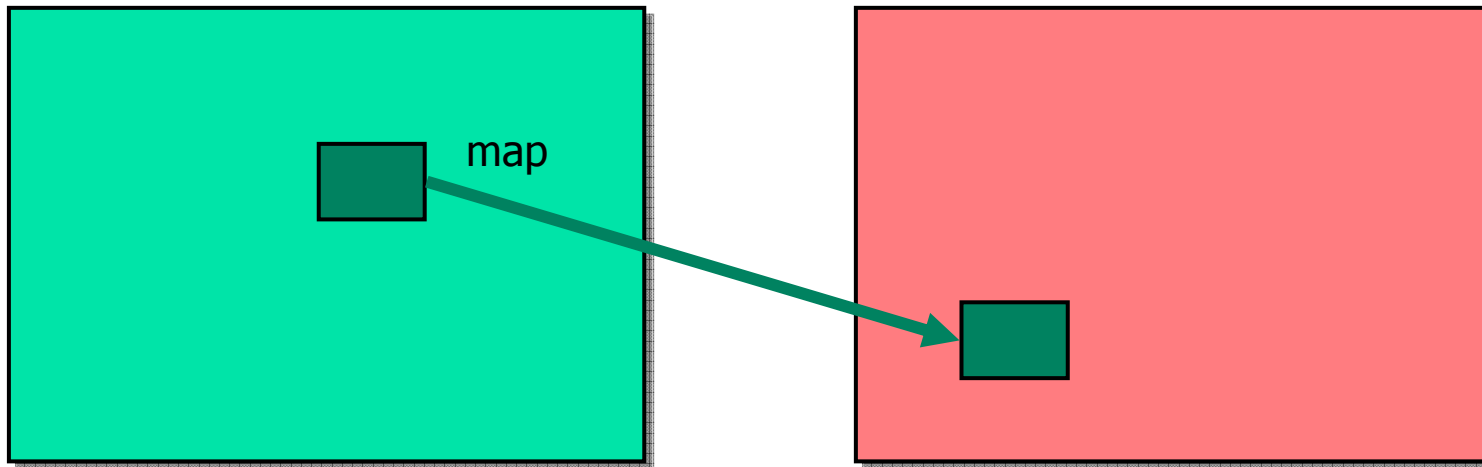


User-Level Device Drivers





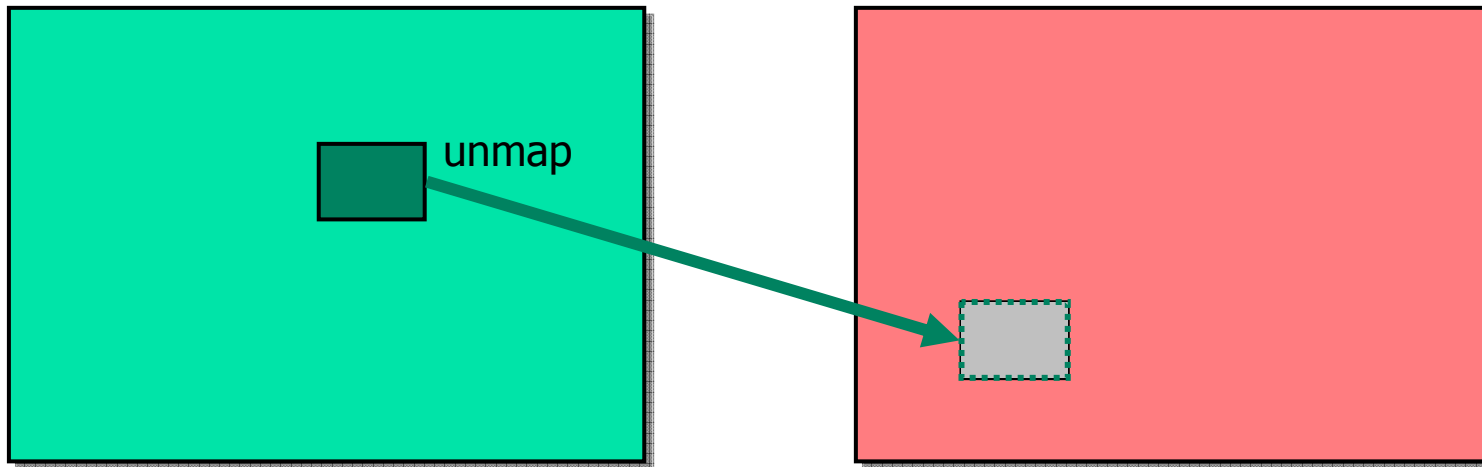
Address Spaces – Mapping



- Setup shared memory regions



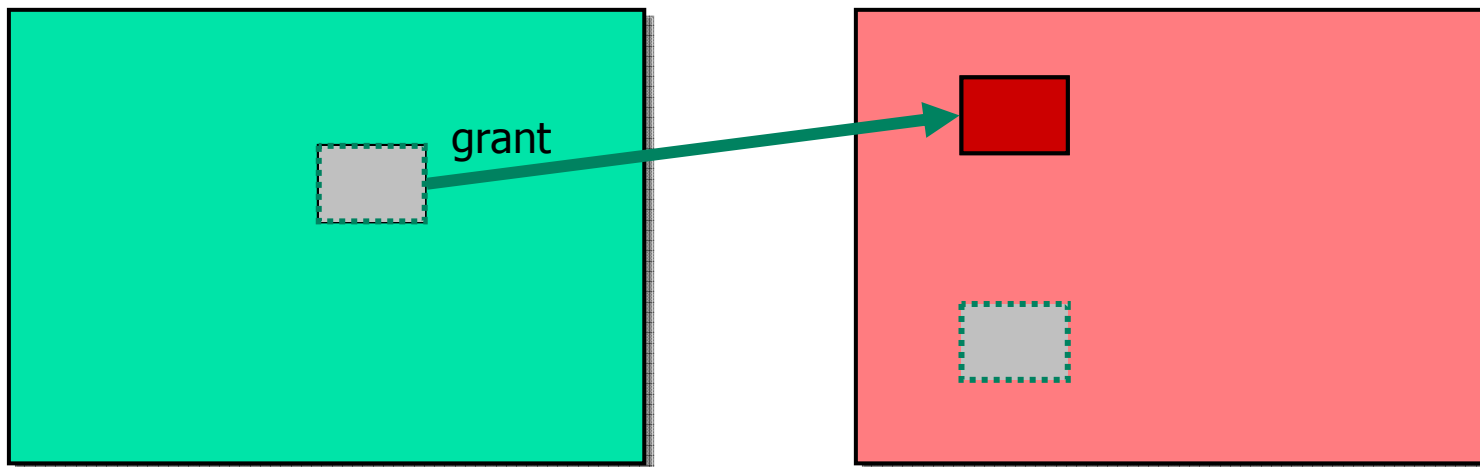
Address Spaces – Mapping



- Revoke shared memory regions



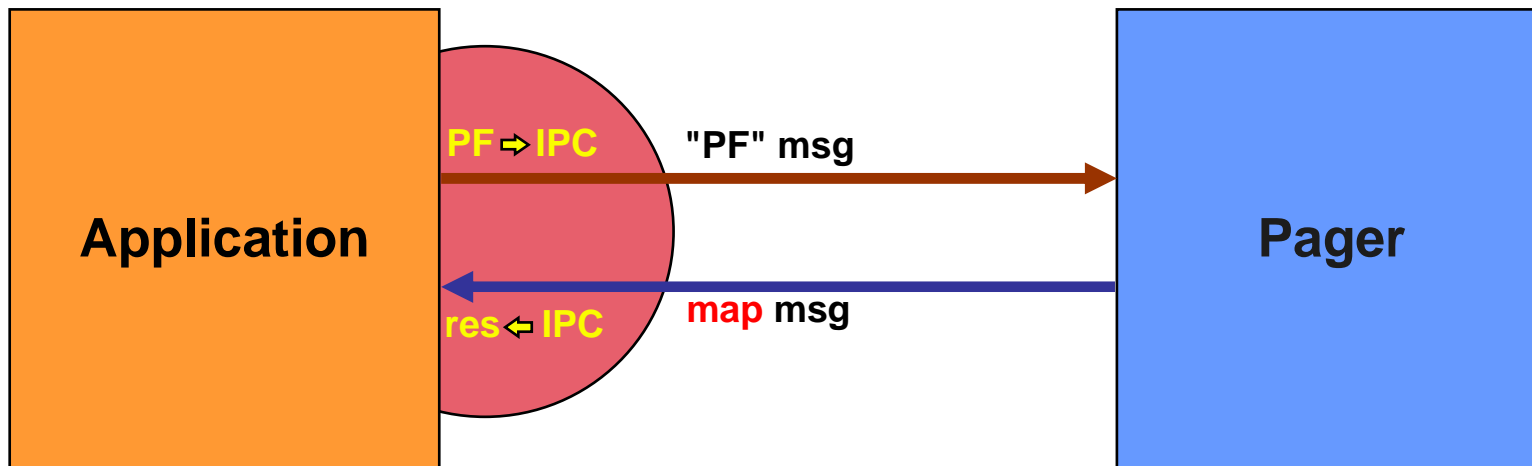
Address Spaces – Mapping



- Donate memory regions to others
- Frees up virtual memory in the granting space
 - Required for file servers, ... (at least useful)

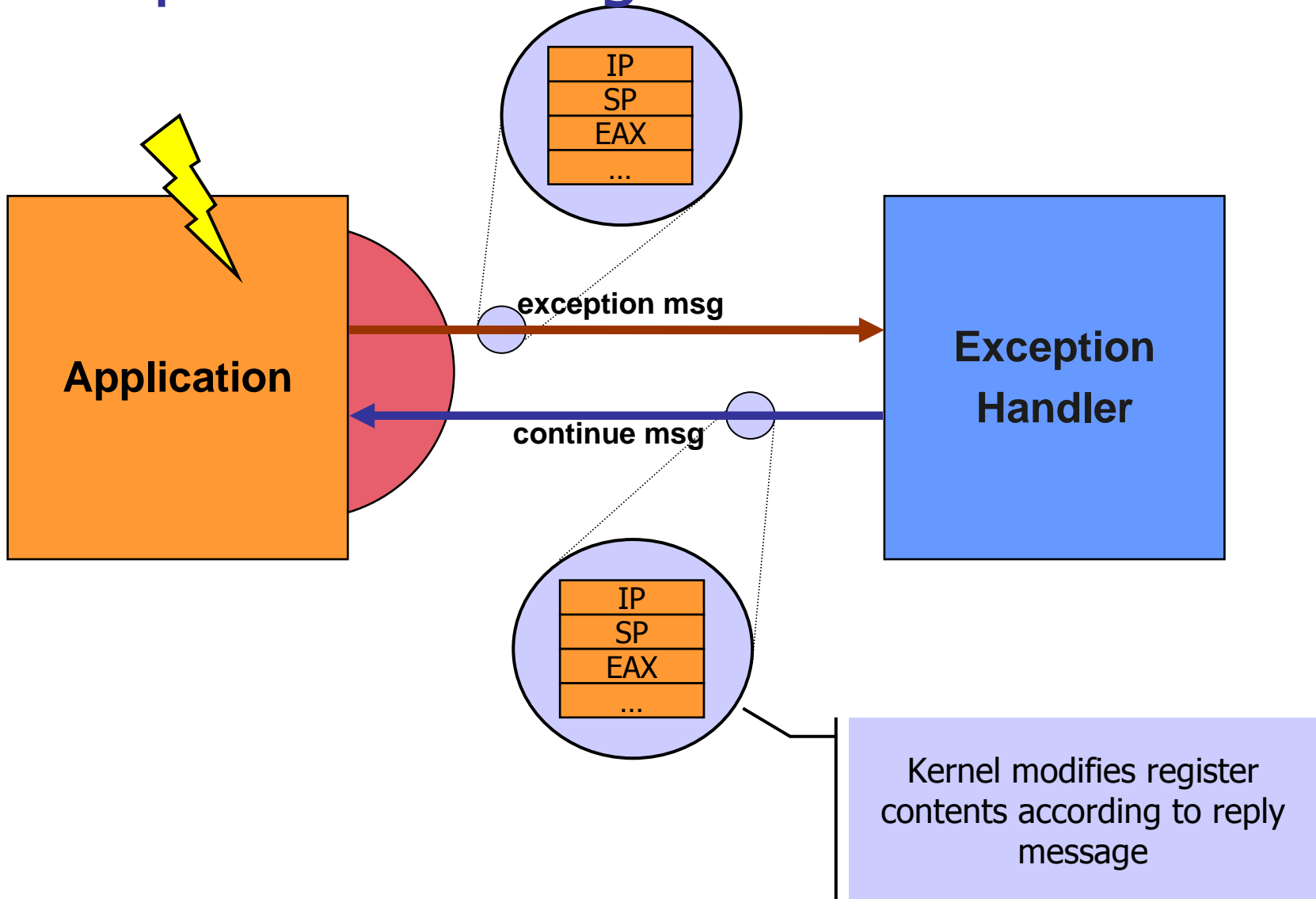


Page Fault Handling



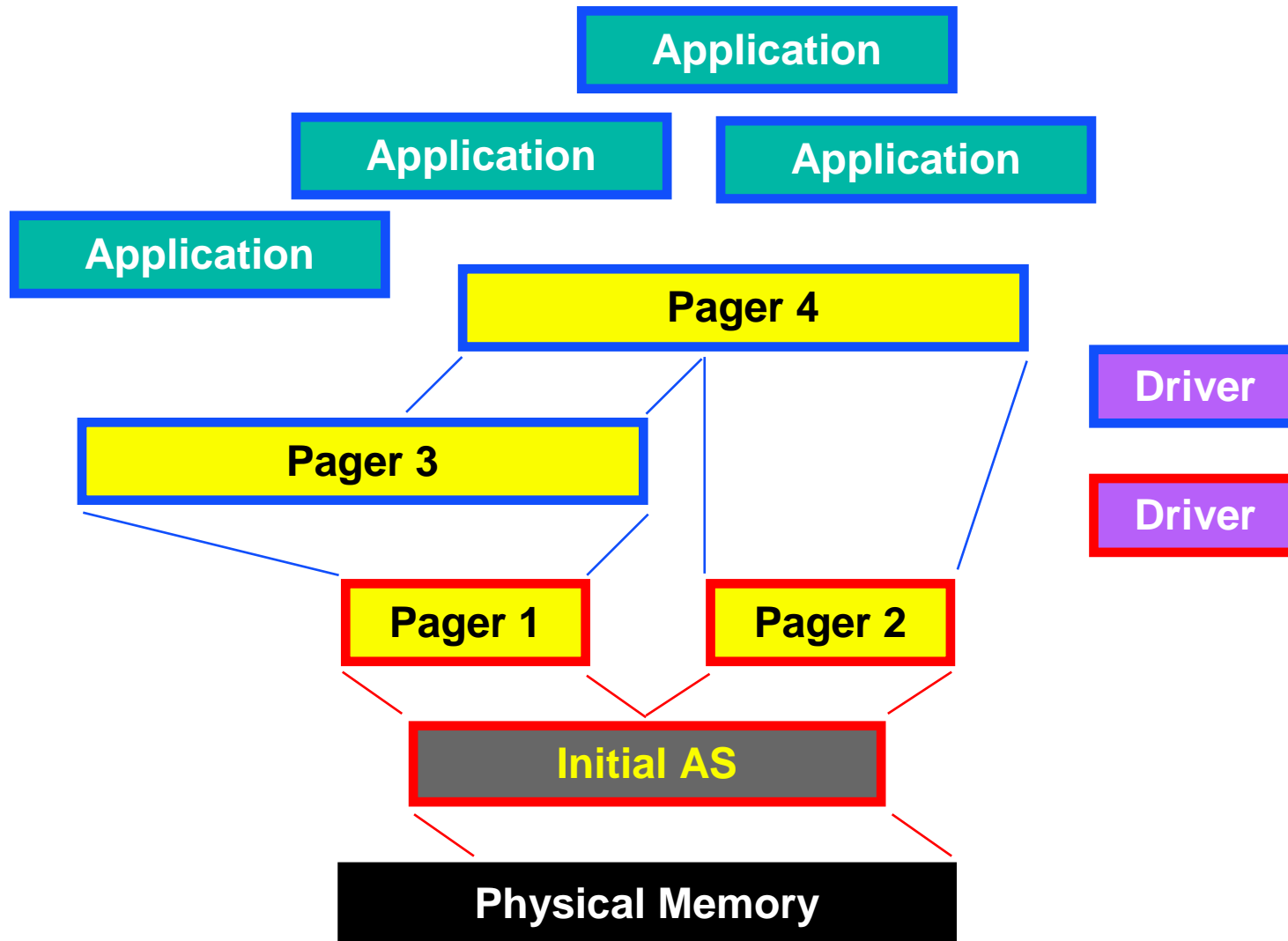


Exception Handling





Recursive Address Spaces





Mach Virtual Memory

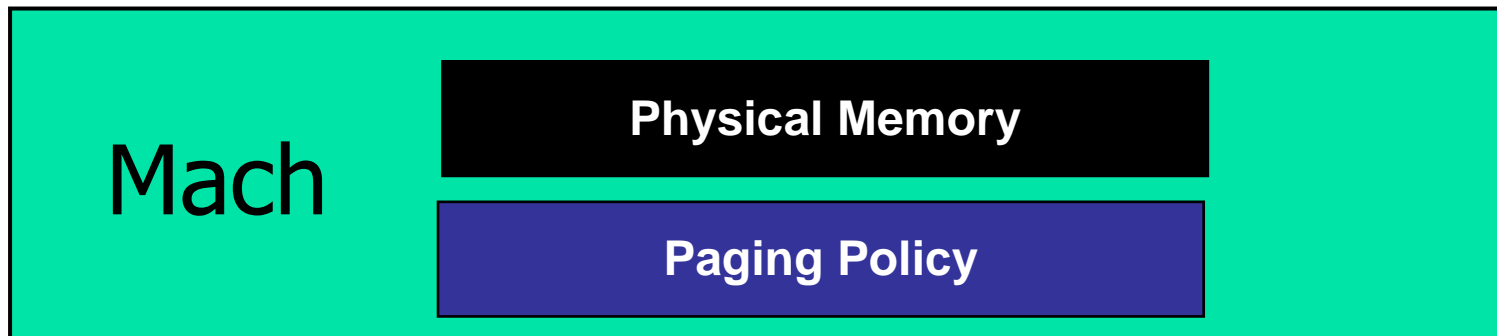
Application

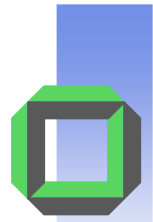
Application

Application

Inflexible

External Page





Tolerated Concepts

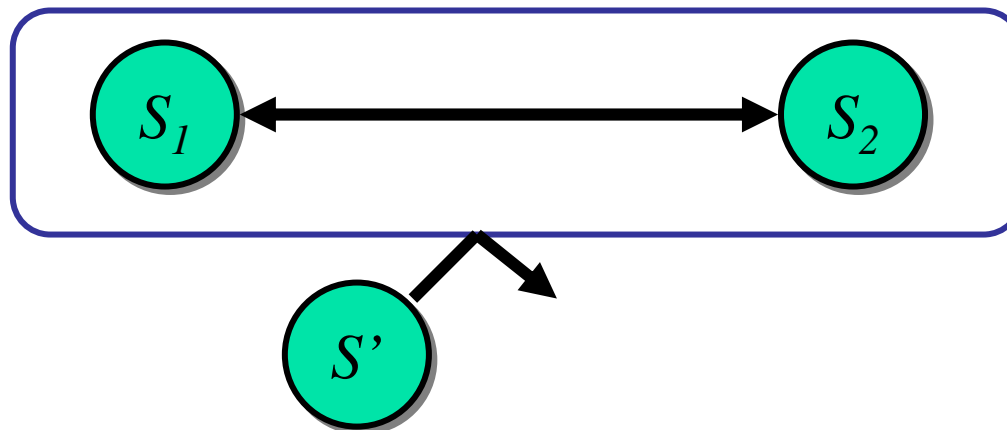
A concept is tolerated in the μ -kernel if ...

competing user-level implementations
would violate system requirements.



Functional Requirements

- Principle of **independence**
 - Subsystem S provides guarantees independent of S'
- Principle of **integrity**
 - Other subsystems can rely on independence guarantees





Requirement: Address Spaces

- μ -kernel must hide hardware address spaces
 - Otherwise violates integrity principle
- But must permit arbitrary protection schemes [and non-protection]
- **Solution**: recursive construction of address spaces outside the kernel



Requirement: Threads

- A **thread** τ is an activity inside an address space with
 - registers
 - instruction pointer
 - stack pointer
 - state information
- $\sigma(\tau) :=$ address space of thread τ



Why Tolerate Threads in μ -kernels?

- The decisive reason: $\sigma(\tau)$
- Modifications to address spaces
$$\sigma(\tau) := \sigma'$$
must be controlled by the kernel
- Thus a notion of τ that represents the above activity
- Additionally: **concurrency**



Requirement: IPC

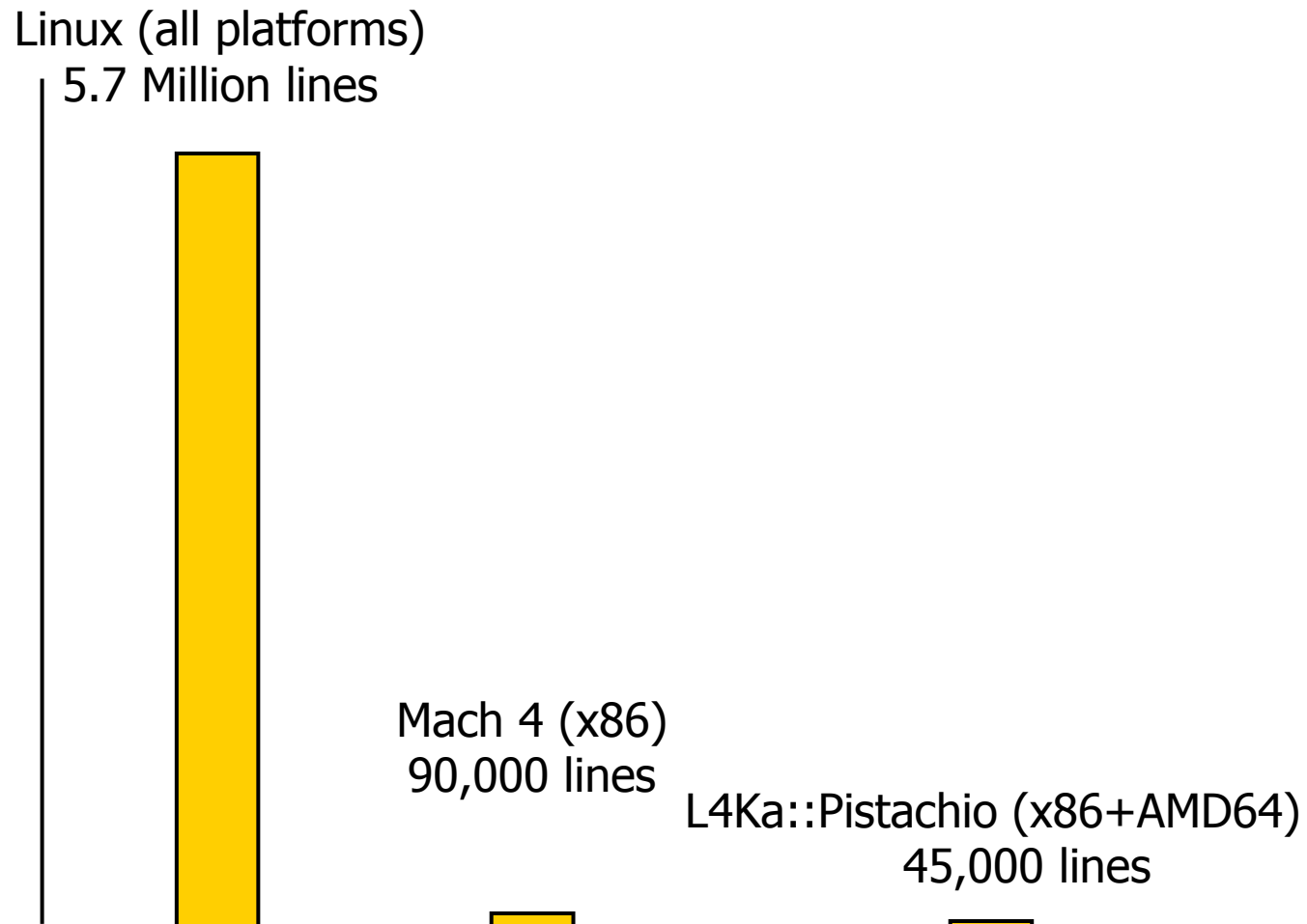
- IPC = inter-process communication
- Inherently required in μ -kernel

- Classical approach
 - Transfer messages between threads

- Contractual
 - Sender determines what to send
 - Receiver agrees to receive the information



Size Comparison





In this course ...

You **do** learn

- How to design a μ K
- Why L4 is sooooo fast
- Reasons why others failed
- Nitty-gritty details about IA32
- Some OS bashing ...
- More cool stuff ...

You **don't** learn

- How to construct a system on a μ K (\rightarrow SDI)
- Linux dos and don'ts
- Operating system X is better than Y



Plan

- Overview, Motivation, Problems
- Threads, System-calls, and Thread Switching
- TCBs and Address Space Layout
- IPC Functionality and Implementation
- Dispatching
- Small Address Spaces - IPC
- Virtual Memory and Mapping Database
- Interrupts, Exceptions and CPU Virtualization
- Security

Many algorithms, often influencing the system design.