

# Operating System Kernels

Presenter: Saikat Guha

Cornell University

CS 614, Fall 2005

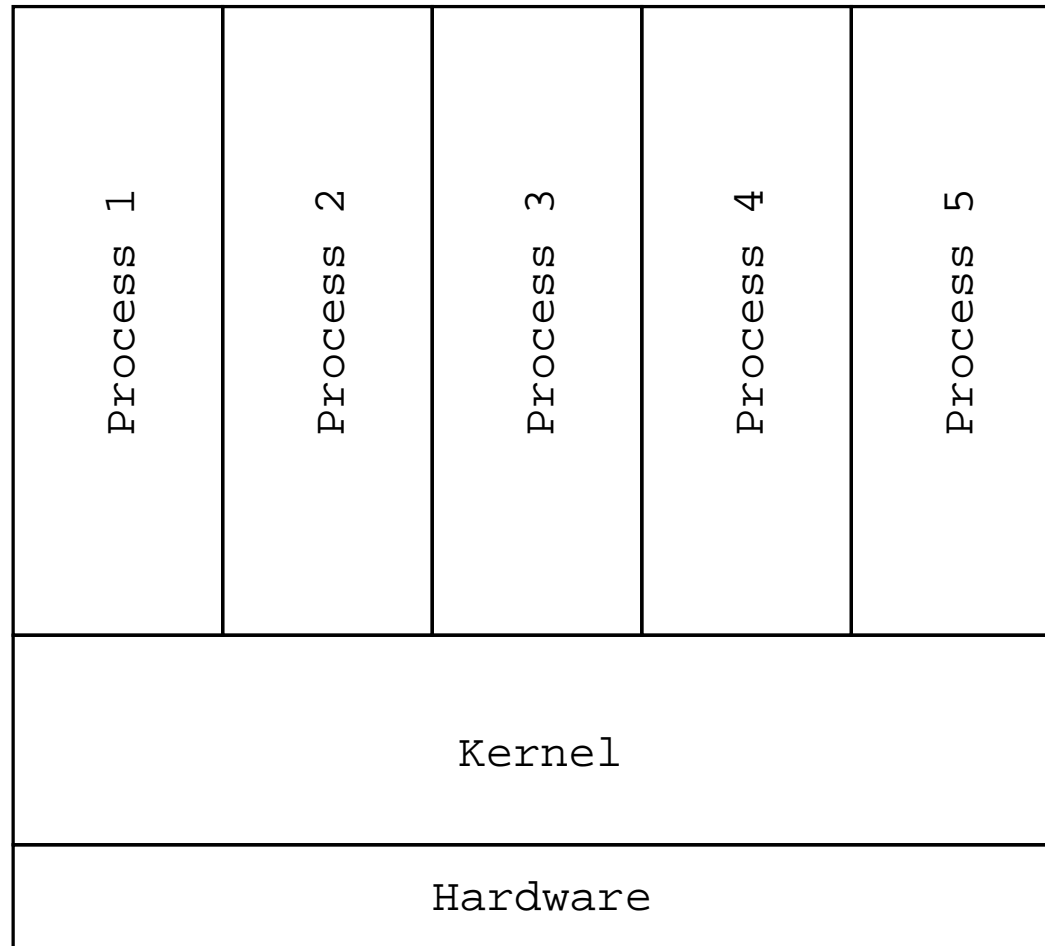
# Operating Systems

- ▶ Initially, the OS was a run-time library
- ▶ Batch ('55–'65): Resident, spooled jobs
- ▶ Multiprogrammed (late '60): Multiple jobs
- ▶ Time-sharing ('70s): Interactive jobs
  - ▶ Multics, UNIX
- ▶ Networked OS, Distributed OS, Parallel OS, Real-time OS

# UNIX

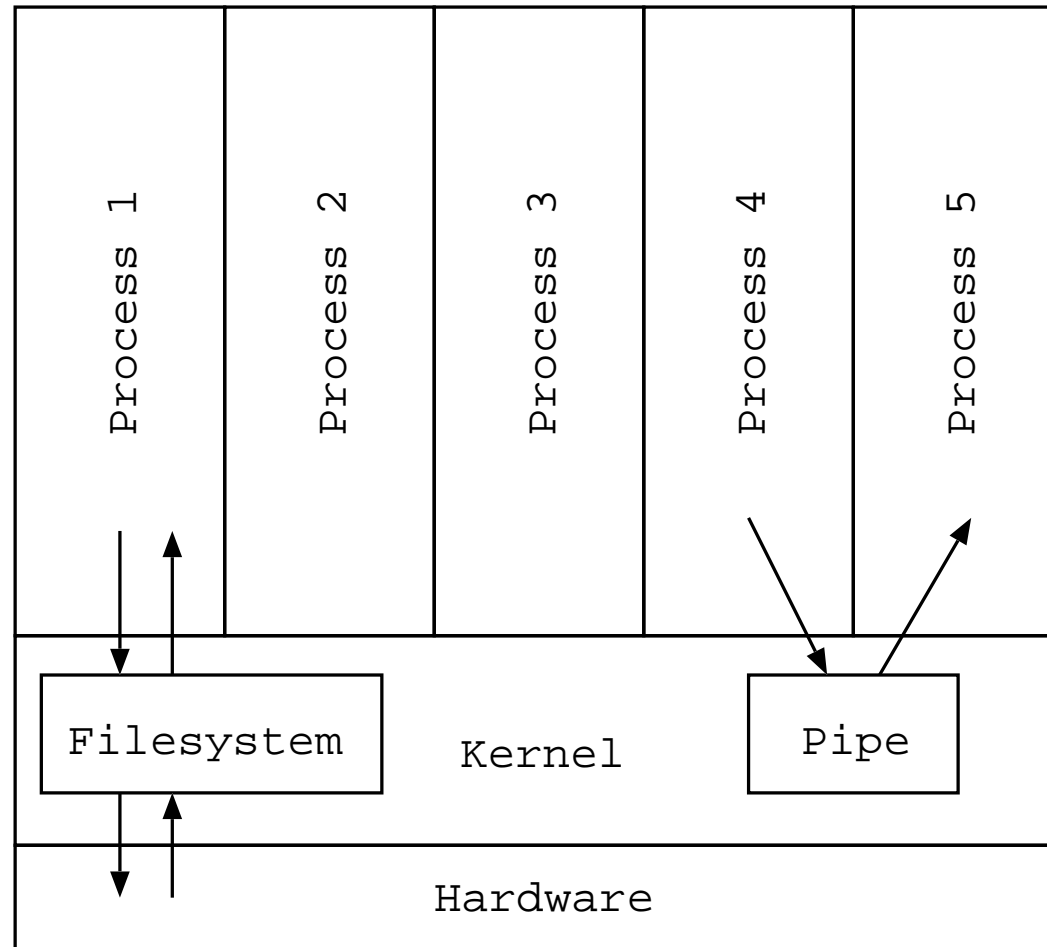
- ▶ **THE** operating system
  - ▶ Dennis Ritchie, Ken Thompson at AT&T
- ▶ “File” Abstraction
- ▶ Kernel
  - ▶ Processes, IPC
  - ▶ Filesystem
  - ▶ Networking (eventually)
  - ▶ Graphics (Windows)
- ▶ Userspace
  - ▶ Shell
  - ▶ Commands

# Operating Systems



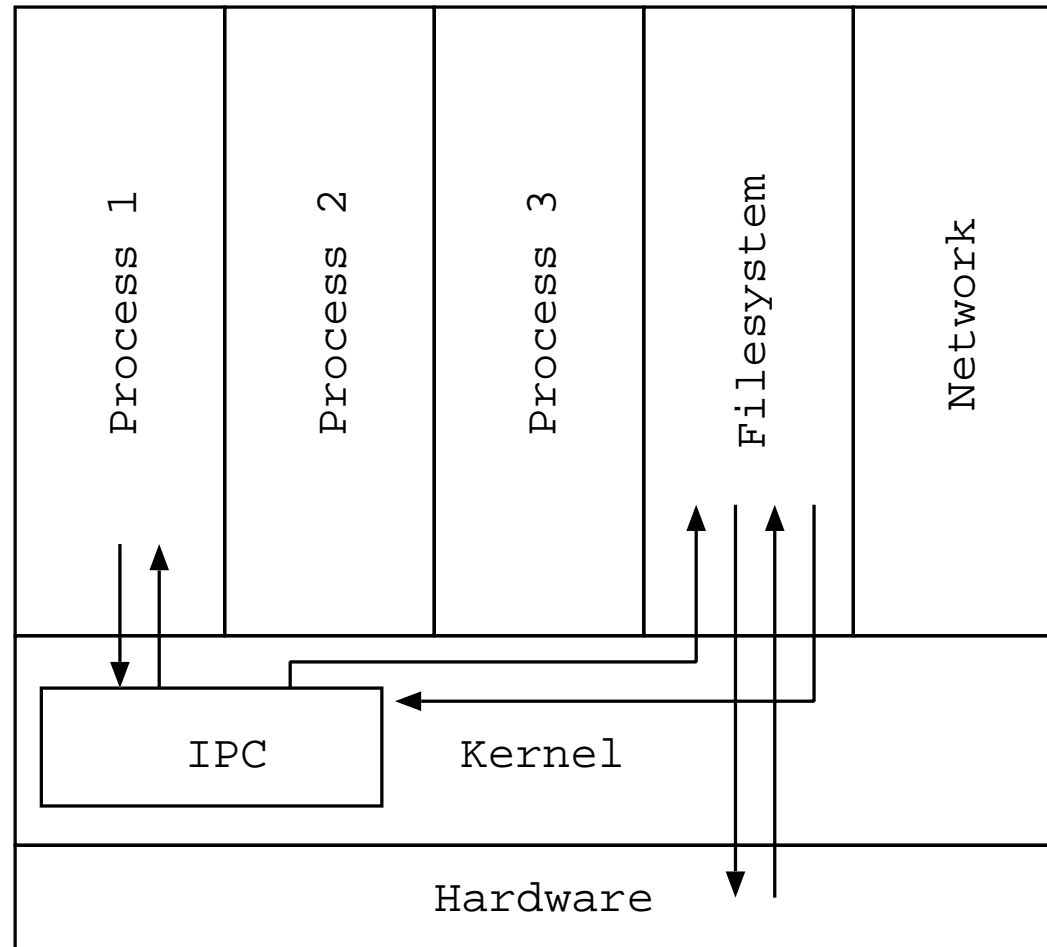
## User-Kernel Split

# Operating Systems



Monolithic Kernel

# Operating Systems



Microkernel

# $\mu$ -Kernels

- ▶ Minimal services
- ▶ Usually threads or processes, address space and inter-process communication (IPC)
- ▶ User-space Filesystem, Network, Graphics, even device drivers sometimes.

# Monolithic Kernels: Advantages

- ▶ Kernel has access to everything
  - ▶ All optimizations possible
  - ▶ All techniques/mechanisms/concepts can be implemented
- ▶ Extended by simply adding more code
  - ▶ Linux at 3.3M lines of code
- ▶ Tackle complexity
  - ▶ Layered kernels
  - ▶ Modular kernels
  - ▶ Object oriented kernels. Do C++, Java, C# help?



# $\mu$ -Kernels: Advantages

- ▶ Minimal
  - ▶ Smaller trusted base
  - ▶ Less error prone
  - ▶ Server malfunction easily isolated
- ▶ Elegant
  - ▶ Enforces modularity
  - ▶ Restartable user-level services
- ▶ Extensible
  - ▶ Different servers/API can coexist

- ▶ 1st generation  $\mu$ -kernels
  - ▶ Mach (CMU)<sup>1</sup>
  - ▶ Chorus (Inria, Chorus systems)
  - ▶ Amoeba (Vrije University)
  - ▶ L3 (GMD)<sup>2</sup>

---

<sup>1</sup>External pager

<sup>2</sup>User-Level Driver

# $\mu$ -Kernels: Problems

- ▶ Overheads
  - ▶ Chen and Bershad, '93
  - ▶ Impact of caches, locality, TLB collisions
  - ▶ Up 66% degradation in Mach
- ▶ Co-located servers for performance
- ▶ Can be optimized to be fast on an architecture
  - ▶ But, performance not preserved on other architectures

# $\mu$ -Kernels

- ▶ 2nd generation  $\mu$ -kernels
  - ▶ Spin (UWash)
  - ▶ Exokernel (MIT)
  - ▶ L4 (GMD/IBM/UKa)<sup>3</sup>

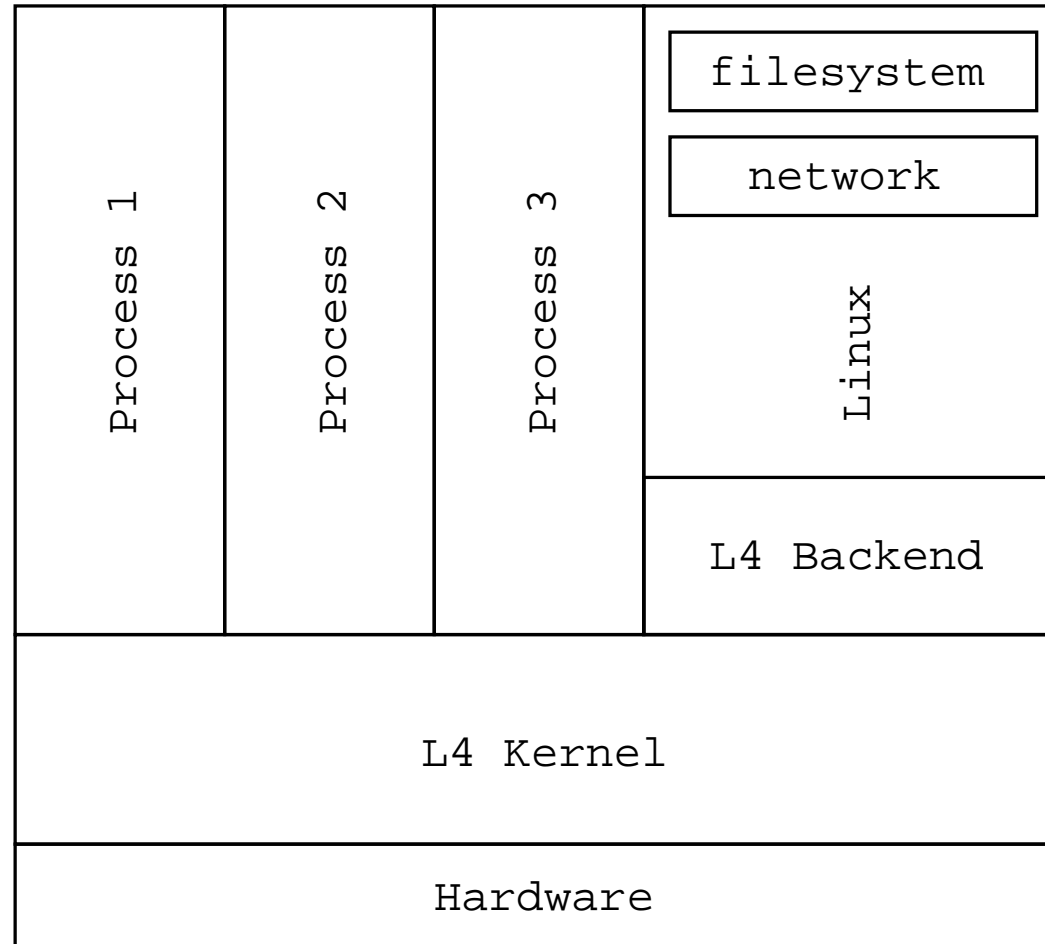
---

<sup>3</sup>User-Level Address Space

# Summary of First Paper

- ▶ The Performance of  $\mu$ -Kernel-Based Systems (Härtig et al., SOSPP '97)
- ▶ Evaluates a L<sup>4</sup>  $\mu$ -kernel based system
- ▶ Ports Linux to run on top of L<sup>4</sup>
- ▶ Suggests improvements

# L<sup>4</sup>-Linux



L<sup>4</sup>-Linux

# L<sup>4</sup>-Linux

- ▶ 2 basic concepts
  - ▶ Threads
  - ▶ Address Spaces (AS)
- ▶ Recursive construction of AS
  - ▶ Grant - Give a page to another AS
  - ▶ Map - Share a page with another AS
  - ▶ Demap - Revoke a mapped or granted page
- ▶ I/O ports treated as AS
- ▶ Hardware interrupts treated as IPC

# L<sup>4</sup>-Linux

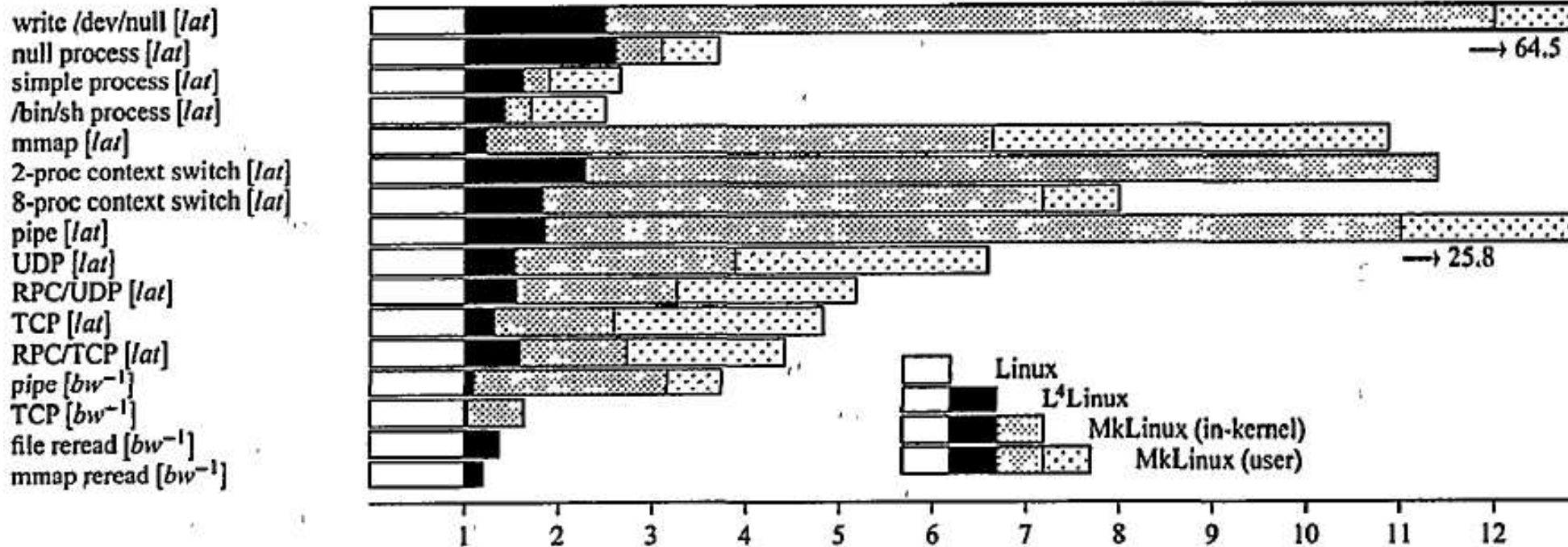
- ▶ TLB caches page-table lookups
  - ▶ Flushed during context switch
  - ▶ Flushing not necessary for tagged TLBs
- ▶ L<sup>4</sup>-Linux avoids frequent flushes
  - ▶ Pentium CPU's emulate tagged TLBs for small address spaces
- ▶ syscall time
  - ▶ Unix – 20 $\mu$ s
  - ▶ Mach – 114 $\mu$ s
  - ▶ L<sup>4</sup> – 5 $\mu$ s



# Performance

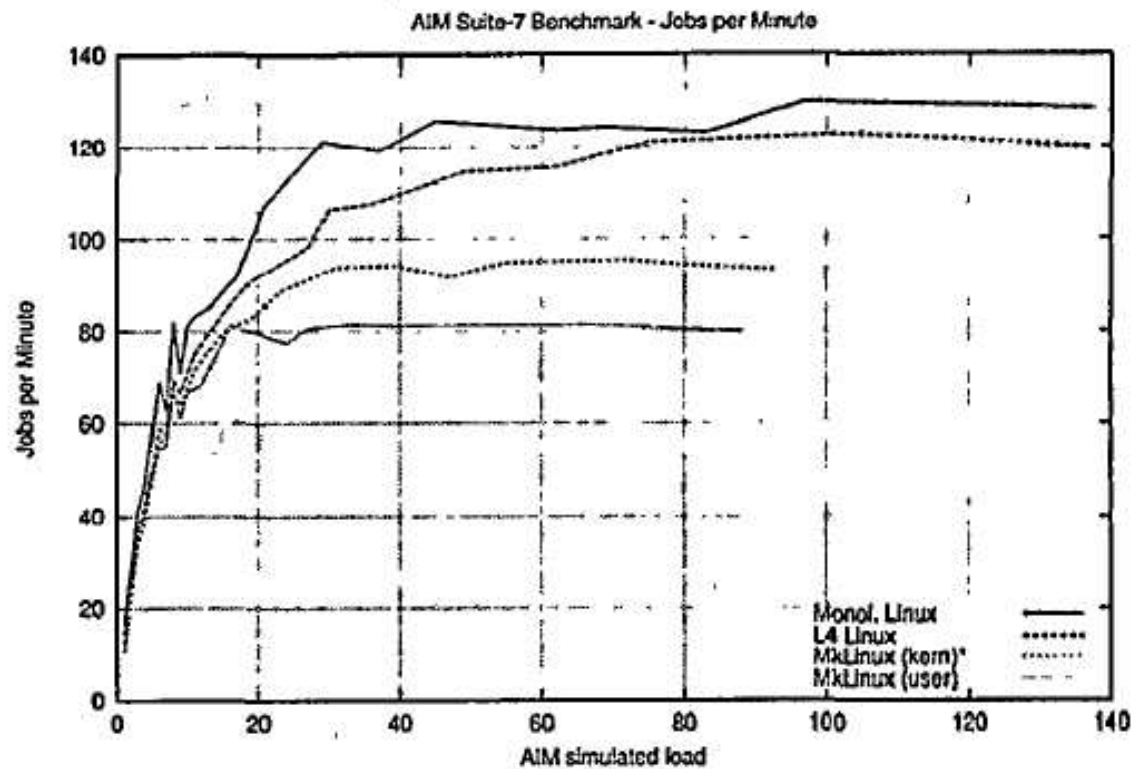
Is L<sup>4</sup>-Linux a practical system?

# Performance



Is L<sup>4</sup>-Linux a practical system? **Yes**

# Performance



Is L<sup>4</sup>-Linux a practical system? **Yes**

# Performance

- ▶ L<sup>4</sup> incurs 5%–10% overhead
- ▶ Collocation alone does not solve performance problems
  - ▶ What about L<sup>4</sup> without collocation?
- ▶ L<sup>4</sup>-Linux is proof-of-concept
  - ▶ Pipes can be made faster
  - ▶ Better VM in non-legacy mode
  - ▶ Can benefit from cache partitioning

# Comparison

## L<sup>4</sup>-Linux

- ▶ Highly optimized (for x86)
- ▶ Functionality limited by Linux
- ▶ Untrusted components isolated

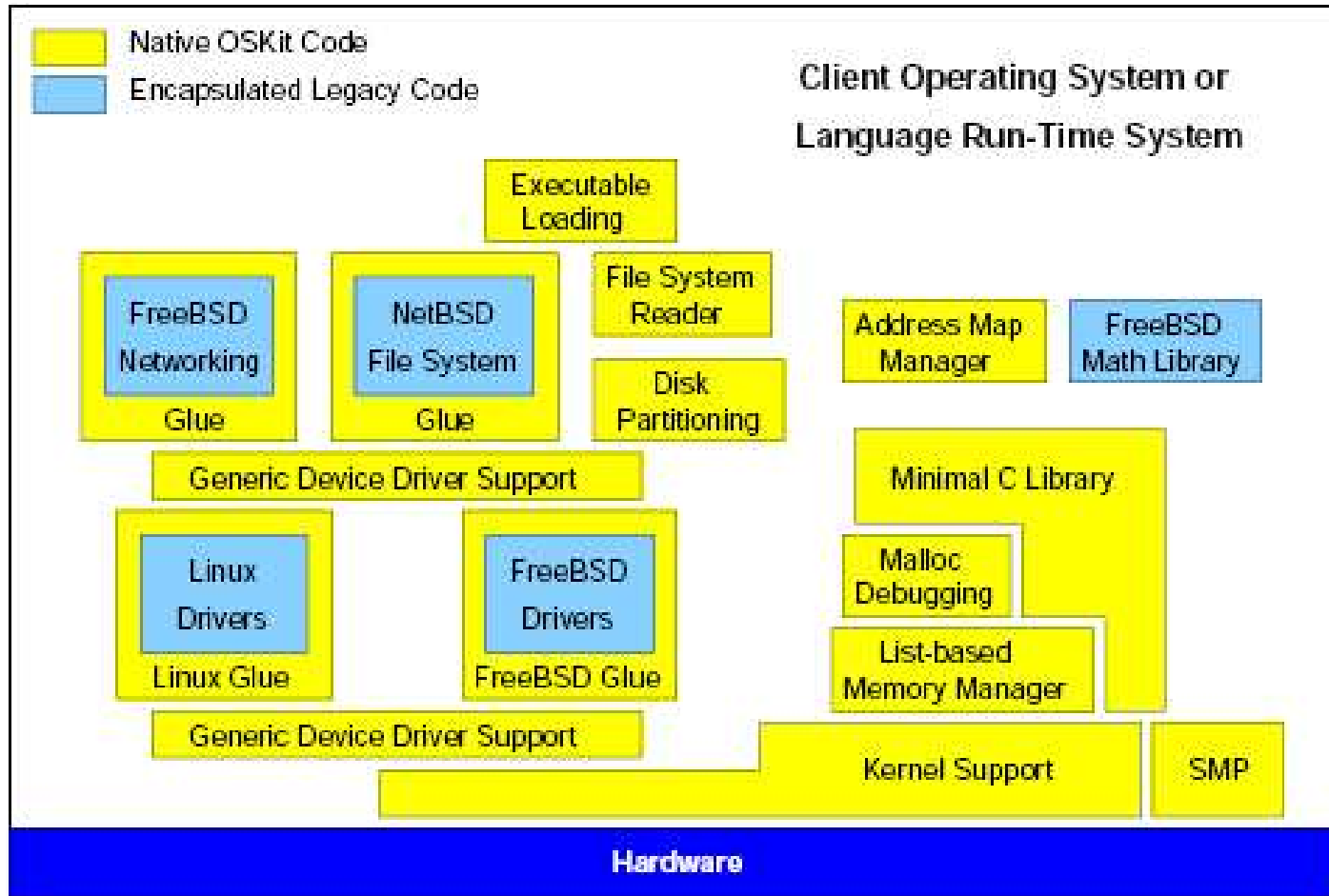
## Flux OSKit

- ▶ Tons of functionality  
(Linux, BSD, Java, SML, ...)
- ▶ Not tuned for high performance
- ▶ Implementation details exposed

# Flux OSKit

- ▶ Framework and reusable OS components
- ▶ Focus on component of research-interest
- ▶ Reuse other existing components for functionality

# Flux OSKit



## Flux OSKit Components

# Flux OSKit

- ▶ Bootloader
  - ▶ Multiboot compliant
- ▶ Kernel Support Library
  - ▶ Architecture specific
- ▶ Memory Management Library
  - ▶ `kmalloc()`, alignment etc.
- ▶ Minimal libc
  - ▶ non-buffered `read()`, `write()` etc.
  - ▶ minimizes dependencies

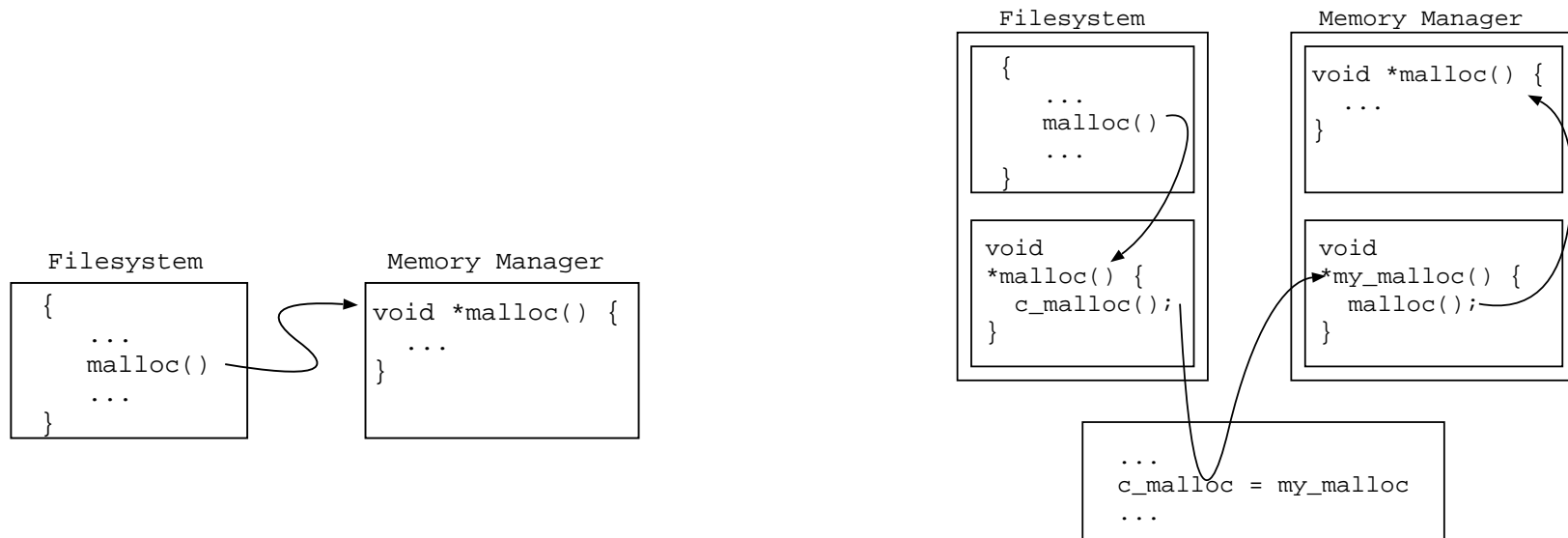


# Flux OSKit

- ▶ Debugging Support
  - ▶ GDB over serial line
- ▶ Device Driver Support
  - ▶ Drivers from Linux, FreeBSD inside wrappers
- ▶ Protocol Stacks
  - ▶ “Wrapped” FreeBSD network stack
- ▶ File System
  - ▶ “Wrapped” NetBSD code

# Flux OSKit

- ▶ OSKit components are separable, no dependence
  - ▶ Other OS: Modularity does not imply independence



Very little overhead

- ▶ Provides abstractions
- ▶ Doesn't hide implementation

# Case Studies

- ▶ ML/OS
  - ▶ SML: Static Typing, Concurrency through continuations, No stack, Aggressive heap usage, Interpreted.
  - ▶ ML/OS: 2 people, one semester using OSKit
- ▶ Java
  - ▶ Existing JVM
  - ▶ Java/OS: 3 weeks using OSKit
- ▶ SR
  - ▶ Concurrent programming language

# Summary

- ▶ L<sup>4</sup>-Linux:  $\mu$ -Kernels can be fast
  - ▶ Full system binary-compatible with Linux runs 5%–10% slower.
- ▶ FluxOSKit: Kernels from reusable components
  - ▶ Write fully-functional research OS in weeks